

---

# Supplementary

---

1	<b>Contents</b>	
2	<b>1 PrinCut</b>	<b>2</b>
3	1.1 How to use PrinCut . . . . .	2
4	1.2 Methodology . . . . .	3
5	1.2.1 Boundary-refine algorithm and adding a cell . . . . .	3
6	1.2.2 Merging two cells and splitting a cell . . . . .	5
7	1.2.3 Smart click correction . . . . .	5
8	<b>2 Label fusion detail</b>	<b>6</b>
9	<b>3 Experiment detail</b>	<b>7</b>
10	3.1 Experiment setting . . . . .	7
11	3.1.1 Cellpose . . . . .	7
12	3.1.2 QCAnet . . . . .	7
13	3.1.3 StarDist . . . . .	7
14	3.1.4 3Dsuite . . . . .	7
15	3.1.5 Vaa3D . . . . .	7
16	3.2 Other methods we tried . . . . .	8
17	<b>4 Other questions</b>	<b>9</b>
18	4.1 Compared with STABLE . . . . .	9
19	4.2 Why three annotators . . . . .	9
20	4.3 Supervised learning . . . . .	9
21	4.4 Evaluation metrics design . . . . .	10

## 22 1 PrinCut

### 23 1.1 How to use PrinCut

24 The PrinCut GUI is shown in Figure 1. PrinCut is a MATLAB app, and its package is also provided  
25 in the supplementary. The app is tested and used on MATLAB 2022b.

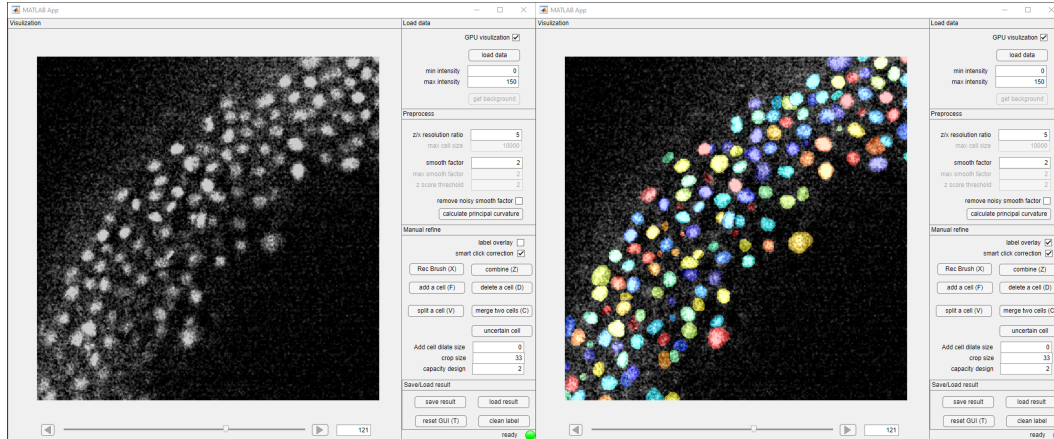


Figure 1: The GUI of PrinCut. The left shows raw data without annotation. The right shows both raw data and annotation overlay.

26 After loading the raw data (the data should be in tif format), users can adjust the contrast, zoom  
27 in/out, and change different z slices to visualize the data. Users then need to give the z/x resolution  
28 ratio and the smooth factor to calculate the principal curvature. For example, if the voxel size of  
29 data is  $1\mu m \times 1\mu m \times 5\mu m$ , then the z/x resolution ratio should be 5. The smooth factor is the  
30 standard deviation (by pixel) of the Gaussian filter used to smooth data before calculating the principal  
31 curvature. After adjusting the z/x resolution ratio and the smooth factor, click the "calculate principal  
32 curvature" button.

33 When the principal curvature is calculated, users can add a cell label by simply clicking the cell,  
34 and a 3D suggestive boundary will be automatically generated. We request the annotators check the  
35 suggestive boundary on every z slice to make sure it's correct. If the suggestive boundary contains  
36 more than one cell, users can split the cell by clicking the center of each cell and PrinCut can give a  
37 suggestive boundary to each cell. If the suggestive boundary only contains part of a cell, users can  
38 merge two labels and PrinCut will give a suggestive boundary including the original labels. Users  
39 can also use the "combine" button to combine two existing boundaries without generating suggestive  
40 boundaries. Examples are shown in Figure 2. Four additional parameters are used to control the  
41 suggestive boundary, which are "smart click correction", "Add cell dilate size", "crop size", and  
42 "capacity design". In the methodology section, we will discuss how those parameters influence the  
43 suggestive boundary.

44 However, there is a chance that the gap between two cells is too weak that the principal curvature is  
45 still negative. In this case, users need to use the brush to draw the expected boundary. The size of the  
46 brush can be adjusted by users, as is shown in Figure 5.

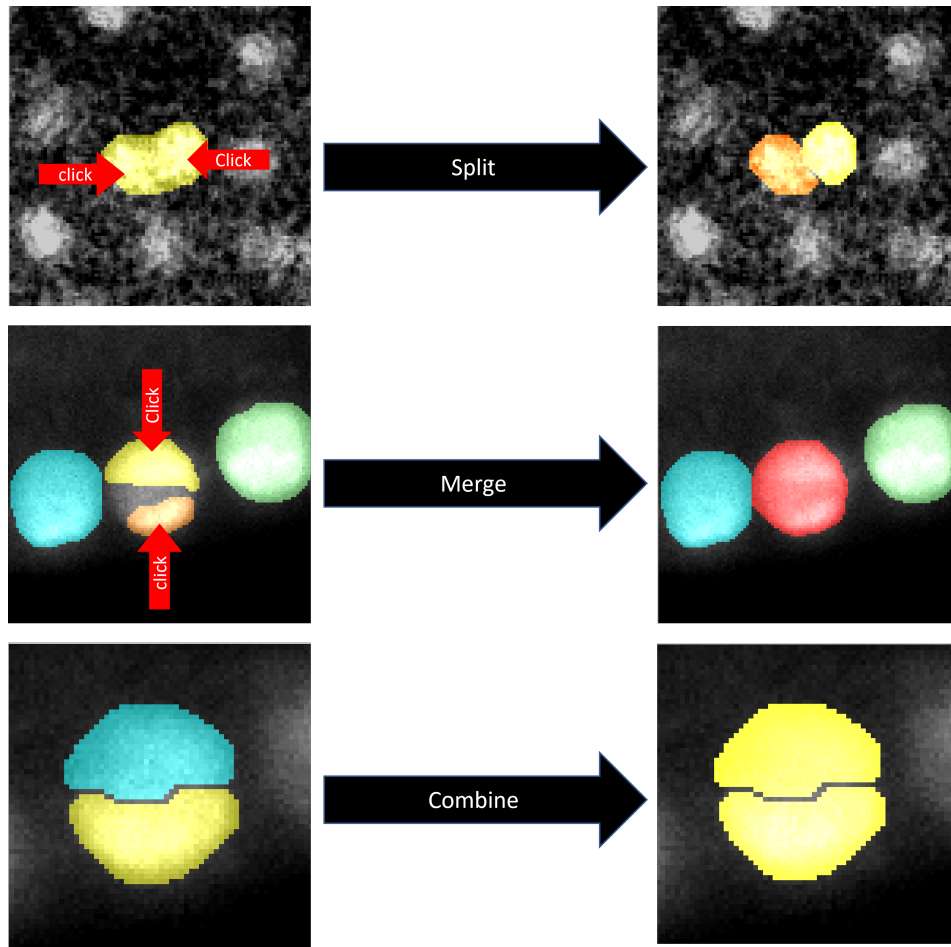


Figure 2: Examples of splitting, merging, and combining existing cells by clicking. PrinCut can automatically generate a suggestive boundary after split or merge process.

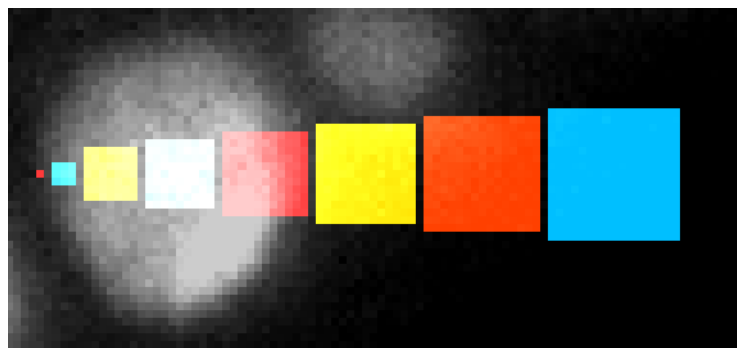


Figure 3: The brushes with different sizes, from a single pixel to  $17 \times 17$  pixels.

## 47 1.2 Methodology

### 48 1.2.1 Boundary-refine algorithm and adding a cell

49 The boundary-refine algorithm can grow a given seed to the suggestive boundary in a given foreground.  
 50 For adding a cell process, the seed is a spherical region using the pixel user clicked as the center and  
 51 the "add cell dilate size" parameter in GUI as the radius, and the foreground is a spherical region with  
 52 the same center and using the "crop size" parameter in GUI as the radius. The seed should be within  
 53 the cell while the foreground should be larger than the cell.

54 The optimization criteria of the boundary-refine algorithm are to grow the seeds such that the principal  
55 curvature along the boundary pixels is maximized. However, this constraint alone is not enough,  
56 since the seed might grow to an excessively large region to meet the constraint, so we also need to  
57 add the constraint that the boundary should be as short as possible. Furthermore, since there can be  
58 multiple seeds inside one region, we need to make sure the regions grown from each seed do not  
59 interfere with each other, or the regions are not intersected with each other.

60 Marker-based watershed inside the foreground is one common solution to solve the problem. But it  
61 is solely based on the score map and will always grow the marker so that all the pixels are used for  
62 segmentation. But in our case, we don't want the seed to grow and fill the whole foreground, and the  
63 grown regions of all the seeds do not need to be in contact with each other. To align the boundary to  
64 the optimal position, we formulate the problem into an optimization problem.

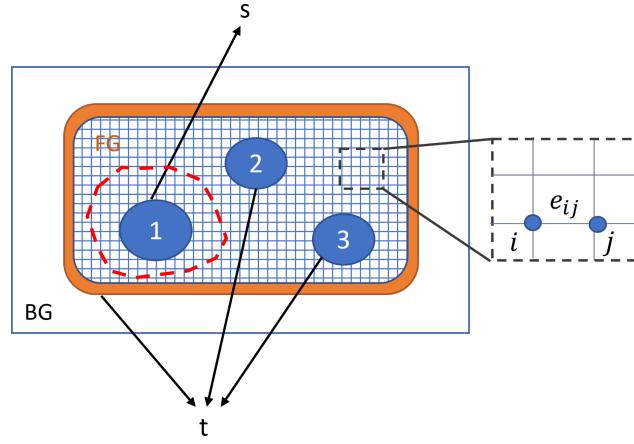


Figure 4: The graph for solving the refined boundary based on the seed

65 Consider the labeling of each pixel  $x_i^n$  in the foreground as  $u_i^n$ , if  $u_i^n = 1$ , then  $x_i$  will be assigned to  
66 the seed  $n$ , otherwise if  $u_i^n = 0$ , the corresponding pixel will be assign as outer region. And we can  
67 define the group of pixels belonging to the seed  $n$  as  $S_n$  and the group for the rest of pixels as  $\bar{S}_n$ .  
68 And the boundary between the two groups can be represented using the pairs of pixels  $C_n$ . In each  
69 pair, one belongs to the grown seed region, the other belongs to the background region.

$$C_n = \{(x_i, x_j) | x_i \in S_n, x_j \in \bar{S}_n\} \quad (1)$$

70 Then the objective function for seed  $n$  can be expressed as,

$$\arg \min_{C_n} \sum_{(i,j) \in C_n} (\hat{G}(i) + \hat{G}(j)) \quad (2)$$

71 where  $\hat{G} = \frac{1}{\max(G, T)^p}$  and  $G$  is the map of the principal curvature for each pixel,  $p$  is the "capacity  
72 design" parameter in GUI,  $T$  is a constant as 0.001. In this way, we transform the problem of  
73 finding the boundary pixels maximizing in the map  $G$  into minimizing in the map  $\hat{G}$ . In addition,  
74 this objective function also implicitly minimizes the boundary length. In order to solve 2, we can  
75 reformulate the problem into a min-cut problem that can be efficiently solved. First, we will introduce  
76 graph construction. As shown in figure 4, the graph for optimizing based on seed 1 is built based on all  
77 the pixels in the 3D foreground, and each node  $i$  represents a pixel  $i$ . Between each pair of neighbor  
78 nodes, an edge  $e_{ij}$  is linked between them. The edge weight is defined as  $weight(e_{ij}) = \hat{G}(i) + \hat{G}(j)$ .  
79 In addition, one pair of pseudo nodes are added, one is the pseudo source node  $s$  and the other the  
80 pseudo sink node  $t$ . All the pixels in the corresponding seed will be connected to the source node,  
81 and the boundary of the foreground will be connected to the sink (labeled in orange in figure 4). The  
82 weights of these two types of edges are set as infinite or very large.



83 With this graph design, the labeling of the nodes can induce a cut set  $C = (i, j) | u(i) \neq u(j)$  and the  
84 following optimization of the labeling  $U$  for seed  $n$  is the same as 2

$$u^n = \arg \min_{u_i, i \in [1, N]} \sum cut(u^n) \quad (3)$$

85 where  $cut(u) = cut(C) = \sum_{i, j \in C} weight(i, j)$ , which corresponds to the sum of the weights for  
86 the edges that are cut (shown in red dashed line in figure 4). By solving the min-cut of the graph, we  
87 can obtain the region grown from seed 1, which are the pixels connected to the source after cutting.

### 88 1.2.2 Merging two cells and splitting a cell

89 Once we know how to add a cell, merging and splitting are straightforward. When merging two cells,  
90 we consider the regions with two old labels as the seed and do morphological dilation for the seed to  
91 get the foreground. The filter of morphological dilation is a sphere using "crop size" as the radius.  
92 Then we do the boundary refinement for the given seed and foreground. When splitting a cell, we  
93 delete the old label first and add two seeds simultaneously. Then we calculate the foreground and  
94 refine the boundary for each seed. As the other seed will be considered as the sink when we refine  
95 the current seed, the grown regions of the two seeds will not overlap with each other.

### 96 1.2.3 Smart click correction

97 The boundary-refine algorithm requires the user to consistently click at the center of the cell. Failure  
98 to do so may result in a suggestive boundary that encompasses only the seed region. This limitation  
99 arises from situations where a user clicks at the edge of a cell and all pixels surrounding the seed  
100 exhibit relatively high positive principal curvature since using the edges around the seed as the  
101 boundary in such cases leads to an even smaller  $cut(u)$  compared to the correct boundary.

102 To address this issue, the "smart click correction" algorithm leverages the pixel within the seed as  
103 the source and identifies pixels with negative principal curvature as the sink. Using the same graph  
104 discussed earlier, the algorithm calculates the shortest path between the source and sink. When the  
105 "smart click correction" checkbox is activated, all pixels along the shortest path are used to generate  
106 a new seed, which will replace the original seed in the boundary refinement algorithm.

107 **2 Label fusion detail**

108 After aligning the labels assigned by various annotators with a designated ground truth, the subsequent  
109 step entails amalgamating these annotations to delineate the boundary of the ground truth. This  
110 process is underpinned by the principle of weighted voting, wherein the contribution of each annotator  
111 is influenced by factors such as the reviewer’s expertise and the annotator’s own experience. But there  
112 are some additional rules we applied to solve some overlapping issues. For two ground truth labels  
113 partially overlapping and two ground truth labels having different confidence scores, the ground truth  
114 label with the higher confidence score will keep the same. If more than 66% pixels of the lower  
115 confidence score ground truth are overlapping with the other label, it will be removed. If less than  
116 33% of its pixels are overlapping with the other label, it will be kept but only use the pixels not  
117 overlapping with the other label. Otherwise, the non-overlapping pixels of lower confidence score  
118 ground truth will be considered undefined masks. For two ground truth labels partially overlapping  
119 and two ground truth labels having the same confidence score, we will sort annotators by their  
120 annotation experience, and treat the label created by a more experienced annotator as the one with a  
121 higher confidence score, then do the same thing we did to two ground truths with different confidence  
122 scores.

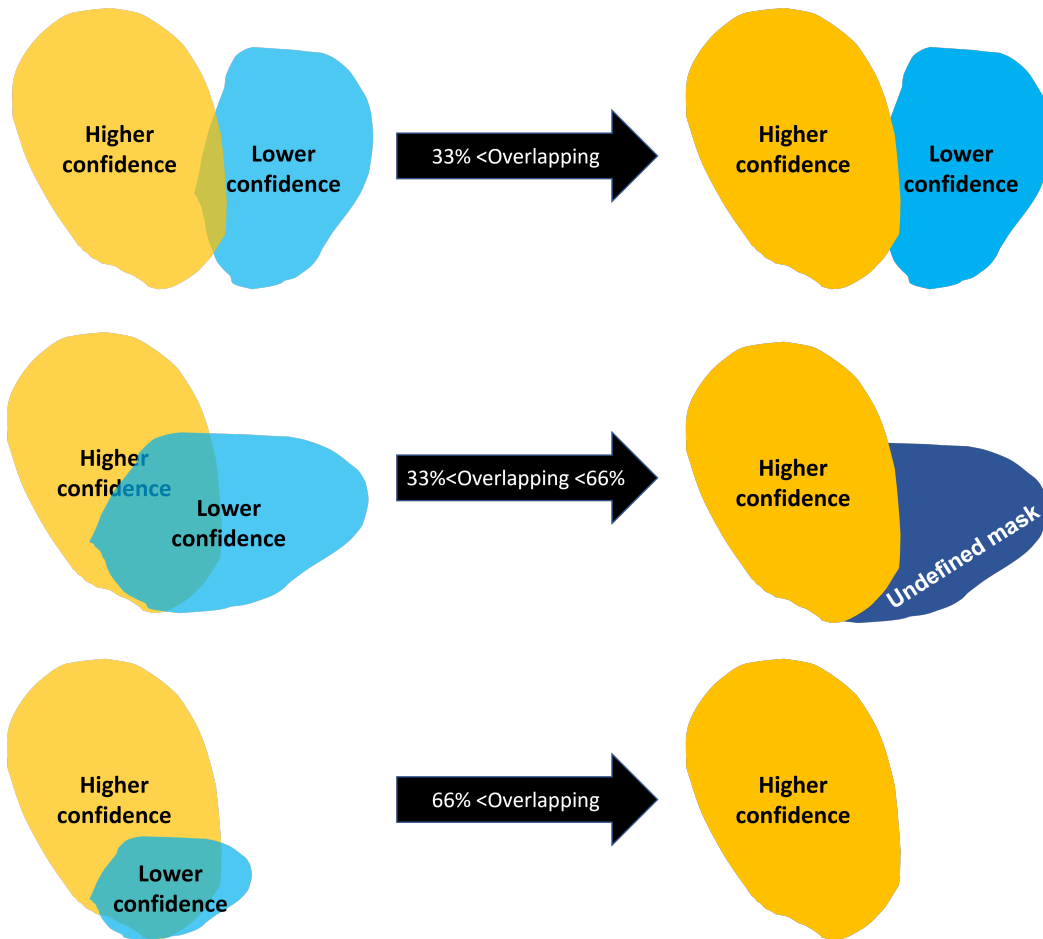


Figure 5: Example of the solution for overlapping problems

## 123 **3 Experiment detail**

### 124 **3.1 Experiment setting**

125 Quite a few methods set the minimum and maximum cell size/volume as hyperparameters. We report  
126 the average cell sizes (each from 10% sample cells in the corresponding image) here so that readers  
127 can copy them directly when testing the existing methods. In Zebrafish 1, the cell size is around  
128 2500. In Zebrafish 2, the cell size is around 400. In Drosophila 1, the cell size is around 4000. In  
129 Drosophila 2, the cell size is around 10000. In Mus Musculus 1, the cell size is around 2500. Other  
130 settings for each method are summarized below.

#### 131 **3.1.1 Cellpose**

132 The developers of Cellpose have provided the pre-trained models to do instance segmentation for  
133 3D images on their GitHub. We adjust the cell diameter and z-aspect for both "cyto" and "nuclei"  
134 models as shown in table 1. Everything else is the same as the default.

Table 1: Cellpose parameters

Name	diameter	z-aspect
zebrafish 1	20	5.4
zebrafish 2	5	1.0
Drosophila 1	30	1.0
Drosophila 2	35	1.0
Mus Musculus	30	1.0

#### 135 **3.1.2 QCAnet**

136 The developers of QCAnet have provided the pre-trained models to do instance segmentation for 3D  
137 images on their GitHub. All settings remained as default.

#### 138 **3.1.3 StarDist**

139 The developers of StarDist have provided the pre-trained models to do instance segmentation for 3D  
140 images on their GitHub. Moreover, they have provided a detailed tutorial on how to use their 2D  
141 segmentation model on GitHub. We referred to their tutorial and only changed the path of the model  
142 from 2D to 3D. Other settings remained as default.

#### 143 **3.1.4 3Dsuite**

144 The developers have provided a Fiji plug-in for the 3d Suite method, which has included several  
145 different segmentation algorithms such as 3D watershed, 3D spot segmentation, 3D iterative thresh-  
146 olding, etc. We select the 3D iterative thresholding algorithm since it yields the best segmentation  
147 performance. The main hyperparameters that need to be tuned are the minimum cell volume, maxi-  
148 mum cell volume (both in terms of the number of pixels). We adjust the minimum and maximum  
149 volumes according to the cell sizes in each image.

#### 150 **3.1.5 Vaa3D**

151 The developers of Vaa3D have provided a software for 3D bioimage processing. In the Vaa3D  
152 software, the only segmentation algorithm is based on gradient vector flow. The hyperparameters that  
153 need to be tuned include: the iterations of diffusion, fusion threshold, and the minimum cell size. We  
154 employ the default setting for the first two hyperparameters while the third one is based on the cell  
155 sizes in each image.

156 **3.2 Other methods we tried**

157 There are many famous semantic segmentation methods, but we didn't compare them with instance  
158 segmentation methods. If the cells are far away from each other, we can use some post-processing to  
split them easily, but for data like NIS3D, they are not very useful. An example is shown in Figure 6.

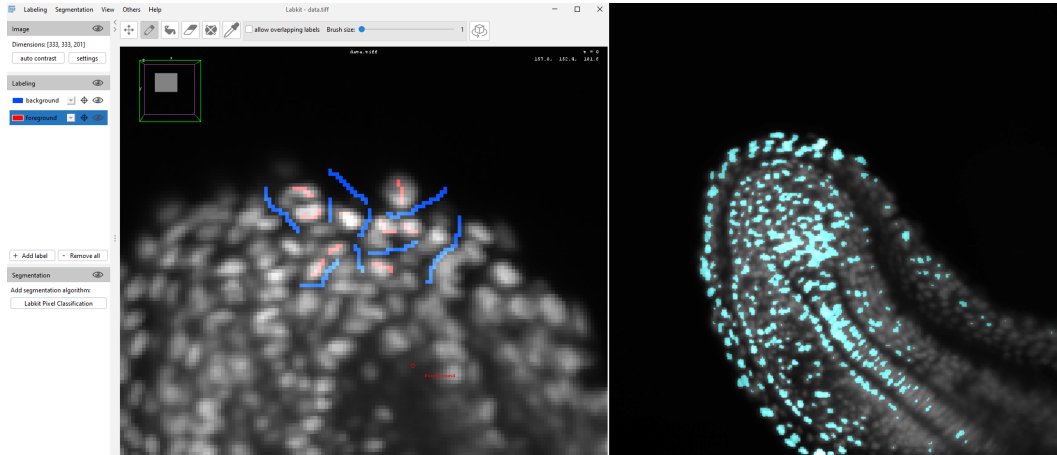


Figure 6: An example of semantic segmentation result (Labkits). The left is the human input, and the right is the segmentation result. The connected component can only give several instances as they are highly connected in 3D space.

159

160 We also tested the methods from MorphoLibJ, which provides a group of unsupervised 3D instance  
161 segmentation methods. We tried a wide range of parameters, but the result is much worse than the  
162 other method we tested and we did not put it into baseline methods. An example is shown in Figure 7.



Figure 7: Segmentation result by morphological segmentation in MorphoLibJ. This result is based on parameters: "object image" and "gradient radius = 3"

163 **4 Other questions**

164 **4.1 Compared with STABLE**

165 An alternative approach involves combining the ground truth created solely by human annotators  
166 using the STAPLE algorithm with conventional unweighted F1 scores or Intersection over Union  
167 (IoU) scores. This particular method attributes uniform weight to all cells identified through the  
168 STAPLE approach, whereas cells exclusively identified by a single annotator are typically classified  
169 to a background. Besides this, STAPLE algorithm will still generate misleading ground truth when  
170 annotators give conflicted labels, rather than giving an undefined mask like our method. Drawing  
171 from our experience, we have observed a preference among users for the identification of strong cells,  
172 and a willingness to accept the omission of inconspicuous cells that might elude human perception.  
173 This differential weighting within our proposed method aligns with this user preference.

174 **4.2 Why three annotators**

175 The intricacies of annotating 3D nuclei are profound, and it is almost inevitable that within this  
176 complex process, discrepancies such as incomplete or erroneous labels may arise (check Figure  
177 Example of incomplete labels and noise label of C.elegans dataset). Without proofreading from  
178 different annotators, the quality of ground truth cannot be guaranteed. In contrast, NIS3D not  
179 only boasts a larger number of annotated instances, but also encompasses a wider array of species,  
180 experimental situations, and developmental stages. Besides, NIS3D carefully designed a strategy to  
181 fuse the ground truth from different annotators, which increase the boundary accuracy and significantly  
reduces the chance of incomplete labels.

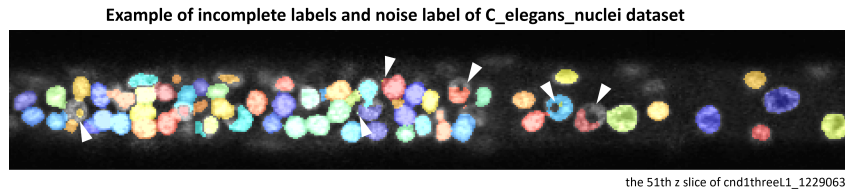


Figure 8: Example of incomplete labels and noise labels

182

183 **4.3 Supervised learning**

184 We fine-tuned Stardist model to show the ability of the dataset to improve the existing model. Here  
185 we use in-image split (50% of the image as the training set and the other 50% as the test set for all  
186 images) and default train parameters provided by Stardist. We didn't use the confidence score and  
187 treated the undefined mask as pixels without cells. With the default setting, Stardist model shows  
significant improvement under the measure of both W-F1 score and W-SEG score.

Stardist pre-trained model result					
data	W-F1	W-Precision	W-Recall	W-IoU	W-SEG
Zebrafish 1	0.586	0.770	0.473	0.263	0.192
Zebrafish 2	0.529	0.977	0.363	0.380	0.183
Drosophila 1	0.924	0.881	0.972	0.691	0.586
Drosophila 2	0.684	0.526	0.979	0.563	0.325
Mus Musculus	0.594	0.789	0.476	0.256	0.191

Stardist fine-tuned model results					
data	W-F1	W-Precision	W-Recall	W-IoU	W-SEG
Zebrafish 1	0.942	0.969	0.915	0.800	0.687
Zebrafish 2	0.804	0.750	0.867	0.707	0.378
Drosophila 1	0.971	0.995	0.947	0.860	0.772
Drosophila 2	0.985	0.993	0.978	0.873	0.822
Mus Musculus	0.793	0.982	0.664	0.597	0.475

Figure 9: Stardist before and after a simple fine-tuning for in-image split

188

#### 189 **4.4 Evaluation metrics design**

190 Choosing the correct metric that adequately reflects the biological nature is important but usually  
191 neglected. When employing the existing metric of the cell tracking challenge and the 2018 Data  
192 Science Bowl, three distinct issues encountered deviate result from human intuition, thereby affecting  
193 the evaluation process. Firstly, some methods can only detect the central portions of cells. Considering  
194 the true positive criteria in those existing benchmarks requiring 50% or more IoU score, even  
195 when a cell is successfully identified by 3D suite, its precision and recall may fall short of human  
196 expectations. This disparity becomes particularly relevant for biological studies that prioritize cell  
197 location, rendering the IoU score less informative. Secondly, approaches like Vaa3D yield favorable  
198 foreground detection outcomes but are plagued by a pronounced under-segmentation dilemma,  
199 existing metrics can't show their advantage for foreground detection. Thirdly, users express a  
200 preference for the detection of robust cells and can tolerate the omission of inconspicuous ones that  
201 may be overlooked by a human. To address these issues, we have reformulated the evaluation metric  
202 to align more closely with our specific objectives. For instance, a high W-F1 score coupled with  
203 a low W-SEG score now indicates successful cell detection while indicating room for boundary  
204 enhancement. Similarly, a high W-IoU score combined with a low W-SEG score signifies accurate  
205 foreground detection, while highlighting potential over-segmentation or under-segmentation concerns.