

Scellseg: a style-aware cell instance segmentation tool with pre-training and contrastive fine-tuning

Authors

Dejin Xun¹, Deheng Chen², Yitian Zhou³, Volker M. Lauschke^{3,4,5}, Rui Wang^{2*}, Yi Wang^{1,6,7*}

Affiliations

1. Pharmaceutical Informatics Institute, College of Pharmaceutical Sciences, Zhejiang University, Hangzhou, Zhejiang 310058, China

2. State Key Lab of CAD&CG, Zhejiang University, Hangzhou, Zhejiang 310058, China

3. Department of Physiology and Pharmacology, Karolinska Institutet, Stockholm, Sweden

4. Dr Margarete Fischer-Bosch Institute of Clinical Pharmacology, Stuttgart, Germany

5. University of Tuebingen, Tuebingen, Germany

6. Innovation Institute for Artificial Intelligence in Medicine of Zhejiang University, Hangzhou, Zhejiang 310018, China

7. State Key Laboratory of Component-based Chinese Medicine, Tianjin University of Traditional Chinese Medicine, Tianjin 300193, China

*Corresponding author. Email: zjuwangyi@zju.edu.cn (Y.W.); ruiwang@zju.edu.cn (R.W.)

Abstract

Deep learning-based cell segmentation is increasingly utilized in cell biology and molecular pathology, due to massive accumulation of diverse large-scale datasets and excellent progress in cell representation. However, the development of specialized algorithms has long been hampered by a paucity of annotated training data, whereas the performance of generalist algorithm was limited without experiment-specific calibration. Here, we present a deep learning-based tool called Scellseg consisted of novel pre-trained network architecture and contrastive fine-tuning strategy. In comparison to four commonly used algorithms, Scellseg outperformed others in average precision and Aggregated Jaccard Index on three disparate datasets. Interestingly, we found that eight images are sufficient for model tuning to achieve satisfied performance based on a shot data scale experiment. We also developed a graphical user interface integrated with functions of annotation, fine-tuning and inference, that allows biologists to easily specialize their self-adaptive segmentation model for analyzing images at the single-cell level.

MAIN TEXT

Introduction

Image-based single cell profiling is widely used in biological, pharmaceutical and medical applications, including in quantitative cytometry¹, spatial transcriptomics², high-content drug screening³ and cancer metastasis analysis⁴. However, due to a lack of robust and facile algorithms for single-cell analysis, average profiling remains as the most commonly used method which may cause loss of information and mislead interpretation of feature associations⁵. In recent years, deep learning has revolutionized the field of computer vision⁶ and catalyzed the advancement of single cell segmentation methods.

44 Differences across cell types, microscopy instruments, treatment methods, imaging modalities,
45 and staining protocols can generate cell images with considerable diversity. As a consequence, cell
46 segmentation algorithms were mostly developed for specific datasets⁷⁻⁹ and these methods
47 performed poorly when applied to other styles of cell images. To overcome this limitation,
48 generalist algorithms have been developed. In 2018, a data science bowl challenge tried to segment
49 nuclei from a large number of different styles of microscopy images using 841 diverse images
50 containing 37,333 nuclei¹⁰. Inspired by this competition, Stringer et al. annotated 608 images
51 containing more than 70,000 segmented objects and developed a generalist algorithm named
52 Cellpose, which exhibited excellent performance in segmenting cell bodies from many image
53 styles¹¹. Although deep-learning based generalist algorithms outperformed compared to traditional
54 machine learning approaches like logistic regression and Random Forest (RF), the state-of-art
55 segmentation tools still lack of capability to be self-adaptive for all kinds of cellular images.
56 Therefore, transfer learning of segmentation models from certain source domain to in-house
57 datasets remains an important challenge for biologists with few computational knowledges.

58 Fine-tuning of pre-trained models has been successfully used in computer vision¹²⁻¹⁵ and
59 natural language processing¹⁶⁻¹⁷ due to its lower input requirements and more rapid convergence to
60 a better performance. For cell instance segmentation, only few preliminary attempts were reported,
61 such as fine-tuning of a nuclear segmentation model to satisfy different needs from distinct
62 laboratories¹⁸, or transferring a pre-trained model of in vitro images to in situ tissue images¹⁹.
63 However, these studies used only nuclei images for pre-training and performed binary or multiclass
64 classification instead of instance topological maps, hence are hard to capture enough prior
65 knowledge for fine-tuning on different kinds of cell images. Besides, they primarily tested model
66 transferability on different nuclei images, specialized evaluation datasets for various cell-like
67 instances such as *C. elegans*²⁰, are far from well developed and studied. Hence, the development of
68 a high-performance universal computational pipeline based on the fine-tuning of pre-trained models
69 remains a challenging but important objective in automated image analysis.

70 In this work, we established a fine-tuning pipeline for cell segmentation algorithms and present
71 a style-aware cell segmentation architecture named Scellseg based on attention mechanisms and
72 hierarchical information to improve the extraction and utilization of style features. We furthermore
73 incorporate a contrastive learning strategy to leverage information from unlabelled and pre-trained
74 data. To evaluate the generalizability of the pipeline, we benchmarked our model on three
75 fundamentally different styles of data from *C. elegans*, label-free phase-contrast cell images, and
76 sub-cellular organelles. Furthermore, it is our first effort to estimate the minimal extent of data
77 required for a satisfying fine-tune model and to demonstrate how instance representation and pre-
78 trained datasets can influence model transferability. To facilitate uptake of this pipeline, we
79 developed a graphical user interface (GUI) which can conduct annotation, fine-tuning and inference,
80 thus making the model accessible for a wide range of users without coding experience. The model
81 can be found at <https://github.com/cellimnet/scellseg-publish>.

84 Results

85 Design of Scellseg with pre-trained architecture and contrastive fine-tuning

86 Firstly, we established a pre-trained and fine-tuning pipeline for the cell segmentation model.
87 For initial training, we utilized a dataset containing various cell types to build a generalist model.
88 Generally, segmentation of untrained images by this model will exhibit limited performance, such
89 as fail of detecting instances or boundary of segmentation instance. To improve the specificity of
90 model and avoid time-consuming re-training, several images from novel data pool can be annotated
91 for fine-tuning established model using few new labelled data (shot data) (Fig. 1a). This workflow
92 generated a style-aware structure to better extract and comprehend style-related information and

93 developed a new fine-tuning strategy based on contrastive learning to better make use of diverse
94 data features, including the unlabelled data (query data) and pre-trained data. The resulting model,
95 which we named Scellseg, contains two branches, a mask branch to compute the segmentation map
96 of input and a contrast branch to explore the information between three types of data (Fig. 1b). The
97 mask branch is utilized during pre-training, fine-tuning and inference, whereas the contrast branch
98 worked only during fine-tuning.

99 The mask branch was based on Cellpose model which is a member of the U-Net²¹ family of
100 algorithms that consist of a downsampling pass that extracts features from input data, an
101 upsampling pass that organizes different features to fit for the final task, and a concatenation
102 operation that relays the information extracted from downsampling process to the upsampling pass.
103 For convenient adjustment, we re-divided this structure, the last Conv Unit was split from
104 upsampling pass and named as Tasker, and the left was named as Extractor, containing the
105 downsampling, upsampling and concatenation parts. To improve the sensitivity of model for
106 different styles, we added attention gates (AGs) when passing the features extracted from
107 downsampling to the upsampling pass. These AGs give the feature map weights to highlight salient
108 features useful for a specific task and suppress feature activation in irrelevant regions. We used
109 dense units to consider the information from early upsampling layers, aiming to delineate accurate
110 object boundaries. To consider different-level style information, we also fed corresponding
111 hierarchical style embeddings into different-level dense units (Extended Data Fig. 1).

112 Unlike conventional fine-tuning strategies only use labelled data, to augment data utilization,
113 we developed a contrastive fine-tuning (CFT) method to employ information from either labelled
114 data or unlabelled and pre-trained data based on contrastive learning. Seven common cellular styles
115 of images were chosen from pre-trained data to form contrast data (Extended Data Fig. 2). In the
116 contrastive fine-tuning process, the contrast branch is used to compute the respective style
117 embeddings of these three data and then a contrast loss function was designed to minimize the
118 difference between embeddings of shot and query data from the same experiment while maximizing
119 the difference between embeddings of shot and contrast data (Fig. 1b). This contrast loss was added
120 into the segmentation loss function, then the total loss optimizes the model via backpropagation.

121 **Transferability of Scellseg with contrastive fine-tuning strategy on three evaluation datasets**

122 To compare the performance of Scellseg in the transferability of models with other algorithms,
123 we adopted three difference datasets named BBBC010_elegans²⁰, LIVECell_bv2²², and mito (in-
124 house prepared dataset containing mitochondrial images), representing three levels of cell-like
125 images from model organism, cells to subcellular structures (Fig. 2a). In total, the datasets contained
126 230 images and 91,024 segmentation objects. We visualized the distribution of areas and numbers
127 for cells per image (Extended Data Fig. 3). The average areas for three datasets are about 1000, 150
128 and 100,000, and numbers of instance in each image ranged from 2 to 2,815. We used t-distributed
129 stochastic neighbor embedding²³ (t-SNE) to visualize the style embeddings (see definition in ref.
130 ¹¹) of these evaluation datasets together with pre-trained datasets and noted that the style of data in
131 each dataset was determinant in major cluster (Fig. 2b).

132 To compare the influence of different instance representation, we benchmarked Scellseg
133 against four other models, U-Net²¹, U-Net³²⁴, HoVer²⁵ and Cellpose¹¹. These four models were
134 set with identical network structure and pre-trained with the same datasets and training strategies.
135 We used the training data of each dataset to fine-tune the model at ten different random states, most
136 models achieved great improvements after fine-tuning, for the BBBC010_elegans dataset, all
137 models yielded at least 35% higher average precision. For U-Net3 model, fine-tuning strategy even
138 yielded a dramatic increase of 62.1% in average precision. The different models differed drastically
139 in performance, and representation of Cellpose (used in Cellpose and Scellseg model) outperformed
140 other methods. Scellseg with contrastive fine-tuning achieved the best performance on all three
141

142 datasets, especially on the BBBC010_elegans dataset. At the universally used intersection over
143 union (IoU) threshold of 0.5, our Scellseg and Cellpose both achieved high average precision when
144 segmenting *C. elegans* (0.882 for Scellseg; 0.868 for Cellpose), microglial cell BV-2 (0.783 for
145 Scellseg; 0.784 for Cellpose), and mitochondria in cardiomyocytes (0.927 for Scellseg; 0.922 for
146 Cellpose). In contrast, Scellseg performed considerably better at higher thresholds, such as 0.75 on
147 all three data sets ([0.670, 0.493, 0.634] for Scellseg compared to [0.587, 0.475, 0.571] for Cellpose,
148 respectively; Fig. 2c-e).

149 We also compared the performance of Scellseg with or without contrastive fine-tuning strategy.
150 As shown in Fig. 3, it is clear that the fine-tuned model exhibited considerably better capability of
151 instance detection. Importantly, fine-tuning strategy improved the ability of distinguishing adjacent
152 cells, which allowed the segmentation of scattered mitochondria around the nuclei in mito dataset.
153 Furthermore, our contrastive fine-tuning strategy outperformed the classic method on the Cellpose
154 test set after fine-tuning on the three evaluation datasets. However, as expected, all re-trained
155 models suffered a sharp decline compared with the initial generalization ability (Extended Data Fig.
156 4).

158 **Pre-trained dataset scale experiments**

159 To explore how the pre-trained dataset can influence model transferability, we used different
160 subsets of the Cellpose training set. The initial subset (Sneuro) only contains one style of images
161 from the Cell Image Library²⁶, then additional styles of images were sequentially added, such as
162 fluorescent cells (Sfluor), non-fluorescent and membrane-labelled cells (Scell), other microscopy
163 data (Smicro), as well as non-microscopy images (“Sgeneral”, corresponding to the full Cellpose
164 train set; Fig. 4a).

165 We pre-trained Scellseg with Sneuro, Sfluor, Scell, Smicro and firstly tested the generalization
166 ability of each model by applying it directly without any adaptation on three evaluation datasets.
167 For *C. elegans*, the model trained with Sneuro, Sfluor, and Scell does not result in successful
168 recognition until the pre-trained dataset contains microscopy instances with structures beyond cells.
169 For small and round BV-2 cells, the generalization ability also increased with the richness of the
170 dataset from Sneuro to Smicro, and model trained on Smicro even outperformed the model trained
171 on Sgeneral. For the segmentation of mitochondria, surprisingly the model trained with Sfluor and
172 Smicro outperformed all others (Fig. 4b).

173 Next, we tested the transferability of each model (Fig. 4c). As expected, the transferability of
174 Scellseg increased with the richness of the dataset from Sneuro to Smicro, and Scellseg pre-trained
175 with Smicro achieved similar transfer performance on three evaluation datasets compared to
176 Scellseg pre-trained with Sgeneral ([0.555, 0.451, 0.564] and [0.554, 0.454, 0.557], mean average
177 precision [mAP] of Smicro and Sgeneral respectively). On the BBBC010_elegans dataset, a model
178 pre-trained on Sneuro achieved only very poor transfer performance (0.013 mAP) and performance
179 increased substantially only after the addition of different styles of fluorescent images (0.436 mAP).
180 As more cell-like images added, performance increased further to 0.554 with the full pre-trained
181 dataset.

183 **Shot data scale experiments and ablation experiments**

184 To explore the extent of annotated data required for fine-tuning, we made a shot data scale
185 experiment on these evaluation datasets. We set 10 scale levels, and for each shot number, we
186 randomly sampled 10 times from the training pool to fine-tune the model, followed by testing of
187 transferability. For this evaluation, we focused on Scellseg with CFT and Cellpose with classical
188 fine-tuning because these models clearly outperformed the other three algorithms. For all datasets,
189 we observed that initial performance is relatively low with large variance (Fig. 5a). As the number

190 of shot images increases, the performance improves drastically and variance becomes smaller while
191 Scellseg significantly outperformed Cellpose. For BBBC010_elegans, LIVECell_bv2 and mito
192 dataset, Scellseg with CFT get [2.0%, 4.8%, 18.5%] final improvement respectively compared to
193 [4%, 6%, 12.2%] for Cellpose with classical fine-tuning. We conducted curve fitting using
194 Hyperbola function for each method per dataset for further inspection of the transferability across
195 different shot numbers. The results show that, when increasing shot number, different methods
196 converged to different values and the mAP converged differently across datasets. For mito dataset,
197 mAP persistently increased while for BBBC010_elegans, the rate of its convergence is relatively
198 fast, and whatever the dataset is, performance plateaued at eight shots. Similar results were obtained
199 using the mean Aggregated Jaccard Index²⁷ as a means to evaluate transfer performance (Extended
200 Data Fig. 5). Therefore, it is suggested that at least of eight images is required to achieve satisfied
201 transfer learning based on generalized model.

202 To verify the function of our contrastive fine-tuning strategy, we conducted ablation
203 experiments. Importantly, our contrastive fine-tuning strategy outperformed Scellseg using the
204 classic fine-tuning method at different shot-number experiments on all three evaluation datasets
205 (Fig. 5b). Due to the similar performance of the model when “only” trained on Smicro, we
206 conducted the same shot data scale experiments and ablation experiments as Scellseg pre-trained
207 on the Sgeneral dataset and, excitingly, our style-aware pipeline worked and again outperformed
208 Cellpose with classic fine-tuning strategy (Extended Data Fig. 6).

209 **Graphical user interface**

211 To facilitate Scellseg accessible for scientists without coding experience, we designed a GUI
212 (Fig. 6) with three functional parts, i) view and draw, ii) fine-tune, and iii) inference. For basic
213 annotations, users can modify the mask of an instance directly at single-pixel resolution without
214 deleting the whole mask. We also developed a cell list management system to help users edit the
215 corresponding mask and provide annotations, thereby allowing to provide ground truth for
216 segmentation and cell class prediction simultaneously. Furthermore, users can easily save or load
217 cell lists.

218 In the second module, users can fine-tune the pre-trained model with their own manually
219 labelled data. The system allows users to choose a pre-trained model from Scellseg, Cellpose and
220 HoVer. Furthermore, each model can be combined with either the contrastive or classic fine-tuning
221 strategy, presented above. It will not only give biologists and pathologists more flexibility and
222 versatility for their image analysis tasks, but also help algorithm engineers to easily conduct
223 experiments to study such pre-trained and fine-tuning pipelines. Finally, users can use the fine-
224 tuned model to conduct inference either for one image or use batch inference. After annotation or
225 inference, users can also acquire images of each single instance for further analysis.

227 **Discussion**

229 Accurate cell instance segmentation is still a challenging task for many laboratories. Although
230 generalization models have been developed, these typically require large annotated datasets, which
231 is time- and labor-consuming in data collection, particularly when a large number of segmented
232 objectives are supposed to be covered. To augment data utilization, we firstly established a pipeline
233 for the fine-tuning of pre-trained cell segmentation algorithms. On this basis, we proposed a style-
234 aware pipeline, yielding the best transferability on three different benchmarking datasets.
235 Specifically, the work achieved three main innovations: Firstly, we refined the architecture of
236 Cellpose through introducing attention mechanisms and hierarchical information, making the
237 model more sensitive to different styles. Secondly, we implemented a contrastive fine-tuning

strategy to leverage the information from both unlabelled and pre-trained data based on contrastive learning, which have also achieved great success in other deep learning applications²⁸⁻³². Finally, we organized three benchmarking datasets containing three levels of cell images for further use in segmentation algorithm development.

Importantly, after fine-tuning, AP for BBBC010_elegans and mito can reach up to about 0.9 at a threshold of 0.5, more than 37% and 36% improvements respectively, resulting in model performance that is generally acceptable for researchers to conduct reliable downstream analysis. Our benchmarking showed that among four different instance representations, topological maps generated by the Cellpose model constituted the best way to introduce rich instance information. Results could be further improved by our style-aware pipeline, which exhibited the best transferability on all three evaluation datasets, indicating that introducing such style relevant information can benefit the fine-tuning process. Notably, while generalization ability overall declined after fine-tuning to a specific task, our contrastive fine-tuning considerably improved generalizability. Furthermore, emerging methods like continual learning³³ are also worth to investigate in this context.

In the pre-trained dataset scale experiments, we observed that transferability of Scellseg-CFT increased with the richness of the pre-trained dataset, suggesting that our Scellseg-CFT pipeline can also benefit from large-scale and high-diversity datasets. Notably, the generalization ability of Scellseg pre-trained with Sfluor containing different styles of fluorescent images, outperformed all other models on the mito evaluation dataset, indicating that model pre-training on diverse but specialized data may yield greater performance than both low-diversity specialized dataset (such as Sneuro) or high-diversity generalized datasets (such as Scell that also contains nonfluorescent images). However, we did not observe similar phenomena on transferability. We also noted that the generalization ability increased with the addition of more different microscopy instances beyond cells to other non-cell instances like *C. elegans*, again demonstrating the success of our style-aware pipeline.

For the shot data scale experiment, it is not surprising that performance increases along with shot number. However, what excited us is that we observed a large payoff when increasing shot number from 1 to 3, whereas performance plateaued after approximately eight shots. These results are of high practical relevance as they indicate that the annotation of only about eight images is sufficient to yield a sufficiently fine-tuned model. Few-shot³⁴, one-shot³⁵ and zero-shot³⁶ learning strategies can be studied to further reduce the number of annotated images needed. Notably, at small shot numbers, different shot data can have very large impacts on the fine-tuning process, whereas we observed that as the shot number increases variance becomes substantially smaller. In the future, active learning³⁷ on cell instance segmentation promises to refine shot data selection for fine-tuning.

In this work, we did not research the influence of basic model backbone and all models were based on convolutional neural networks (CNN). In recent years, self-attention architectures (such as Transformer³⁸) have shown great success and there have been studies attempting to apply them to computer vision³⁹. Such transformer architectures have better expressive ability but require more data for accurate training. Nevertheless, we believe that such approaches will eventually provide an important improvement in computer vision compared to CNN.

By integrating attention mechanisms and hierarchical information for style-aware segmentation with a contrastive fine-tuning strategy, Scellseg features the highest transferability when benchmarked on three diverse imaging datasets against currently used segmentation methods. Scellseg optimizes cell and object recognition in diverse microscopy data and, combined with an easy-to-use GUI, can make advanced parallelized segmentation accessible also to researchers and histologists without coding experience. Moreover, the Extractor and Tasker design can facilitate the adaption to other computer vision tasks, such as segmentation and simultaneous class

prediction²⁵, or conducting feature extraction for phenomics analysis⁴⁰. We anticipate Scellseg will serve not only for cell segmentation, but also other for a wide range of other applications in cell biology and biomedicine.

Methods

The code was written in Python programming language v.3.7.4. All experiments were conducted on NVIDIA GeForce RTX 2080Ti. The deep learning framework used Pytorch⁴¹ v.1.7.1.

Datasets

Pre-training datasets. We used the Cellpose dataset published by Stringer et al¹¹ which contains a total of 608 images and over 70,000 segmented instances. 540 images were used as training set (the last of every 8 images was chosen as validation set) and 68 images were used as test set. Here, the whole training set (also named as Sgeneral) was used to pre-train the models and the test set was used to evaluate the generalization ability in Extended Data Fig. 4. Furthermore, a subset of the training set containing a total of seven styles of images with five images per style was used as the contrast data (Extended Data Fig. 2).

Evaluation datasets. Three datasets were used to evaluate the transferability of different models, here called BBBC010_elegans, LIVECell_bv2 and mito. BBBC010_elegans was downloaded from the Broad Bioimage Benchmark Collection⁴², containing 100 images of *C. elegans* in a screen to find novel anti-infectives. There are two phenotypes in this dataset, for worms treated with ampicillin, they appear curved in shape and smooth in texture, while untreated worms appear rod-like in shape and slightly uneven in texture. Only the brightfield channel was used. We discarded images with heavily crossed instances because it is not the focus of our work, the problem may be solved by some special postprocessing algorithm²⁰ or introducing the z-axis information when designing the ground truth. Finally, 49 images were reserved, 10 were used as the training set and 39 were used for testing.

The LIVECell_bv2 dataset²² consists of 536 phase-contrast images and over 330,000 segmented instances. These images were achieved using label-free phase-contrast imaging and cells in this dataset have small spherical morphology and are homogeneous across populations. Of the available images, 386 were used as the training set and 152 were used for testing.

We also generated a novel dataset called Mito dataset, which consisted of 49 fluorescent images of mitochondria from high content screening studies. The images were acquired by ImageXpress Micro Confocal (Molecular Devices). Each image contains two distinct channels, a nuclear channel stained with Hoechst-33342 (Sigma) and a mitochondria channel stained with tetramethylrhodamine methyl ester (TMRM, Sigma). All these images were manually annotated by a single human operator (D.J.X.), 10 images were used as training set and 39 were reserved for testing. Because there was no clear boundary between individual cells, the Mito dataset was used to compare the performance of different algorithms regarding mitochondrial segmentation at the single cell level.

All these three datasets were organized in Cellpose format. The summary information can be seen in Extended Data Table. 1.

Models

When training a cell instance segmentation model, we usually provide raw images and the corresponding masks which label the individual instances with different positive integers per image.

332 Although different values can represent different instances in these masks, it is impractical to
333 directly predict such masks because the max value in each mask represents the number of instances,
334 which are different across images. Thus, the model has to pre-set a very large shape of last conv
335 unit in Pytorch⁴¹ tensor shape format to cover all instances. However, such approaches can result
336 in inefficient memory usage and may not learn well in such a high dimensionality. It is challenging
337 to find an excellent representation of instances and until now there have been four main methods:
338 U-Net²¹, U-Net³²⁴, HoVer²⁵ and Cellpose¹¹.

339 For the classic U-Net model (usually called as U-Net2), we directly map the annotated masks
340 to 2-classes, zero represents background and one represents instance. This method usually performs
341 poorly on touching cells because instance information was completely discarded. In 2018, Fidel et
342 al²⁴ introduced cell borders as the third class to make the network notice the original gap between
343 cells (usually called as U-Net3), they yield a significant improvement compared with U-Net2. In
344 2019, Simon et al²⁵ further developed the model on multiple independent multi-tissue histology
345 image datasets. For each cell per image, they generated horizontal and vertical distance maps to
346 bring in rich instance information when inference, marker-controlled watershed⁴³ was used as the
347 postprocessing to create the final masks. In 2020, Stringer et al¹¹ generated topological maps
348 through a process of simulated diffusion from masks, and when on a test image, they used gradient
349 tracking⁴⁴ to recover individual cells.

350 Here, we wanted to compare how expressive power different methods can provide, so we used
351 the same architecture as Cellpose, only changed the final shape of convolutional layer, loss function
352 and postprocessing to adapt to each method. Scellseg model adopted the representation of Cellpose
353 because the best performance in the experiments. Except the different parts of Scellseg, all other
354 architecture sets were same as Cellpose too.

355 Two-channel 224×224 images were set as input for all 5 models in this work. The primary
356 channel contains instances to segment and the second optional channel can provide extra
357 information such as nuclei channel to support model learning. The hierarchical level of Conv Units
358 was set as [32, 64, 128, 256]. We computed the style embeddings through applying the average
359 pooling on feature map of last Conv Unit and the dimensionality of each level style embeddings
360 after being concatenated in upsampling pass is [256, 384, 448, 480].

361 **Pre-train segmentation models**

363 **Pre-train different models with Sgeneral.** We trained five models (U-Net2, U-Net3, HoVer,
364 Cellpose, Scellseg) with Sgeneral, which contains totally 540 images, 64 of which were reserved
365 for validation.

366 For U-Net2 and U-Net3, a learning rate of 0.002 was selected to achieve good model
367 convergence. For HoVer and Scellseg, the loss function is same as Cellpose, which was defined as:

$$368 \quad L_{segmentation} = BCE(y_{b,2}, lbl_{b,0}) + 0.5 \times MSE(y_{b,0:2}, 5 \times lbl_{b,1:3}) \quad (1)$$

369 Where BCE represents the binary cross-entropy loss, MSE represents the mean square error loss, y
370 represents the ultimate output of model, lbl represents the ground truth, and subscripts corresponded
371 to respective dimensions in y or lbl , b represents the batch size, which we here set to 8.

372 All models were trained for 500 iterations with stochastic gradient descent, the mean diameter
373 was set to 30, all other training hyper-parameters were same as Cellpose.

374 **Pre-training Scellseg for pre-trained dataset scale experiments.** We trained four other
375 models across different subsets of Cellpose mentioned above: Sneuro, Sfluor, Scell, Smicro. For
376 each subset, the last of every 8 images was reserved for validations (11, 36, 47, 56, respectively).
377 All other training hyper-parameters were the same as for Scellseg pre-trained with Sgeneral.

378 Training logs of all models are shown in Extended Data Fig. 7.

379 380 **Fine-tune segmentation models**

381 **Classic fine-tuning strategy (FT).** When fine-tuning, batch size was set to 8, epoch was set
382 to 100, the optimizer was Adam, the initial learning rate was set to 0.001 and every quarter of
383 epochs it was reduced by 50%. Before being fed to the network, the image-mask pairs were resized,
384 randomly rotated and reshaped with the ultimate shape of input as (8, 2, 224, 224).

385 **Contrastive fine-tuning strategy (CFT).** The contrast loss function was defined as:

$$386 L_{contrast} = \frac{MSE(shot, query)}{MSE(shot, contrast) + 10^{-5}} \quad (2)$$

387 Where MSE represents the mean square error and was used to compute the difference between
388 embeddings, and 10^{-5} was added to prevent divisions by zero. This contrast loss was added into the
389 segmentation loss function during contrastive fine-tuning, so the final loss function was defined as:

$$390 L_{total} = L_{segmentation} + L_{contrast} \times Sigmoid(\alpha) \quad (3)$$

391 Where α is a scalar to control the weight of contrast loss, also learnt during the fine-tuning process.
392 A sigmoid function was used to assure that the coefficient of contrast loss changed smoothly
393 between zero and one. The initial α in the contrast loss function was set to 0.2, initial learning rate
394 of α was set to 0.1, with reductions by 50% every quarter of epochs. For query and contrast data,
395 they were resized, randomly cropped, randomly rotated and reshaped before being fed to the
396 network with identical input shapes. Other parameters were the same as for the classic fine-tuning
397 strategy.

398 For both classic and contrastive fine-tuning strategies, we fine-tuned all layers because this
399 method performed best compared with downsampling part or the whole extractor (Extended Data
400 Fig. 8). For each dataset, we computed the instance diameter using shot data without using the
401 automated method provided by Cellpose, which was used in resizing the current mean diameter of
402 instances to the mean diameter used for model pre-training.

403 404 **Benchmarking**

405 **Metrics.** We used average precision (AP) and the Aggregated Jaccard index (AJI) to evaluate
406 segmentation performance (See ref. ¹¹ for detailed definitions). Except in Fig. 2c-e, we averaged
407 the AP or AJI over IoU from 0.50 to 0.95 with a step size of 0.05 for convenient comparison and
408 reserving the overall performance information simultaneously as detailed below:

$$409 mAP = (AP_{0.50} + AP_{0.55} + \dots + AP_{0.90} + AP_{0.95}) / 10 \quad (4)$$

$$410 mAJI = (AJI_{0.50} + AJI_{0.55} + \dots + AJI_{0.90} + AJI_{0.95}) / 10 \quad (5)$$

411 **Shot data scale experiments.** We set a total of 10 scale levels; for BBBC010_elegans and
412 mito we used [1, 2, 3, 4, 5, 6, 7, 8, 9, 10] and for LIVECell_bv2 [1, 2, 4, 8, 16, 32, 64, 128, 256,
413 386]. For each shot number experiment, we randomly sampled 10 times from the training set to
414 fine-tune the pre-trained model. To eliminate issues due to different training data, the random state
415 was kept identical across models. For example, we sampled the 9th, 5th, and 2nd image from the
416 total of 10 images in the training set of a 3-shot experiment for the mito dataset, and then used the
417 same images as training data for all five models.

418 419 **Statistical Analysis**

420 All figures were made using GraphPad PRISM 8.0 software (GraphPad Software, Inc., CA,
421 USA). All graphs display mean values, and the error bars represent the standard deviation (SD).
422 Statistical analyses were conducted with two-way repeated measures analysis of variance (ANOVA)
423 followed by Sidak's multiple comparisons test in Fig. 2c-e and two-way ANOVA in Fig. 5a. A
424 nonlinear regression curve fit was performed using a hyperbolic function in Fig. 5a, given as:

$$425 \quad Y = \frac{B_{max} * X}{K_d + X} \quad (6)$$

426 Where B_{max} and K_d are constants.
427
428

429 Data availability

430 Three evaluation datasets used in this work will be made available at [https://scellseg-
431 data.s3.cn-northwest-1.amazonaws.com.cn/evaluation_datasets.zip](https://scellseg-data.s3.cn-northwest-1.amazonaws.com.cn/evaluation_datasets.zip) upon publication.
432
433

434 Code availability

435 The Scellseg software, including detailed tutorial, can be freely available at GitHub
436 (<https://github.com/cellimnet/scellseg-publish>).
437
438

439 References

- 440 1. Cheng, S., Fu, S., Kim, Y. M., Song, W., Li, Y., Xue, Y., Yi, J. & Tian, L. Single-cell
441 cytometry via multiplexed fluorescence prediction by label-free reflectance microscopy. *Sci.*
442 *Adv.* **7**, eabe0431 (2021).
- 443 2. Petukhov, V., Xu, R. J., Soldatov, R. A., Cadinu, P., Khodosevich, K., Moffitt, J. R. &
444 Kharchenko, P. V. Cell segmentation in imaging-based spatial transcriptomics. *Nat.*
445 *Biotechnol.* (2021).
- 446 3. Moen, E., Bannon, D., Kudo, T., Graf, W., Covert, M. & Van Valen, D. Deep learning for
447 cellular image analysis. *Nat. Methods.* **16**, 1233–1246 (2019).
- 448 4. Pan, C., Schoppe, O., Parra-Damas, A., Cai, R., Todorov, M. I., Gondi, G., von Neubeck, B.,
449 Böğürçü-Seidel, N., Seidel, S., Sleiman, K., Veltkamp, C., Förstera, B., Mai, H., Rong, Z.,
450 Trompak, O., Ghasemigharagoz, A., Reimer, M. A., Cuesta, A. M., Coronel, J., Jeremias, I.,
451 Saur, D., Acker-Palmer, A., Acker, T., Garvalov, B. K., Menze, B., Zeidler, R. & Ertürk, A.
452 Deep learning reveals cancer metastasis and therapeutic antibody targeting in the entire body.
453 *Cell.* **179**, 1661-1676.e19 (2019).
- 454 5. Rohban, M. H., Abbasi, H. S., Singh, S. & Carpenter, A. E. Capturing single-cell
455 heterogeneity via data fusion improves image-based profiling. *Nat. Commun.* **10**, 2082
456 (2019).
- 457 6. Chai, J., Zeng, H., Li, A. & Ngai, E. W. T. Deep learning in computer vision: A critical
458 review of emerging techniques and application scenarios. *Machine Learning with*
459 *Applications.* **6**, 100134 (2021).
- 460 7. Tareef, A., Song, Y., Huang, H., Feng, D., Chen, M., Wang, Y. & Cai, W. Multi-pass fast
461 watershed for accurate segmentation of overlapping cervical cells. *IEEE Trans. Med.*
462 *Imaging.* **37**, 2044–2059 (2018).

- 463 8. Krasowski, N. E., Beier, T., Knott, G. W., Kothe, U., Hamprecht, F. A. & Kreshuk, A. Neuron
464 segmentation with high-level biological priors. *IEEE Trans. Med. Imaging.* **37**, 829–839
465 (2018).
- 466 9. Loewke, N. O., Pai, S., Cordeiro, C., Black, D., King, B. L., Contag, C. H., Chen, B., Baer,
467 T. M. & Solgaard, O. Automated cell segmentation for quantitative phase microscopy. *IEEE*
468 *Trans. Med. Imaging.* **37**, 929–940 (2018).
- 469 10. Caicedo, J. C., Goodman, A., Karhohs, K. W., Cimini, B. A., Ackerman, J., Haghghi, M.,
470 Heng, C., Becker, T., Doan, M., McQuin, C., Rohban, M., Singh, S. & Carpenter, A. E.
471 Nucleus segmentation across imaging experiments: the 2018 Data Science Bowl. *Nat.*
472 *Methods.* **16**, 1247–1253 (2019).
- 473 11. Stringer, C., Wang, T., Michaelos, M. & Pachitariu, M. Cellpose: a generalist algorithm for
474 cellular segmentation. *Nat. Methods.* **18**, 100–106 (2021).
- 475 12. Chu, B., Madhavan, V., Beijbom, O., Hoffman, J. & Darrell, T. Best practices for fine-tuning
476 visual classifiers to new domains. *Proc. Eur. Conf. Comput. Vis.* 435–442 (2016).
- 477 13. Guo, Y., Shi, H., Kumar, A., Grauman, K., Rosing, T. & Feris, R. SpotTune: transfer learning
478 through adaptive fine-tuning. *Proc. IEEE Conf. Comput. Vis. Pattern Recog.* 4805–4814
479 (2019).
- 480 14. Cetinic, E., Lipic, T. & Grgic, S. Fine-tuning convolutional neural networks for fine art
481 classification. *Expert Syst. Appl.* **114**, 107–118 (2018).
- 482 15. You, K., Kou, Z., Long, M. & Wang, J. Co-Tuning for transfer learning. *Conference and*
483 *Workshop on Neural Information Processing Systems.* 17236–17246 (2020).
- 484 16. Howard, J. & Ruder, S. Universal language model fine-tuning for text classification. *Annual*
485 *Meeting of the Association for Computational Linguistics.* 328–339 (2018).
- 486 17. Lee, C., Cho, K. & Kang, W. Mixout: effective regularization to finetune large-scale
487 pretrained language models. *International Conference on Learning Representations.* (2020).
- 488 18. Zaki, G., Gudla, P. R., Lee, K., Kim, J., Ozbun, L., Shachar, S., Gadkari, M., Sun, J., Fraser,
489 I. D. C., Franco, L. M., Misteli, T. & Pegoraro, G. A Deep Learning Pipeline for Nucleus
490 Segmentation. *Cytom. Part A.* **97**, 1248–1264 (2020).
- 491 19. Jin, Y., Toberoff, A. & Azizi, E. Transfer learning framework for cell segmentation with
492 incorporation of geometric features. Preprint at
493 <https://biorxiv.org/lookup/doi/10.1101/2021.02.28.433289> (2021).
- 494 20. Wählby, C., Kamensky, L., Liu, Z. H., Riklin-Raviv, T., Conery, A. L., O’Rourke, E. J.,
495 Sokolnicki, K. L., Visvikis, O., Ljosa, V., Irazoqui, J. E., Golland, P., Ruvkun, G., Ausubel,
496 F. M. & Carpenter, A. E. An image analysis toolbox for high-throughput *C. elegans* assays.
497 *Nat. Methods.* **9**, 714–716 (2012).
- 498 21. Ronneberger, O., Fischer, P. & Brox, T. U-Net: Convolutional Networks for Biomedical
499 Image Segmentation. *Medical Image Computing and Computer-Assisted Intervention.* **9351**,
500 234–241 (2015).
- 501 22. Edlund, C., Jackson, T. R., Khalid, N., Bevan, N., Dale, T., Dengel, A., Ahmed, S., Trygg, J.
502 & Sjögren, R. LIVECell—A large-scale dataset for label-free live cell segmentation. *Nat.*
503 *Methods.* **18**, 1038–1045 (2021).
- 504 23. Maaten, Lvd. & Hinton, G. Visualizing data using t-SNE. *J. Machine Learning Res.* **9**, 2579–
505 2605 (2008).

- 506 24. Guerrero-Pena, F. A., Marrero Fernandez, P. D., Ing Ren, T., Yui, M., Rothenberg, E. &
507 Cunha, A. Multiclass Weighted Loss for Instance Segmentation of Cluttered Cells. Proc.
508 IEEE Int. Conf. Image Process. 2451–2455 (2018).
- 509 25. Graham, S., Vu, Q. D., Raza, S. E. A., Azam, A., Tsang, Y. W., Kwak, J. T. & Rajpoot, N.
510 Hover-Net: Simultaneous segmentation and classification of nuclei in multi-tissue histology
511 images. *Med. Image Anal.* **58**, 101563 (2019).
- 512 26. Yu, W., Lee, H. K., Hariharan, S., Bu, W. Y. & Ahmed, S.
513 <https://doi.org/10.7295/W9CCDB6843>.
- 514 27. Kumar, N., Verma, R., Sharma, S., Bhargava, S., Vahadane, A. & Sethi, A. A Dataset and a
515 Technique for Generalized Nuclear Segmentation for Computational Pathology. *IEEE Trans.*
516 *Med. Imaging.* **36**, 1550–1560 (2017).
- 517 28. Cho, I., Huo, Y. & Yoon, S. Weakly-supervised contrastive learning in path manifold for
518 monte carlo image reconstruction. *ACM Trans. Graph.* **40**, 1-14 (2021).
- 519 29. Gunel, B., Du, J., Conneau, A. & Stoyanov, V. Supervised contrastive learning for pre-
520 trained language model fine-tuning. *Proc. Int. Conf. Learn. Represent.* (2021).
- 521 30. Yu, Y., Zuo, S., Jiang, H., Ren, W., Zhao, T. & Zhang, C. Fine-tuning pre-trained language
522 model with weak supervision: a contrastive-regularized self-training approach. Annual
523 Meeting of the Association for Computational Linguistics. Proceedings of the 2021
524 Conference of the North American Chapter of the Association for Computational Linguistics:
525 Human Language Technologies, 1063–1077 (2021).
- 526 31. Tian, Y., Sun, C., Poole, B., Krishnan, D., Schmid, C. & Isola, P. What Makes for Good
527 Views for Contrastive Learning? Conference and Workshop on Neural Information
528 Processing Systems. (2020).
- 529 32. Chen, T., Kornblith, S., Norouzi, M. & Hinton, G. A Simple Framework for Contrastive
530 Learning of Visual Representations. *Proceedings of Machine Learning Research.* **119**, 1597-
531 1607 (2020).
- 532 33. Zeng, G., Chen, Y., Cui, B. & Yu, S. Continual learning of context-dependent processing in
533 neural networks. *Nat. Mach. Intell.* **1**, 364–372 (2019).
- 534 34. Das, D. & Lee, C. S. G. A two-stage approach to few-shot learning for image recognition.
535 *IEEE Trans. Image Process.* **29**, 3336–3350 (2020).
- 536 35. Michaelis, C., Ustyuzhaninov, I., Bethge, M. & Ecker, A. S. One-Shot Instance Segmentation.
537 Preprint at <https://arxiv.org/abs/1811.11507> (2018).
- 538 36. Xu, X., Tsang, I. W. & Liu, C. Complementary Attributes: A New Clue to Zero-Shot
539 Learning. *IEEE T. Cybern.* **51**, 12 (2021).
- 540 37. Zhou, Z., Shin, J. Y., Gurudu, S. R., Gotway, M. B. & Liang, J. Active, continual fine tuning
541 of convolutional neural networks for reducing annotation efforts. *Med. Image Anal.* **71**,
542 101997 (2021).
- 543 38. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł. &
544 Polosukhin, I. Attention is all you need. Conference and Workshop on Neural Information
545 Processing Systems. 6000–6010 (2017).
- 546 39. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S. & Guo, B. Swin Transformer:
547 Hierarchical Vision Transformer using Shifted Windows. Preprint at
548 <https://arxiv.org/abs/2103.14030> (2021).

- 549 40. Cuccarese, M. F., Earnshaw, B. A., Heiser, K., Fogelson, B., Davis, C. T., McLean, P. F.,
550 Gordon, H. B., Skelly, K.-R., Weathersby, F. L., Rodic, V., Quigley, I. K., Pastuzyn, E. D.,
551 Mendivil, B. M., Lazar, N. H., Brooks, C. A., Carpenter, J., Probst, B. L., Jacobson, P.,
552 Glazier, S. W., Ford, J., Jensen, J. D., Campbell, N. D., Statnick, M. A., Low, A. S., Thomas,
553 K. R., Carpenter, A. E., Hegde, S. S., Alfa, R. W., Victors, M. L., Haque, I. S., Chong, Y. T.
554 & Gibson, C. C. Functional immune mapping with deep-learning enabled phenomics applied
555 to immunomodulatory and COVID-19 drug discovery. Preprint at
556 <https://biorxiv.org/content/10.1101/2020.08.02.233064v2> (2020).
- 557 41. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z.,
558 Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M.,
559 Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J. & Chintala, S. PyTorch: an
560 imperative style, high-performance deep learning library. Conference and Workshop on
561 Neural Information Processing Systems. (2019).
- 562 42. Ljosa, V., Sokolnicki, K. L. & Carpenter, A. E. Annotated high-throughput microscopy
563 image sets for validation. *Nat. Methods.* **9**, 637–637 (2012).
- 564 43. Cheng J. & Rajapakse, J. C. Segmentation of clustered nuclei with shape markers and
565 marking function. *IEEE Trans. Biomed. Eng.* **56**, 741–748 (2009).
- 566 44. Li, G., Liu, T., Nie, J., Guo, L., Chen, J., Zhu, J., Xia, W., Mara, A., Holley, S. & Wong, S.
567 T. C. Segmentation of touching cell nuclei using gradient flow tracking. *J. Microsc.* **231**, 47–
568 58 (2008).
- 569 45. Oktay, O., Schlemper, J., Folgoc, L. L., Lee, M., Heinrich, M., Misawa, K., Mori, K.,
570 McDonagh, S., Hammerla, N. Y., Kainz, B., Glocker, B. & Rueckert, D. Attention U-Net:
571 Learning pytttWhere to Look for the Pancreas. The International Conference on Medical
572 Imaging with Deep Learning. (2018).

573
574

575 **Acknowledgments**

576 The authors are grateful for the support from ZJU PII-Molecular Devices Joint Laboratory
577 and support from "Medicine + X" interdisciplinary Center of Zhejiang University.

578 **Funding**

579 Y.W. is supported by National Key R&D Program of China (2021YFC1712905), National
580 Natural Science Foundation of China (No. 82173941), the Innovation Team and Talents Cultivation
581 Program of National Administration of Traditional Chinese Medicine (No. ZYYCXTD-D-202002).
582 R.W. is supported by National Natural Science Foundation of China (No. 61872319) and Natural
583 Science Foundation of Zhejiang Provincial (No. LR18F020002).

584 **Author contributions**

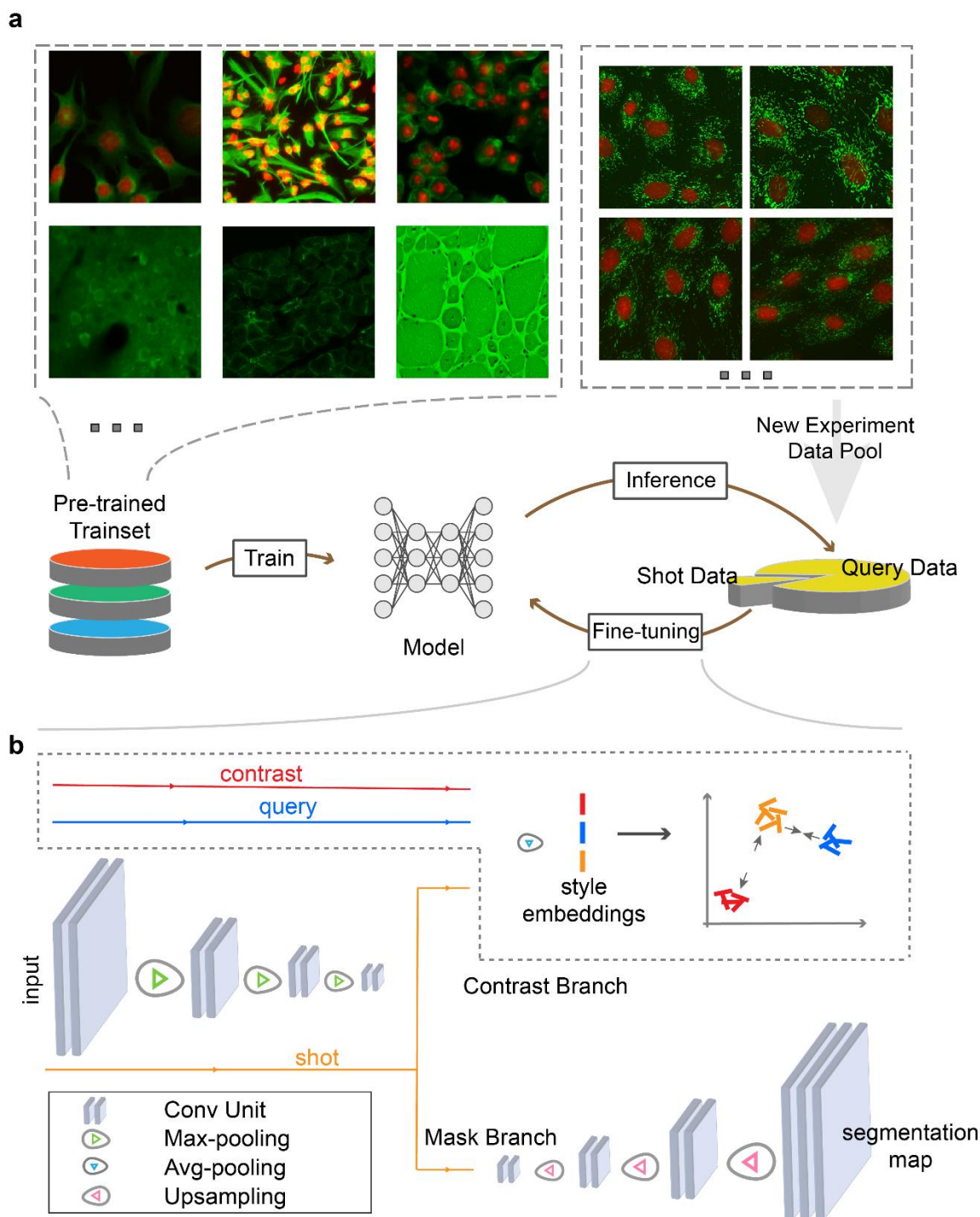
585 Y.W. and R.W. proposed the concept and supervised the overall project. D.J.X. established
586 the pipeline, organized datasets, written the code of pipeline, conducted experiments, performed
587 data analysis, designed and finished the code writing of graphical user interface. D.H.C. participated
588 in part of the code writing of graphical user interface. D.J.X., V.M.L., Y.W., R.W., Y.T.Z.
589 participated in the preparation of the manuscript.

590 **Competing interests**

591 The authors declare no competing interests.

592

593 **Figures and Tables**



594 **Fig. 1 | Pipeline of pre-trained and fine-tuning strategy.** **a**, Overview of fine-tuning a
 595 pre-trained model for a new experiment. The shot data means hand-labelled data while
 596 query data means unlabelled data. **b**, Diagram of the proposed contrastive fine-tuning
 597 strategy. The contrast data is a subset of pre-trained data. The network is a representation
 598 of U-Net family, the detailed architecture of our proposed model is shown in Extended
 599 Data Fig. 1. Different colored lines and arrows mark the flow of data.

600

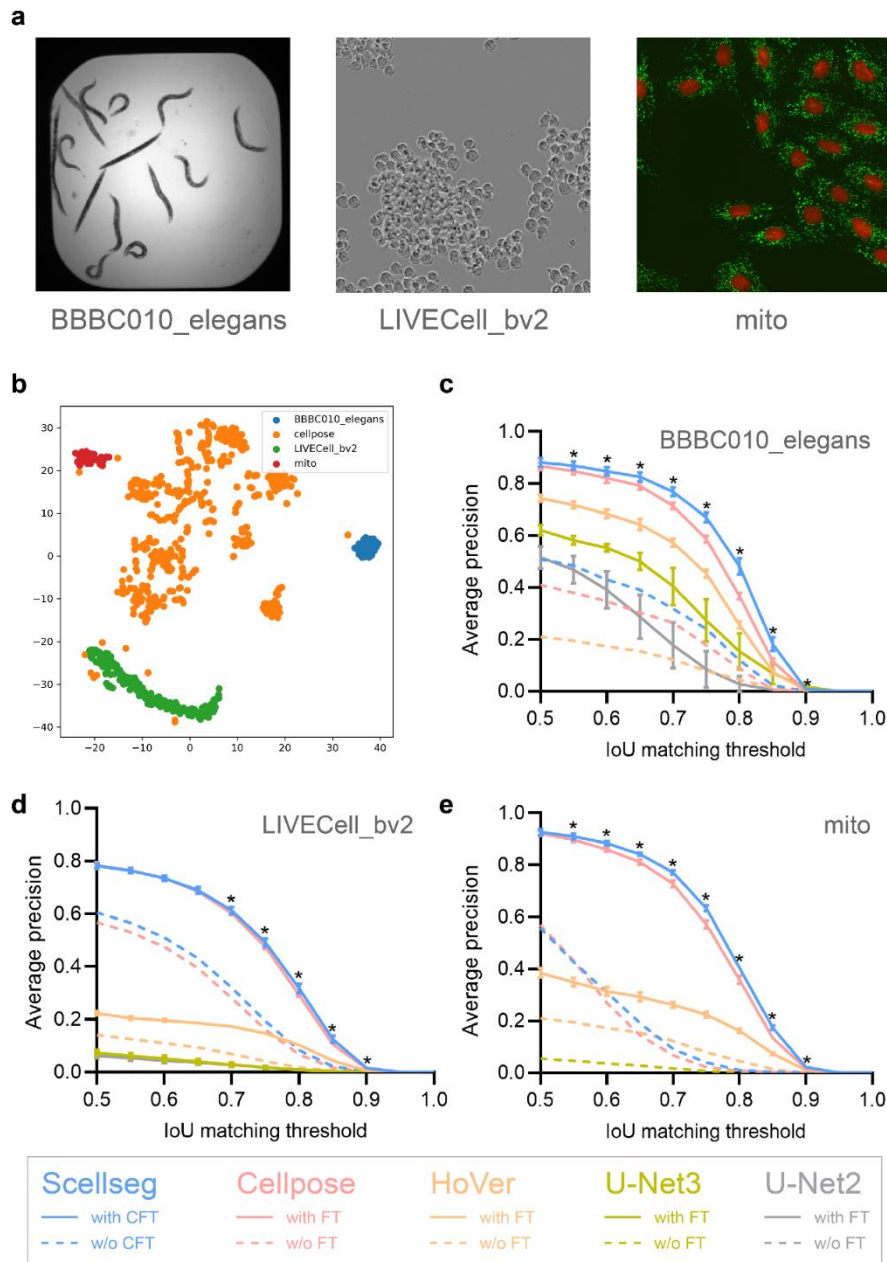
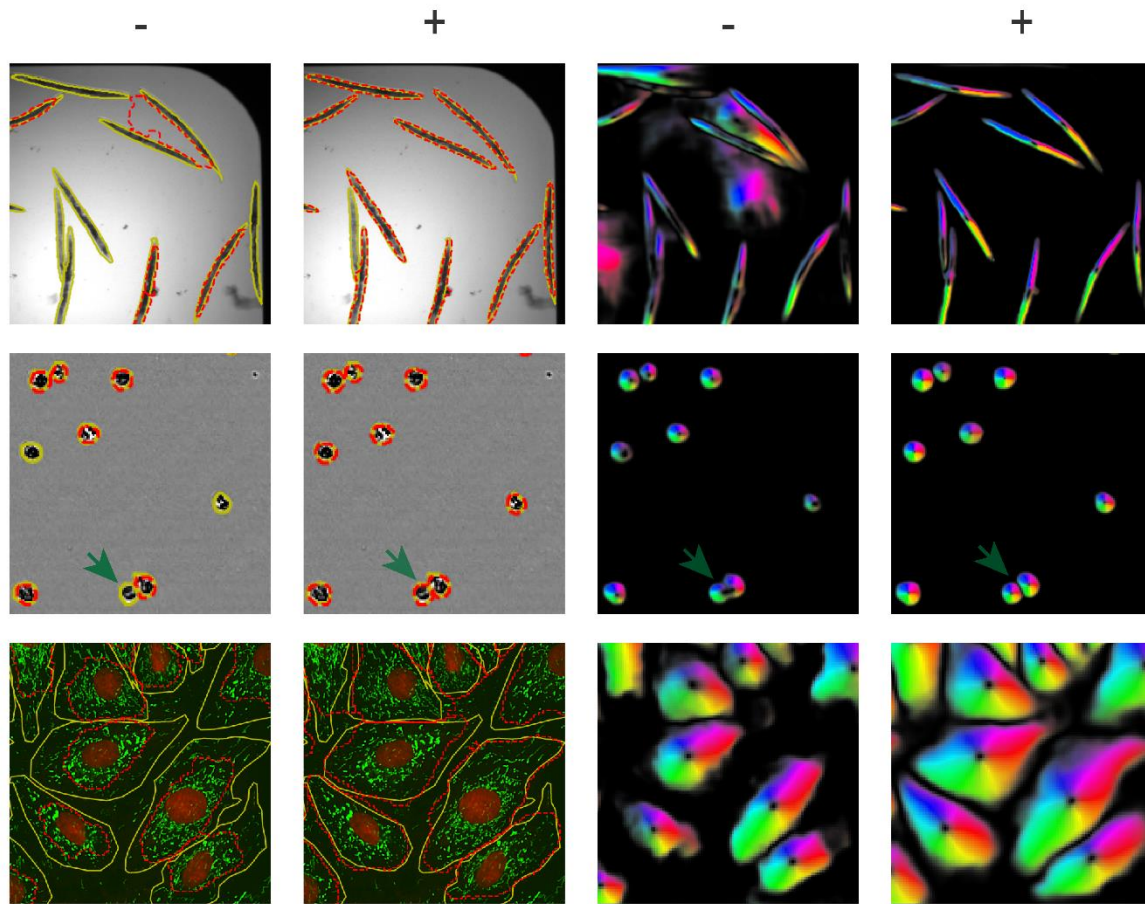


Fig. 2 | Transferability of Scellseg with contrastive fine-tuning strategy on three evaluation datasets. **a**, Example images of three datasets. **b**, Visualization for style embeddings of three datasets and pre-trained dataset using t-SNE. **c-e**, Performance of different models on BBBC010_elegans (**c**), LIVECell_bv2 (**d**) and mito (**e**) dataset. Different colors correspond to different models, the dotted lines denote the performance of applying models directly and the solid lines denote the performance after fine-tuning. For Scellseg, we use contrastive fine-tuning (CFT) and for others is classic fine-tuning strategy (FT). We did not plot the line which corresponding performance is less than 0.01. Each pre-trained and fine-tuning pipeline was conducted 10 times at various random states, error bars represent the mean \pm SD. * means P -value <0.05 , determined by two-way ANOVA followed by Sidak's multiple comparisons test for Scellseg with CFT and Cellpose with FT.



614 **Fig. 3 | Example Scellseg segmentation results with and without contrastive fine-**
615 **tuning on three evaluation datasets.** The right two columns show the direct topological
616 maps outputted by Scellseg model. The left two columns show ultimate masks of different
617 datasets, the ground truth masks are shown in yellow solid line, and the predicted masks
618 are shown in dotted red line. Symbol “-” represents results of applying models directly and
619 symbol “+” represents results after contrastive fine-tuning. Green arrows emphasize the
620 segmentation of adherent cells.
621

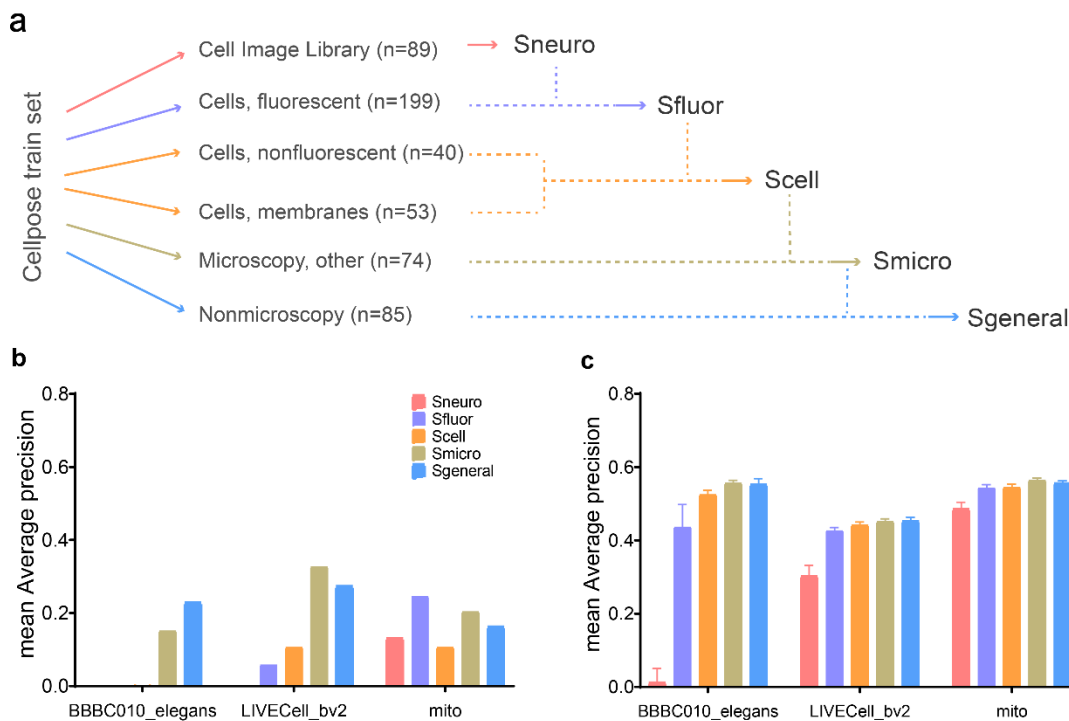


Fig. 4 | Pre-trained dataset scale experiments. **a**, Composition of different subsets from Cellpose train set. **b**, Generalization ability of different pre-trained Scellseg models on three evaluation datasets. Generalization ability means the performance of employing pre-trained model directly. **c**, Transferability of different pre-trained Scellseg models on three evaluation datasets. Transferability means the performance of employing the pre-trained model after fine-tuning. Each pre-trained and fine-tuning pipeline was conducted 10 times at various random states, error bars represent the mean \pm SD.

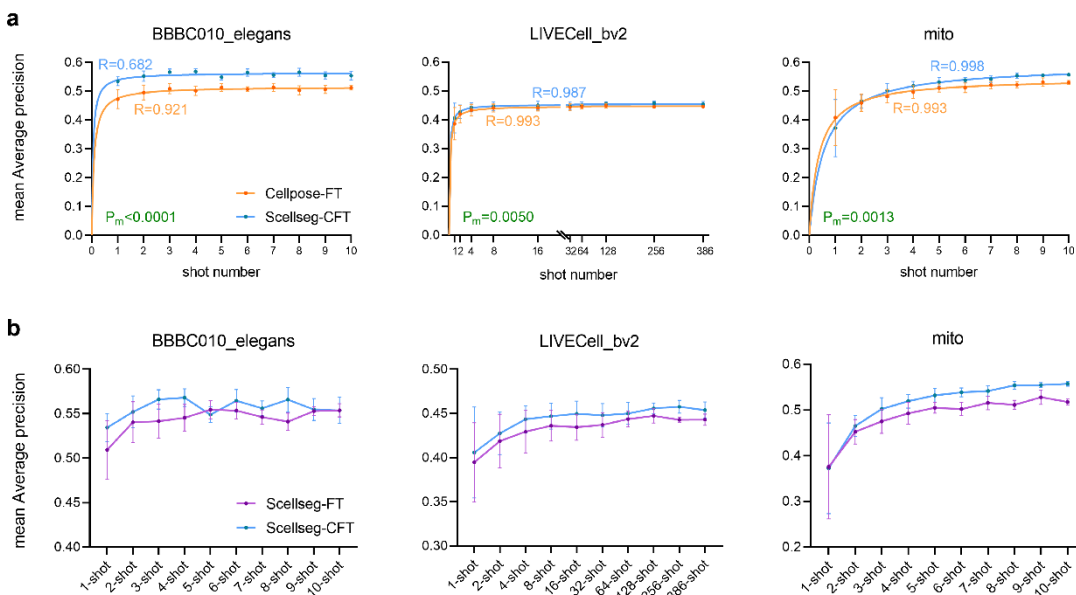


Fig. 5 | Shot data scale experiments (a) and ablation experiments (b). Each pre-trained and fine-tuning pipeline was conducted 10 times at various random states, error bars represent the mean \pm SD. **a**, Performance of Cellpose-FT and Scellseg-CFT on three evaluation datasets at 10 shot data scales. Cellpose-FT represents Cellpose with classic fine-tuning and Scellseg-CFT represents Scellseg with contrastive fine-tuning strategy. Another metric to evaluate the models is shown in Extended Data Fig. 5. We performed nonlinear regression based on hyperbola function and corresponding R-value of fitted curve is plotted in the picture. A two-way ANOVA analysis was conducted for group comparison of Scellseg-CFT and Cellpose-FT per dataset and corresponding Pm-value is plotted in the picture. Pm-value <0.05 was considered the performance between Scellseg-CFT and Cellpose-FT is significant. **b**, Ablation experiments for contrastive fine-tuning strategy. Scellseg-FT represents Scellseg with classic fine-tuning strategy, data of Scellseg-CFT is completely same as Scellseg-CFT in (a).

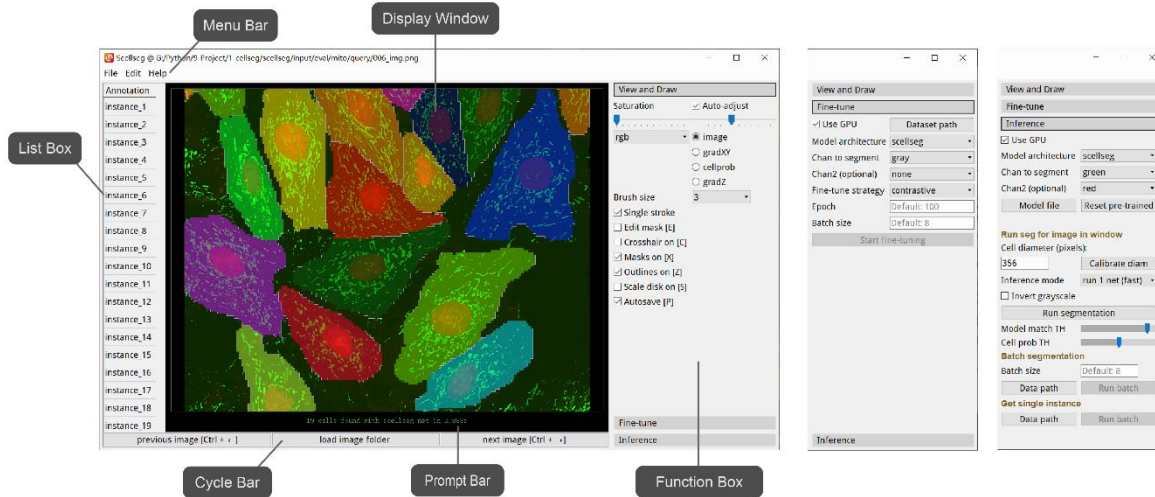
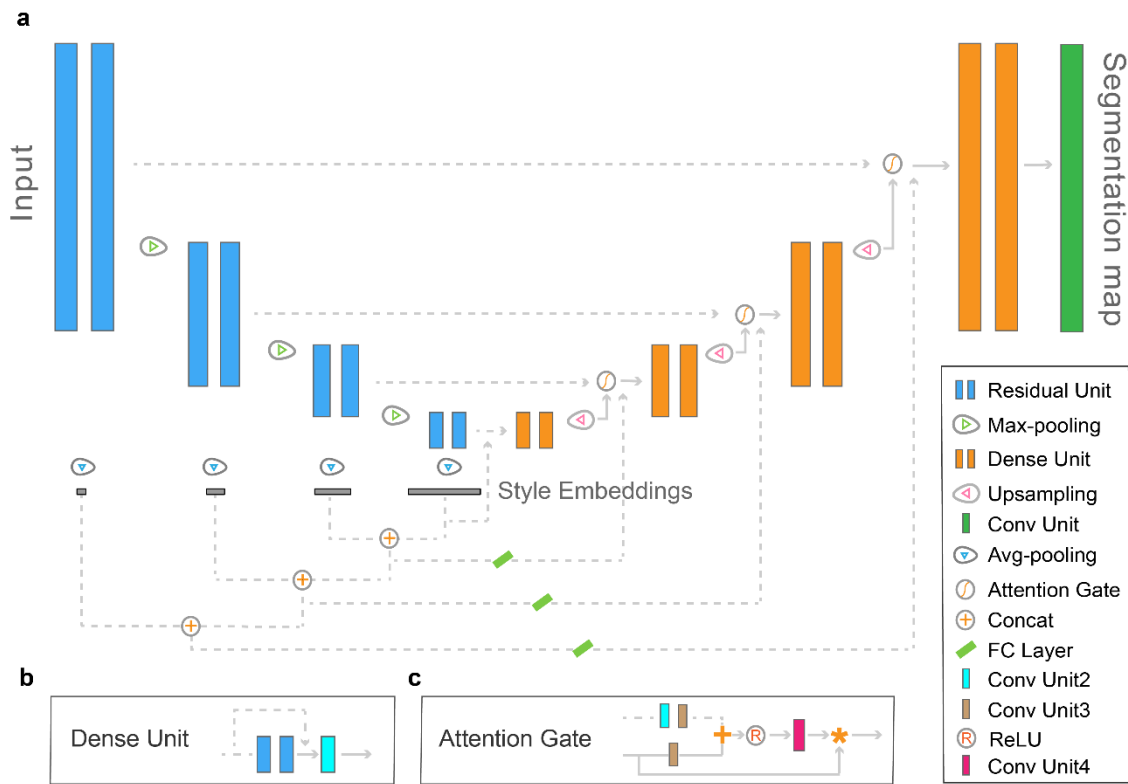


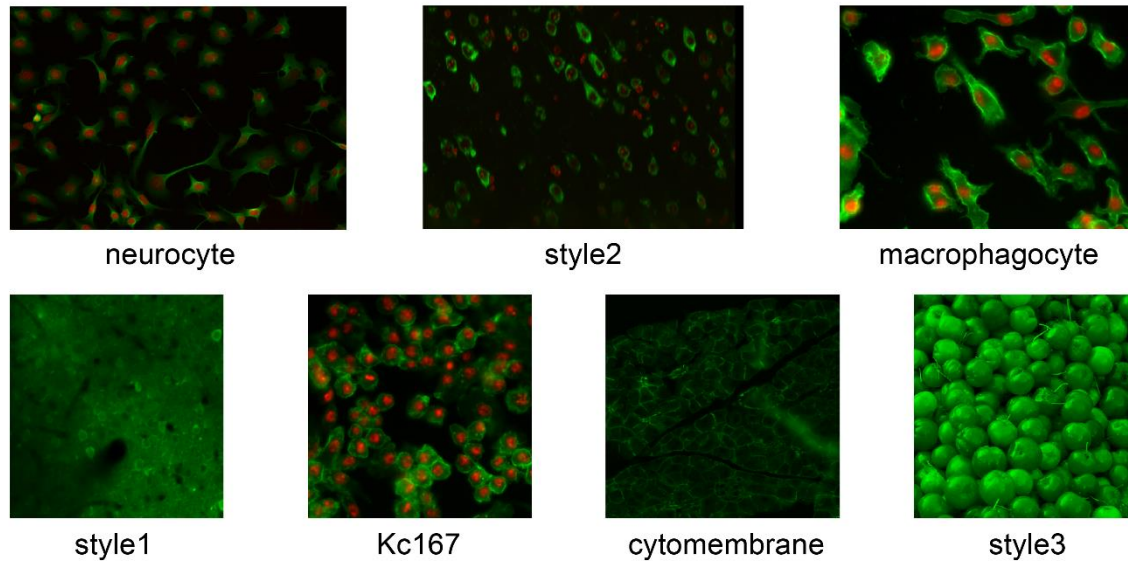
Fig. 6 | Graphical user interface (GUI). This GUI contains six modules: Menu Bar, List Box, Display Window, Prompt Bar, Cycle Bar and Function Box. There are three main functions: “View and draw”, “Fine-tune” and “Inference”.

644
645
646
647

648 **Extended Data**



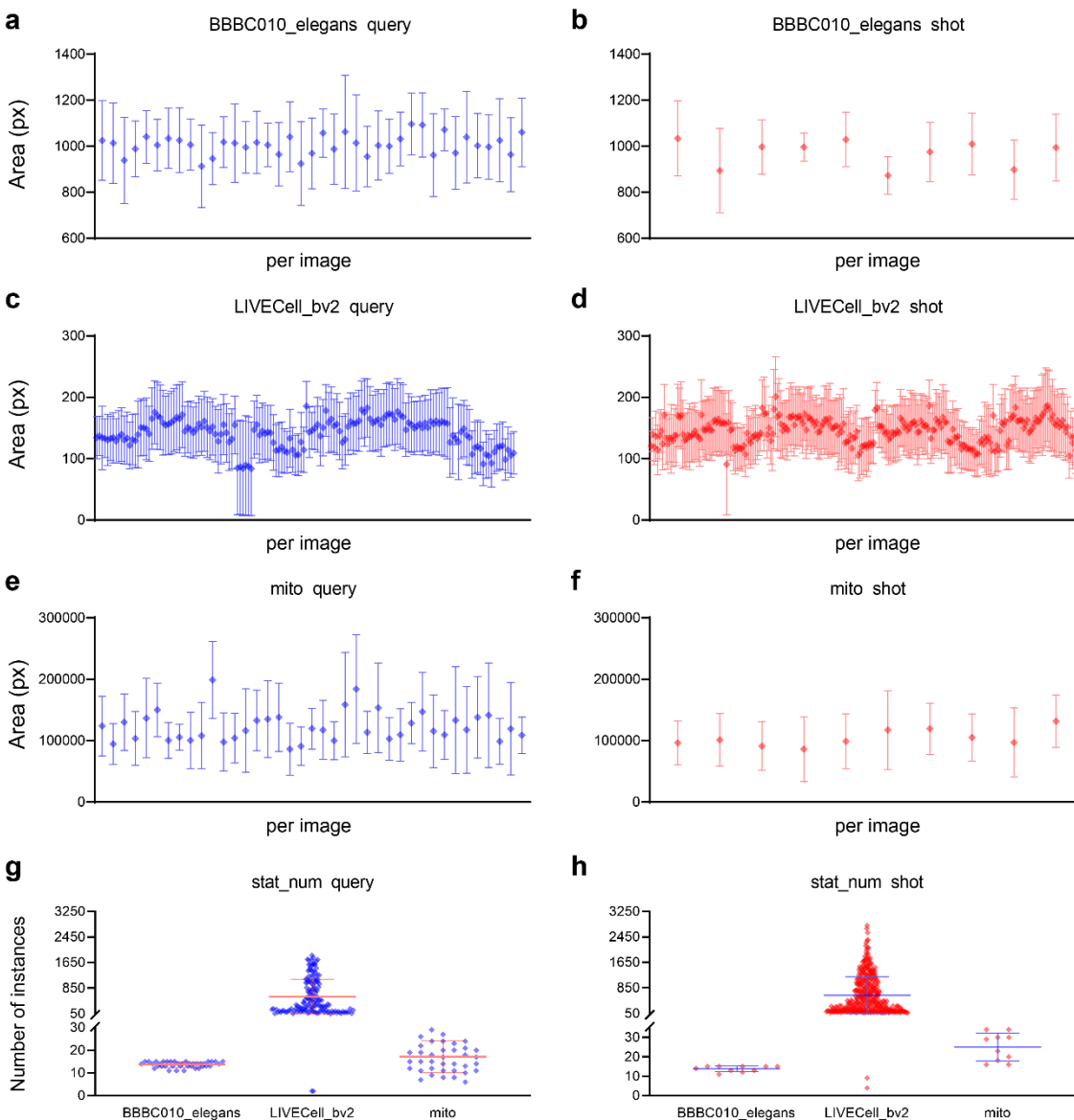
649 **Extended Data Fig. 1 | Architecture of the proposed Scellseg.** **a**, The Residual Unit
 650 refers to Cellpose¹¹, Dense Unit refers to Hover-Net²⁵, Attention Gate refers to Attention
 651 U-Net⁴⁵, Conv Unit represents a set of operations including BatchNorm2d, ReLU and
 652 Conv2d in Pytorch⁴¹. The blue parts (including Max-pooling operations) were called as
 653 downsampling pass, the dot lines (including operations marked on them) were called as
 654 concatenation part, specially, the orange parts (including Avg-pooling operations) were
 655 called as upsampling pass, and the last green Conv Unit was named as Tasker,
 656 downsampling, upsampling together with concatenation part were named as Extractor.
 657 Input images are progressively encoded and decoded to get the ultimate segmentation
 658 map. Style embeddings of each scale is obtained by using global average pooling on
 659 respective convolutional map. **b**, Detail architecture of Dense Unit. Conv Unit2 represents
 660 a set of operations including Conv2d, BatchNorm2d and ReLU. **c**, Detail architecture of
 661 Attention Gate. Conv Unit3 represents a set of operations including Conv2d and
 662 BatchNorm2d. Symbol “+” represents tensor plus and Symbol “*” represents tensor
 663 multiplication. Conv Unit4 represents a set of operations including Conv2d, BatchNorm2d
 664 and Sigmoid.
 665



Extended Data Fig. 2 | Example image per style of Contrast data. There are totally 7 styles of images we used in our contrastive fine-tuning strategy, each style includes five images.

666
667
668
669

670



671

672

673

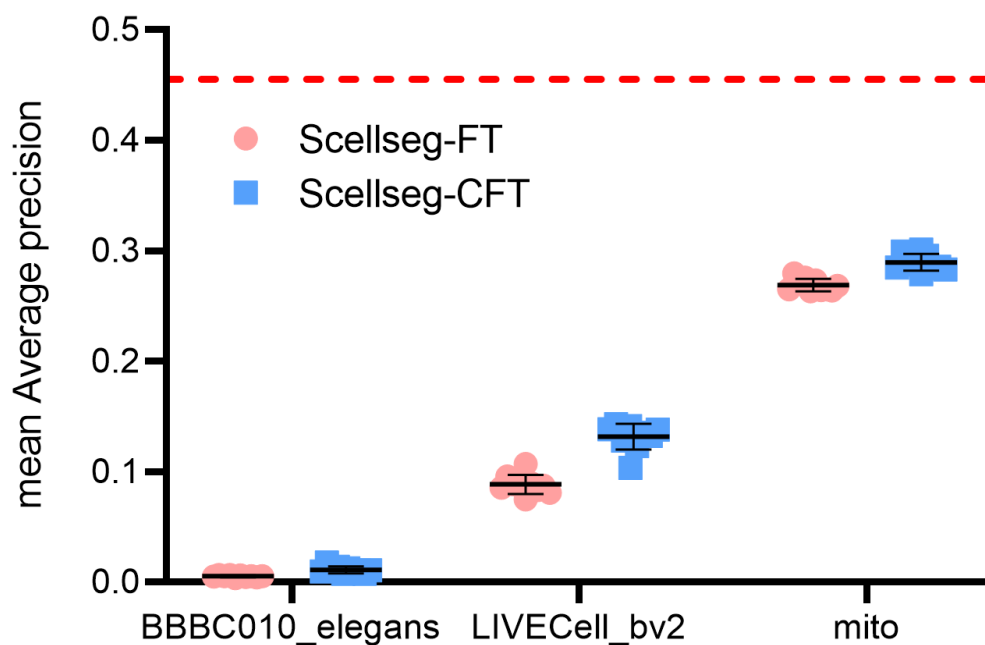
674

675

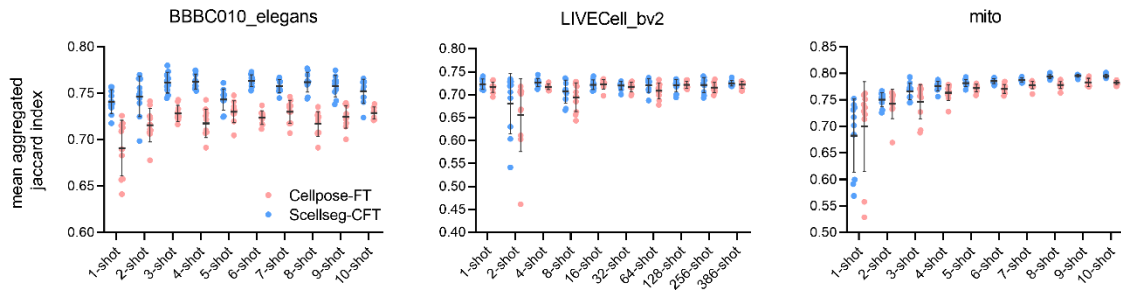
676

677

Extended Data Fig. 3 | Statistics of three evaluation datasets. a-f, Distribution of cell areas in px for each image in query or shot data of BBBC010_elegans (**a**, **b**), LIVECell_bv2 (**c**, **d**), mito (**e**, **f**) dataset, error bars represent the mean \pm SD. **g-h**, Distribution of number of instances for each image in query data (**g**) or shot data (**h**) of three datasets, each dot represents one image in corresponding dataset, error bars represent the mean \pm SD.

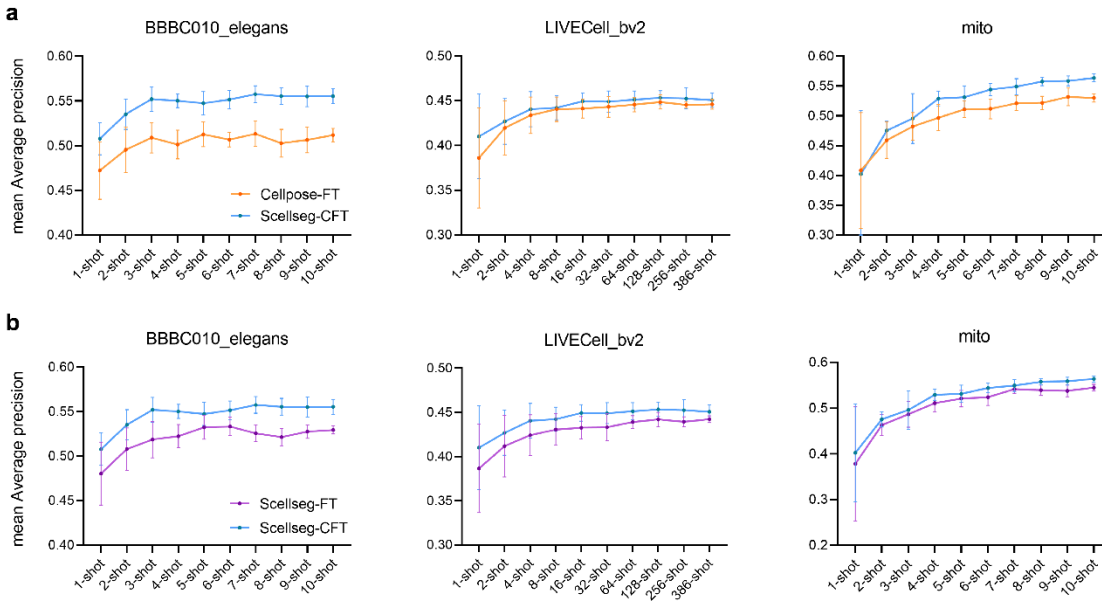


678 **Extended Data Fig. 4 | Generalization ability on Cellpose test set of Scellseg with**
679 **different fine-tuning strategies.** Scellseg-FT represents Scellseg with classic fine-tuning
680 and Scellseg-CFT represents Scellseg with contrastive fine-tuning strategy. Red dot line
681 represents employing Scellseg directly on test set. Each pre-trained and fine-tuning
682 pipeline was conducted 10 times at various random states, error bars represent the mean \pm
683 SD.
684

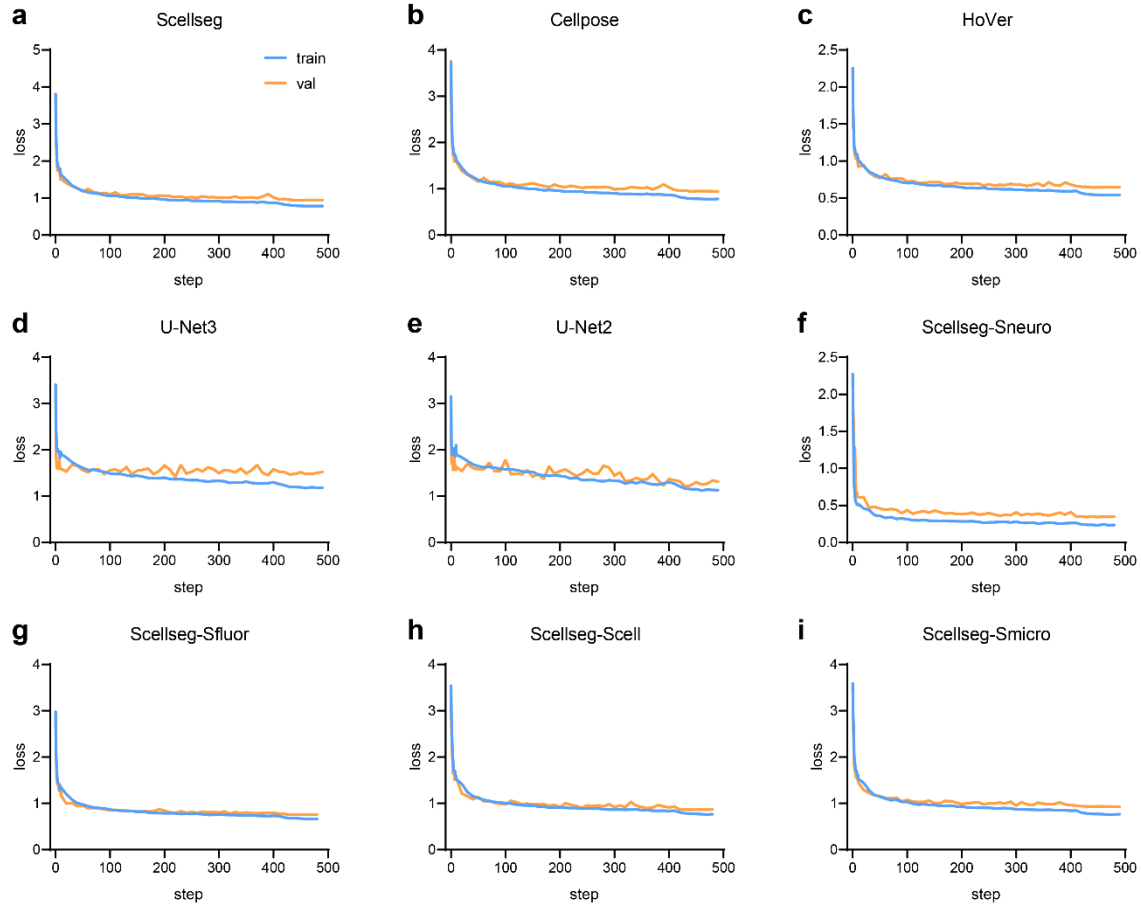


Extended Data Fig. 5 | Using mean Aggregated Jaccard Index metric to evaluate segmentation performance in shot data scale experiments. Error bars represent the mean \pm SD.

685
686
687
688

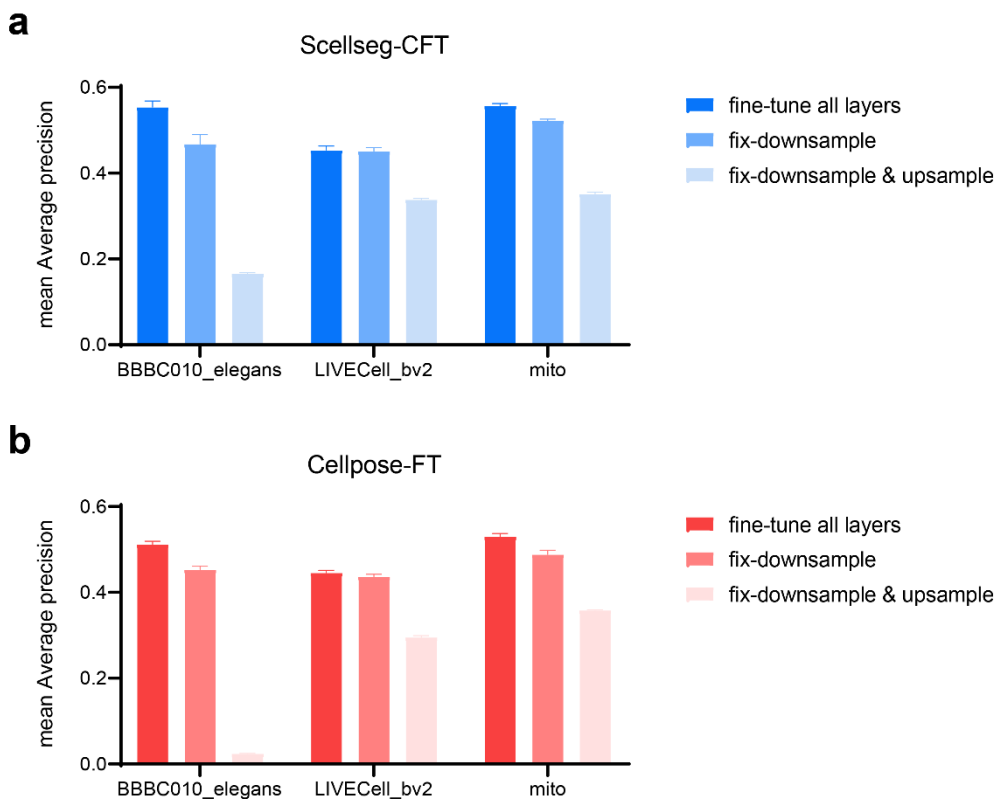


689 **Extended Data Fig. 6 | Shot data scale experiments (a) and ablation experiments (b)**
690 **on Scellseg pre-trained with Smicro.** Each pre-trained and fine-tuning pipeline was
691 conducted 10 times at various random states, error bars represent the mean \pm SD. Data of
692 Scellseg-CFT in (a) is completely same as Scellseg-CFT in (b).
693



Extended Data Fig. 7 | Train logs of models in the paper.

694
695



696 **Extended Data Fig. 8 | Performance of different fine-tuning methods.** We compared
697 three kinds of fine-tuning methods for Scellseg-CFT (**a**) and Cellpose-FT (**b**) on three
698 evaluation datasets, respectively are fine-tuning all layers of the model, fixing the
699 downsampling layers and fixing downsampling-upsampling layers. Each fine-tuning
700 method was conducted 10 times at various random states, error bars represent the mean \pm
701 SD.
702

Dataset	Number of images in train set	Number of images in test set	Number of instances in test set	Acquisition Modality	Shape of images	Characteristics
BBBC010_elegans	10	39	670	Bright Field	696*520*1	worms appear rod-like or curved in shape
LIVECell_bv2	386	152	533	Fluorescent, 10X	704*520*1	small spherical morphology, homogeneous population
mito	10	39	89821	Phase-contrast imaging, 60X	2048*2048*3	mitochondria scatter around the nuclei, no clear boundaries
Total	406	230	91024			

703
704
705

Extended Data Table. 1 | Summary statistics of three evaluation datasets.