

1 **CelloType: A Unified Model for Segmentation and Classification of Tissue Images**

2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46

Minxing Pang<sup>1</sup>, Tarun Kanti Roy<sup>2</sup>, Xiaodong Wu<sup>3,4</sup>, Kai Tan<sup>5,6,7,\*</sup>

<sup>1</sup>Applied Mathematics & Computational Science Graduate Group, University of Pennsylvania, Philadelphia, PA, USA

<sup>2</sup>Department of Computer Science, The University of Iowa, Iowa City, IA, USA

<sup>3</sup>Department of Electrical and Computer Engineering, The University of Iowa, Iowa City, IA, USA

<sup>4</sup>Department of Radiation Oncology, University of Iowa, Iowa City, IA, USA

<sup>5</sup>Department of Pediatrics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA

<sup>6</sup>Division of Oncology and Center for Childhood Cancer Research, Children's Hospital of Philadelphia, Philadelphia, PA, USA

<sup>7</sup>Center for Single Cell Biology, Children's Hospital of Philadelphia, Philadelphia, PA, USA

\*Corresponding Author Contact:

Kai Tan, [tank1@chop.edu](mailto:tank1@chop.edu)

4004 CTRB

3501 Civic Center Blvd

Philadelphia, PA, 19104

(267) 425-0050

47 **Abstract**

48

49 Cell segmentation and classification are critical tasks in spatial omics data analysis. We  
50 introduce CelloType, an end-to-end model designed for cell segmentation and classification of  
51 biomedical microscopy images. Unlike the traditional two-stage approach of segmentation  
52 followed by classification, CelloType adopts a multi-task learning approach that connects the  
53 segmentation and classification tasks and simultaneously boost the performance of both tasks.  
54 CelloType leverages Transformer-based deep learning techniques for enhanced accuracy of  
55 object detection, segmentation, and classification. It outperforms existing segmentation methods  
56 using ground-truths from public databases. In terms of classification, CelloType outperforms a  
57 baseline model comprised of state-of-the-art methods for individual tasks. Using multiplexed  
58 tissue images, we further demonstrate the utility of CelloType for multi-scale segmentation and  
59 classification of both cellular and non-cellular elements in a tissue. The enhanced accuracy and  
60 multi-task-learning ability of CelloType facilitate automated annotation of rapidly growing  
61 spatial omics data.

## 62 Introduction

63  
64 Recent advancements in spatial omics technologies have markedly improved our ability to  
65 analyze intact tissues at the cellular level, revealing unparalleled insights into the link between  
66 cellular architecture and functionality of various tissues and organs<sup>1</sup>. Collaborative efforts, such  
67 as the Human Tumor Atlas Network<sup>2</sup>, the Human Biomolecular Atlas Program<sup>3</sup>, and the BRAIN  
68 initiative, are leveraging these technologies to map spatial organizations of various types of  
69 healthy and diseased tissues. With the anticipated surge in spatial omics data, there is a pressing  
70 need for sophisticated computational tools for data analysis. A typical analysis workflow of  
71 spatial omics data begins with cell segmentation. Following cell segmentation and quantification  
72 of molecular analytes, cell type annotation is the next critical, albeit often time-consuming task  
73 before further analysis can proceed. Conventional analysis pipelines perform these two tasks  
74 sequentially, typically using the segmentation results as the inputs for the classification task. As  
75 representatives of state-of-the-art segmentation methods, Mesmer<sup>4</sup> uses a convolutional neural  
76 network (CNN)<sup>5</sup> backbone and a Feature Pyramid Network with the watershed algorithm for  
77 both nuclear and cell segmentation. Cellpose<sup>6</sup> and Cellpose2<sup>7</sup> use a CNN with a U-net<sup>8</sup>  
78 architecture to predict the gradient of topological map. A gradient tracking algorithm is then used  
79 to obtain the segmentation mask. For cell classification task, CellSighter<sup>9</sup> employs CNN to  
80 predict cell types based on segmentation masks and the tissue images. CELESTA<sup>10</sup> uses an  
81 iterative algorithm to assign cell types based on quantified cell-by-protein matrix.

82  
83 Despite achieving satisfactory performance in certain tissues, conventional approaches have  
84 several limitations. First and foremost, the reliance of cell classification models on segmentation  
85 results hampers their ability to leverage the full spectrum of semantic information present in  
86 tissue images. In fact, these two tasks are interconnected. Segmentation can enhance focus on  
87 relevant signals, thus mitigating noise and enabling more precise learning of class features for  
88 classification. Conversely, information specific to classes aid in the segmentation process, as the  
89 unique texture and morphology of certain object types can enhance segmentation accuracy.  
90 Second, the two-step approach is computationally inefficient, requiring separate training for each  
91 task. Third, the performance of existing segmentation methods also varies significantly across  
92 different tissue types, suggesting substantial room for improvement. Moreover, to our  
93 knowledge, existing methods do not offer a confidence assessment for the segmentation task.

94  
95 Deep learning, especially through the use of CNNs, has gained popularity in biomedical image  
96 analysis, especially in segmentation<sup>11</sup> and classification<sup>9</sup>. Mesmer, for example, has notably  
97 improved cell segmentation accuracy using CNN. However, recent developments in computer  
98 vision has shown that Transformer-based models<sup>12</sup>, such as the Detection Transformer (DETR)  
99<sup>13</sup> and the Detection transformer with Improved deNoising anchOr (DINO)<sup>14</sup>, significantly  
100 outperform CNN-based models in object detection. These Transformer-based models have also  
101 shown superior performance in instance segmentation of histological images<sup>15</sup>. Despite these  
102 breakthroughs, the application of Transformer-based models to cell/nuclear segmentation in  
103 multiplexed images and other spatial omics data type remains unexplored. A unified framework,  
104 MaskDINO<sup>16</sup>, which integrates object detection and segmentation, has shown superior  
105 performance across diverse datasets for multi-class instance segmentation. However, its effective  
106 ness has only been tested on RGB images of natural objects. This leaves a significant gap in  
107 applying Transformer-based models to multiplexed tissue images, which present greater

108 challenges due to their larger number of imaging channels, varying shapes of tightly  
109 apposed/overlapping cellular and non-cellular elements.

110  
111 The limitations of current methodologies and the advent of novel deep learning techniques  
112 motivated us to develop CelloType, an end-to-end method for joint cell segmentation and  
113 classification. CelloType employs a Transformer-based deep neural network architecture with  
114 multiple branches to handle object detection, segmentation, and classification concurrently. We  
115 benchmarked the performance of CelloType against state-of-the-art methods using a variety of  
116 public image datasets, including single-channel, and multiplexed fluorescent tissue and cell  
117 images and bright-field images of nature objects. We further demonstrated a novel feature of  
118 CelloType for multi-scale segmentation and classification to delineate both cellular and  
119 noncellular elements in tissue images.

120

## 121 **Results**

122

### 123 **Overview of CelloType**

124

125 CelloType is a deep neural network (DNN)-based framework (Figure 1) designed for joint multi-  
126 scale segmentation and classification of a variety of biomedical microscopy images, including  
127 multiplexed molecular images, histological images, and bright-field images. The core of  
128 CelloType's functionality begins with the extraction of multi-scale image features through the  
129 use of a Swin Transformer<sup>17</sup>. These features are then fed into the DINO object detection module  
130 that extracts instance-specific latent features and predicts a preliminary object bounding box with  
131 associated class label for each instance. Finally, the MaskDINO segmentation module integrates  
132 the multi-scale image features from the Swin Transformer and DINO outputs to produce the final  
133 refined instance segmentations. The CelloType model is trained using a loss function that  
134 considers segmentation masks, object detection boxes, and classes labels.

135

136 The DINO module's architecture (Figure 1b) includes a Transformer encoder-decoder set-up  
137 with multiple prediction heads. It begins by flattening image features and integrating them with  
138 positional embeddings<sup>18</sup>. By employing a strategy that mixes anchor and content queries, the  
139 module can adapt to various object features. The module refines bounding boxes through a  
140 deformable attention mechanism. A contrastive denoising training (CDN) procedure is used  
141 together with the attention mechanism to improve the robustness of bounding box detection.  
142 Finally, a linear transformation is applied to the denoised bounding box features to predict the  
143 class label of the object.

144

145 CelloType can tackle diverse image analysis tasks including cell/nuclear segmentation, non-  
146 cellular structure segmentation, and multi-scale segmentation (Figure 1c). Different data types  
147 are used to train CelloType for various tasks. For cell or nuclear segmentation, training data  
148 includes one/two-channel images with corresponding cell membrane or nuclear masks. For joint  
149 segmentation and classification, the training data consists of images with segmentation mask,  
150 bounding box, and class label of each object. The images can contain many channels in addition  
151 to the cell membrane and nuclear channels. CelloType is implemented in Python and publicly  
152 available at <http://github.com/tanlabcode/CelloType>.

153

## 154 **Benchmark of cell and nuclear segmentation performance using multiplexed images**

155

156 We first applied CelloType to the TissueNet dataset<sup>4</sup> that includes tissue images generated using  
157 six multiplexed molecular imaging technologies (CODetection by indexing (CODEX)<sup>19</sup>, Cyclic  
158 Immunofluorescence (CycIF)<sup>20</sup>, Imaging Mass Cytometry (IMC)<sup>21</sup>, Multiplexed Ion Beam  
159 Imaging (MIBI)<sup>22</sup>, Multiplexed Immunofluorescence (MxIF)<sup>23</sup>, and Vectra<sup>24</sup>) and six tissue types  
160 (breast, gastrointestinal, immune, lung, pancreas, skin). The images were divided into 2,580  
161 training patches (512 x 512 pixels) and 1,324 test patches (256 x 256 pixels).

162

163 We compared CelloType with two state-of-the-art methods, Mesmer<sup>4</sup> and Cellpose<sup>27</sup>. For object  
164 detection and instance segmentation, we used the Average Precision (AP) metric<sup>25</sup> defined by the  
165 Common Objects in Context (COCO) project and the Intersection over Union (IoU) thresholds  
166 from 0.5 to 0.9 in 0.05 increments (Methods). The precision-IoU curves (Figure 2a) revealed that  
167 CelloType consistently outperformed both Mesmer and Cellpose2 across the entire range of IoU  
168 thresholds on the TissueNet dataset. Additionally, considering that CelloType provides a  
169 confidence score for each segmentation mask and the COCO metric incorporates these  
170 confidence scores in matching predicted and ground truth cell boundaries, we also evaluated a  
171 version of CelloType that outputs confidence scores, CelloType\_C. Overall, performance is  
172 higher for cell segmentation than nuclear segmentation for all methods except for Mesmer. For  
173 cell segmentation, CelloType\_C achieved an average AP of 0.556, significantly surpassing the  
174 basic CelloType (0.450), Cellpose2 (0.354), and Mesmer (0.312). For nuclear segmentation,  
175 CelloType\_C achieved a mean AP of 0.655, outperforming CelloType (0.571), Cellpose2  
176 (0.516), and Mesmer (0.237) by considerable margins. These results underscore CelloType's  
177 superior segmentation accuracy and the added value of incorporating confidence scores.

178

179 To evaluate the effect of imaging technology and tissue type on the segmentation performance,  
180 we next analyzed the mean AP scores stratified by these two factors (Figure 2b). Overall,  
181 performance of all methods is lowest on the IMC data and breast tissue data. CelloType and  
182 CelloType\_C consistently outperformed Mesmer and Cellpose2 across the technology platforms  
183 and tissue types. Figure 2d-e show representative cell and nuclear segmentation results by the  
184 compared methods. These examples illustrate Cellpose2 tends to produce segmentation  
185 boundaries that are larger than the ground truth and thus often under-segmentation. On the other  
186 hand, Mesmer tends to miss more cells or nuclei.

187

## 188 **Benchmark of cell segmentation performance using diverse image types**

189

190 To further evaluate CelloType's performance of cell segmentation across diverse microscopy  
191 images beyond multiplexed fluorescent images, we applied CelloType to the Cellpose Cyto  
192 dataset<sup>6</sup> which include fluorescent, bright-field microscopy images of cells and images of natural  
193 objects. Since most of the images in this dataset contain only one channel and Mesmer was  
194 trained on two-channel image data, we only benchmarked the performance of CelloType,  
195 CelloType\_C, and Cellpose2.

196

197 Across the entire dataset, CelloType\_C achieved an average AP of 0.469, surpassing the  
198 performance of both CelloType (0.368) and Cellpose2 (0.322). This superiority is consistently  
199 observed across 6 diverse image sets (Figure 3b). Figure 3c shows representative segmentation

200 results by Cellpose2 and CelloType for a single-channel image from the “Other microscopy”  
201 category. Consistent with the findings in Figure 2d with multiplexed IMC image, Cellpose2  
202 exhibited a tendency for under-segmentation, while CelloType produced more precise  
203 segmentation boundaries. Additionally, Figure 3d shows the segmentation result for another  
204 single-channel image from the “Non-fluorescent” cell category, where CelloType demonstrated  
205 enhanced accuracy in both identifying the correct number of cells and delineating their  
206 boundaries, in contrast to Cellpose2, which tended to over-segment.

207

## 208 **Joint segmentation and cell type classification of multiplexed images**

209

210 To assess the performance of CelloType for simultaneous cell segmentation and classification,  
211 we applied it to a colorectal cancer CODEX dataset<sup>26</sup>. This dataset consists of 140 images of  
212 tumor tissue sections from 35 patients. Each tissue section was imaged using 56 fluorescent  
213 antibodies plus two nuclear stains, resulting in a total of 58 channels. These images were  
214 processed into 512 x 512 pixels image patches, which were subsequently divided into a training  
215 set of 720 patches and a test set of 120 patches (Supplemental Figure 1). Given the lack of  
216 established methods for simultaneous cell segmentation and classification, we combined  
217 Cellpose2 and CellSighter as a baseline model. This choice was motivated by the reported  
218 superior performance of each method for their respective task.

219

220 Using manual cell type annotation as the ground truth, we computed the AP score at an IoU  
221 threshold of 0.5 (i.e. AP50) for each cell type. CelloType achieved a mean AP50 of 0.84 across  
222 all cell types, markedly exceeding the Cellpose2+CellSighter model’s mean AP of 0.24 (Figure  
223 4a). Furthermore, both CelloType\_C and CellSighter produce a confidence score for their cell  
224 type predictions. To assess the utility of the confidence score, we explored the relationship  
225 between these confidence scores and accuracy of predictions. Notably, CelloType’s confidence  
226 scores demonstrated a strong, nearly linear correlation with prediction accuracy, particularly  
227 within the confidence score range of 0.5 to 0.7. In contrast, the relationship for CellSighter’s  
228 confidence scores appeared flat, indicating a lack of reliable calibration in its confidence  
229 assessment (Figure 4b).

230

231 Figure 4c shows two examples of predictions by CelloType and Cellpose2+CellSighter along  
232 with the ground truth annotations. These predictions encompass cell segmentation masks,  
233 predicted cell types and associated confidence scores. CelloType correctly predicted the  
234 identities of the vast majority of cells of different types with varying morphologies and  
235 abundance. For instance, in the top image, CelloType correctly predicted abundant neoplastic  
236 cells, alongside rare regulatory T cells (Treg), and morphologically irregular macrophages.  
237 Similarly, in the bottom image, CelloType correctly predicted abundant smooth muscle cells and  
238 sparsely distributed CD8+ T cells. In contrast, the Cellpose2+CellSighter model misclassified  
239 several cell types as plasma cells (top image) and granulocytes (bottom image). Moreover, we  
240 found many instances where CellSighter’s predictions, despite being incorrect, were  
241 accompanied by high confidence scores, as indicated by arrows.

242

243 We next evaluated the performance of each component of the Cellpose2+CellSighter model,  
244 focusing on the segmentation function of Cellpose2 and the cell type classification function of  
245 CellSighter. Figure 5a shows the AP-IoU curve for cell segmentation on the colorectal cancer

246 CODEX dataset. CelloType achieved a mean AP of 0.585, significantly exceeding Cellpose2's  
247 mean AP of 0.345. In assessing CellSighter's classification performance, we used the ground  
248 truth segmentation masks as inputs, treating the task purely a classification task. The resulting  
249 confusion matrix revealed the distribution of predictions for each cell type and the accuracy  
250 values displayed along the diagonal (Figure 5b). Furthermore, Figure 5c shows CellSighter's  
251 classification precision for 11 cell types, achieving a mean precision of 0.53, compared to  
252 CelloType's mean AP50 score of 0.81. This comparative analysis underscores CelloType's  
253 superior performance not only as an end-to-end tool for cell type annotation but also in its  
254 individual functions for segmentation and classification, outperforming the two-stage approach  
255 of combining Cellpose2 and CellSighter.

256  
257

### 258 **Multi-scale segmentation and annotation by CelloType**

259

260 Non-cellular components, such as the vasculature, lymphatic vessels, trabecular bone, and extra  
261 cellular matrix, and reticular fibers play important roles in tissue function. These elements are  
262 typically much larger than cells. Moreover, certain cell types like macrophages and adipocytes  
263 are either large or possess irregular shapes. Together, these elements present challenges to  
264 conventional segmentation methods. Furthermore, existing methods are incapable of  
265 simultaneous, multi-scale segmentation of both cellular and non-cellular elements within a tissue  
266 image. To assess the effectiveness of CelloType for multi-scale segmentation and classification,  
267 we applied it to a human bone marrow CODEX dataset<sup>27</sup> (Supplemental Figure 2). This dataset  
268 comprises 12 whole-slide images of bone marrow sections from healthy donors, with each tissue  
269 section imaged using 53 fluorescent antibodies plus one nuclear stain, totaling 54 channels. The  
270 images were divided into 512 x 512 pixels patches with 1600 for training and 400 for testing.  
271 The dataset presents a unique challenge due to the diversity of cell/non-cell types, notably  
272 adipocytes, which are substantially larger than other cell types, and trabecular bone fragments,  
273 which have irregular and complex shapes.

274

275 Using 5-fold cross-validation, we evaluated the performance of CelloType on simultaneous  
276 segmentation and classification of both cell and non-cell elements in the bone marrow, including  
277 small regularly shaped cell types and much larger adipocytes and irregularly shaped trabecular  
278 bone fragments. CelloType achieved average AP50 values of 55.4, 44.3, and 58.9 for adipocytes,  
279 trabecular bone fragments, and the rest of cell types, respectively (Figure 6a). Consistent with  
280 our results with the colorectal cancer CODEX dataset, we observed a strong correlation between  
281 the prediction confidence scores and prediction accuracy (Figure 6b). Figure 6c shows two  
282 representative examples of predictions by CelloType along with the ground truths. In addition to  
283 correctly identifying smaller and regularly shaped cells, CelloType correctly identified most  
284 adipocytes and trabecular bone fragments. This result demonstrates CelloType's efficacy of  
285 analyzing challenging tissue images consisting of tightly packed cells and non-cell elements with  
286 varying sizes and shapes.

287

### 288 **Discussion**

289

290 We present CelloType, an end-to-end method for joint segmentation and classification for  
291 biomedical microscopy images. Unlike existing methods that treat segmentation and

292 classification as separate tasks, CelloType uses a multi-task learning approach. By leveraging  
293 advancements in Transformer-based deep learning techniques, CelloType offers a unified  
294 approach to object detection, segmentation, and classification. It starts with Swin Transformer-  
295 based feature extraction from an image, followed by the DINO object detection and classification  
296 module, which produces latent features and detection boxes that, when combined with the raw  
297 image inputs within the MaskDINO module, culminate in refined instance segmentation and  
298 classification. The shared encoder in the DINO module extracts latent information that is shared  
299 by both tasks, explicitly enhancing the connection between the segmentation and classification  
300 tasks, and simultaneously boosting the performance of both tasks. Moreover, the improved  
301 object detection accuracy of DINO through deformable attention and contrastive denoising  
302 allows the classification task to focus on relevant regions of the image.

303  
304 It should be noted that this work has the following limitations. First, CelloType requires training  
305 for segmentation and classification tasks. In terms of segmentation, there is a rapid growth of  
306 training data, exemplified by resources like TissueNet and Cellpose Cyto databases. Models that  
307 are pre-trained on these public datasets are readily transferable to new images, provided that they  
308 contain nuclear and/or membrane channels. However, for classification, training data is  
309 considerably more limited. As a result, pretrained CelloType classification model cannot be  
310 readily applied to new images unless there is a substantial overlap of cell/structure types between  
311 the training and testing images. To mitigate this need for training data for classification,  
312 methodologies such as few-shot learning<sup>28</sup>, self-supervised learning, and contrastive learning<sup>29</sup>  
313 can be incorporated into the CelloType framework. Additionally, with the rapid growth of spatial  
314 omics data, it is anticipated that high-quality tissue annotations will also grow quickly.  
315 Consequently, CelloType's pre-training process can be broadened to include a wider array of  
316 datasets, thereby facilitating its application in automated annotation of common tissue types.

317  
318 Spatial transcriptomics technologies can profile hundreds to thousands of genes at single-cell  
319 resolution, yielding a much larger number of features compared to spatial proteomics  
320 technologies such as CODEX which typically can only profile fewer than a hundred proteins.  
321 This substantial increase in the feature space, coupled with the distinct spatial distribution  
322 patterns of RNA transcripts versus proteins, introduces new computational challenges for  
323 segmentation and classification. To address the challenge of high dimensionality, a spatially  
324 aware dimensionality reduction step<sup>30</sup> can be integrated into the CelloType framework. To  
325 capture the spatial distribution patterns of RNA transcripts, an additional learnable positional  
326 embedding step can be introduced in the DINO module. These enhancements could significantly  
327 broaden CelloType's applicability to a wide range of spatial omics data.

328  
329

## 330 **Online Methods**

331

### 332 **CelloType**

333

334 A schematic overview of CelloType is depicted in Fig. 1a. The method consists of three  
335 modules: 1) a feature extraction module based on the Transformer deep neural network model to  
336 generate multi-scale image features which are used in the DINO and MaskDINO modules; 2) a  
337 DINO module for object detection and classification; and 3) a MaskDINO module for



338 segmentation. The resulting latent features and detected bounding boxes are then integrated with  
339 the input image in the MaskDINO module to produce instance segmentation results. Both DINO  
340 and MaskDINO modules are integrated in a single neural network model for an end-to-end  
341 learning.

342

### 343 *Feature extraction module*

344

345 Multi-scale image features are generated using the Swin Transformer<sup>17</sup> deep neural network  
346 model. Swin Transformer is a hierarchical version of the original Transformer model that utilizes  
347 shifted window operations for efficient self-attention. It can capture both local and global  
348 features, outperforming conventional convolutional networks in modeling complex image data  
349 with improved computational efficiency. Here we use the Swin-L Transformer model pretrained  
350 on the Common Objects in Context (COCO) Instance Segmentation dataset<sup>25</sup>.

351

### 352 *DINO object detection and classification module*

353

354 The DINO<sup>14</sup> deep neural network architecture, standing for "DETR with Improved DeNoising  
355 Anchor Boxes", is a novel end-to-end object detection model improving upon the DETR  
356 (Detection Transformer) architecture. DINO leverages the strengths of the Transformer  
357 architecture to effectively capture spatial relationships, essential for discerning overlapping or  
358 adjacent cells. On the other hand, DINO incorporates denoising techniques essential for the  
359 precise identification of cells against intricate backgrounds and under-varied imaging conditions.  
360 Major components of the DINO architecture in CelloType are described as follows.

361

362 1. Query Initialization and Selection: To generate the initial anchor box for detecting  
363 objects, the model uses two types of queries: positional queries and content queries. It  
364 initializes anchor boxes only based on the positional information of the selected top-K  
365 features, while keeping content queries unchanged. These queries provide spatial  
366 information of the objects. On the other hand, content queries remain learnable and are  
367 used to extract content features from the image. This mixed query selection strategy helps  
368 the encoder to use better positional information to pool more comprehensive content  
369 features, hence more effectively combines spatial and content information for object  
370 detection. This mixed query selection method is formulated as following:

371

$$Q_{pos} = f_{\text{encoder}}(X), Q_{\text{content}} = \text{learnable}$$

372

373 where  $Q_{pos}$  and  $Q_{\text{content}}$  represent positional and content queries, respectively.  $Q_{pos}$  is a  
374 n-by-4 matrix and  $Q_{\text{content}}$  is a n-by-embed\_dim matrix where n is the number of anchor  
375 boxes and embed\_dim is the embedding dimension.  $X$  represents the flattened image  
376 features and positional embeddings.

376

377 2. Anchor Box Refinement and Contrastive Denoising Training: DINO refines the anchor  
378 boxes step-by-step across decoder layers using deformable attention<sup>31</sup>. The conventional  
379 attention mechanism examines the whole image whereas the deformable attention selects  
380 more important regions of the image and controls the range of self-attention more  
381 flexibly, making the computation more efficient. The conventional denoising training

382 technique<sup>32</sup> involves adding controlled noise to ground truth labels and boxes, formulated  
383 as:

$$384 \quad |\Delta x| < \lambda \frac{w}{2}, |\Delta y| < \lambda \frac{h}{2}, |\Delta w| < \lambda w, |\Delta h| < \lambda h$$

385 where  $(x, y, w, h)$  denotes a ground truth bounding box where  $(x, y)$  is the center  
386 coordinates of the box and  $w$  and  $h$  are the width and height of the box.  $\lambda$  denotes a  
387 hyper-parameter controlling the scale of noise. Contrastive Denoising Training adds both  
388 positive and negative samples of the same ground truth, enhancing the model's ability to  
389 distinguish between objects. DINO involves generating two types of queries (positive and  
390 negative) with different noise scales  $\lambda_1$  and  $\lambda_2$ , where  $\lambda_1 < \lambda_2$ .

391  
392 3. Classification Head and Confidence Score: For the classification of each bounding box, a  
393 linear transformation is applied to the corresponding denoised features. The linear layer  
394 outputs a logit vector  $Z = [z_1, z_2, \dots, z_{K+1}]$ , where  $K$  is the number of classes. The vector  
395 represents the raw predictions for  $K$  classes and the “no object” class. Subsequently, a  
396 SoftMax function is employed on the logit vector to compute the class probabilities:

$$397 \quad \text{SoftMax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^{K+1} e^{z_j}}$$

398 The confidence score for each detected object is taken as the maximum class probability  
399 (excluding the “no object” class) outputted by the model. This score represents the  
400 model's confidence in its prediction of the class for the detected object.

#### 401 ***MaskDINO segmentation module***

402  
403 We use MaskDINO<sup>16</sup> to predict the segmentation masks using outputs from the feature extractor  
404 module and DINO decoder. MaskDINO enhances the DINO architecture by integrating a mask  
405 prediction branch. This mask branch utilizes the DINO decoder's content query embeddings,  $q_c$ ,  
406 to perform dot-product operations with pixel embedding maps, derived from both image and  
407 latent features at high resolution. These operations result in a set of binary masks, where each  
408 segmentation mask,  $m$ , is computed as follows:

$$409 \quad m = q_c \otimes M(T(C_b) + F(C_e))$$

411 where  $q_c$  is the content query embedding,  $M$  is the segmentation head,  $T$  is a convolutional layer  
412 to map the channel dimension to the Transformer hidden dimension,  $C_b$  is the feature map from  
413 the feature extractor module,  $C_e$  is the latent features from the DINO Transformer encoder, and  $F$   
414 is an interpolation-based upsampling function to increase the resolution of latent feature and to  
415 make the result match the size of the image feature.

416 Segmentation task, being a pixel-level classification task, offers more detailed information in the  
417 initial training stages compared to the region-level object detection task. Therefore, MaskDINO  
418 employs the Unified and Enhanced Query Selection technique, which enables the DINO object  
419 detection module to leverage the detailed information from the segmentation task early in the  
420 training process, enhancing the detection task by providing better-initialized queries for

421 subsequent stages. This cooperative task approach between detection and segmentation results in  
422 improved detection performance due to the enhanced box initialization informed by  
423 segmentation mask.

424 During the unified model training, the loss function is calculated by considering three  
425 components: segmentation mask, bounding box prediction and class prediction. The composite  
426 loss function is expressed as follows:

$$427 \quad \text{Loss} = \lambda_{\text{cls}} L_{\text{cls}} + \lambda_{\text{box}} L_{\text{box}} + \lambda_{\text{mask}} L_{\text{mask}}$$

428 where  $L_{\text{cls}}$ ,  $L_{\text{box}}$ ,  $L_{\text{mask}}$  represent classification, bounding box, and segmentation mask losses,  
429 respectively, and  $\lambda_{\text{cls}}$ ,  $\lambda_{\text{box}}$ ,  $\lambda_{\text{mask}}$  are their corresponding weights.

### 430 **Implementation of CelloType for Segmentation tasks**

431 The CelloType software was implemented using the Detectron2 library. Detectron2 is a  
432 Facebook AI Research open source library that provides a high-performance, easy-to-use  
433 implementation of state-of-the-art object detection algorithms written with PyTorch<sup>33</sup>.  
434 Furthermore, it efficiently manages large datasets and features a flexible architecture that  
435 facilitates customization and integration of various image detection or segmentation pipelines.  
436

437 Dataset was randomly divided into 80% for training, 10% for validation, and 10% for testing. All  
438 images and cell/nuclear masks in the training, validation, and testing sets were converted to align  
439 with Detectron2's JSON dictionary schema. For the training dataset, bounding boxes were  
440 derived for each cell using the ground truth segmentation masks. The final dictionary  
441 encompasses the bounding box, segmentation mask, and raw image for each cell.  
442

443 For model training, we initialized the DINO and MaskDINO parameters using the weights  
444 pretrained on the COCO instance segmentation dataset, as this dataset is extensive and diverse,  
445 providing a foundational knowledge for the model. This pretraining helps in better feature  
446 extraction and generalization. We used the Adam optimizer with a learning rate of  $10^{-6}$  and a  
447 batch size of 8. For every 5 training epochs, the trained model was evaluated on the validation  
448 set. The training was terminated when the evaluated AP scores did not improve after 15 epochs.  
449 The model with the best AP scores was used for predicting the cell masks.  
450

451 For evaluation and testing, we set the number of queries to 1,000 which determines the number  
452 of boxes and masks generated by the model. In general, this number should exceed the instance  
453 count in each image yet remain reasonable to reduce computational cost. Considering the  
454 maximum cell count in an image patch in all our datasets does not exceed 1,000, this number  
455 was used as the default parameter. Consequently, the model outputs 1,000 instances per image,  
456 each comprising a segmentation mask and a corresponding confidence score. For testing, a  
457 confidence threshold of 0.3 was used to call predicted instances.  
458

### 460 **Implementation of CelloType for classification task**

461

462 The same training, validation, and testing protocols were used as for the segmentation task.  
463 However, during model training for multiplexed images with over three channels, the `n_channels`  
464 hyperparameter within the Swin Transformer was set to match the input images' dimensionality.

#### 466 **Running of existing methods**

##### 468 *Mesmer*

469  
470 Mesmer was run using the pretrained model detailed by the authors in the “Mesmer-  
471 Application.ipynb” notebook located in the DeepCell-tf GitHub repository. Key parameter  
472 settings included “`image_mpp`”= 0.5, “`compartment`” = “whole-cell” for cell segmentation and  
473 “nuclear” for nuclear segmentation.

##### 475 *Cellpose2*

476  
477 For TissueNet and Cellpose Cyto datasets, Cellpose2 was run using the pretrained model  
478 provided by the authors. For colorectal and bone marrow CODEX datasets, we retrained the  
479 Cellpose2 model following the procedure described by the authors at  
480 <https://cellpose.readthedocs.io/en/latest/gui.html#training-your-own-cellpose-model>.

##### 482 *CellSighter*

483  
484 We trained the CellSighter cell type classification model following the protocol provided by the  
485 authors. Key parameters settings included “`crop_input_size`”=60, “`crop_size`”=128,  
486 “`epoch_max`”=300 epochs, and “`lr`”=0.001.

#### 488 ***Combining Cellpose2 and CellSighter for segmentation and classification***

489  
490 Since there is no existing method for end-to-end joint segmentation and cell type classification,  
491 we devised a baseline model combining Cellpose2 and CellSighter, given their reported high  
492 performance in the respective tasks. Training of the hybrid model comprised two steps, each  
493 optimizing the performance of the individual method. For Cellpose2, CODEX images and  
494 corresponding ground-truth cell segmentation masks were used for model training. For  
495 CellSighter, the same ground-truth cell segmentation masks along with associated cell type  
496 labels were used for training.

497  
498 During the testing phase, a CODEX image was processed with the trained Cellpose2 model to  
499 produce cell segmentation masks, which were subsequently used by the trained CellSighter  
500 model for cell type classification. The final results were the combination of the segmentation  
501 results of Cellpose2 and cell type classification results of CellSighter.

#### 503 **Metrics and procedure for evaluating segmentation accuracy**

504  
505 The Average Precision (AP) metric is a widely adopted standard for evaluating the performance  
506 of instance segmentation methods in computer vision tasks<sup>34,35</sup>. Specifically, for a given  
507 Intersection-over-Union (IoU) threshold,  $t$ , a prediction is considered a true positive if the IoU

508 between the predicted segmentation and the ground truth is greater than  $t$ . The IoU is defined as  
509 the ratio of the area of overlap between the predicted segmentation mask and the ground truth  
510 mask. The AP is calculated at IoU values from 0.50 to 0.9 with a step size of 0.05. The final AP  
511 is the average of the AP values at these different IoU thresholds. This gives a more  
512 comprehensive understanding of a model's performance, from relatively lenient (IoU=0.50) to  
513 stricter overlaps (IoU=0.9).

514

515 In the context of multiple classes, mAP is computed by taking the mean of the AP values  
516 calculated for each individual class. Specifically, if the task only has one class, such as cell  
517 segmentation or nuclear segmentation, the mAP would be the average precision across all the  
518 IoU we evaluated. This gives an overall sense of the method's performance across the various  
519 classes in the dataset, rather than focusing on its efficacy in detecting a single class.

520

521 To evaluate segmentation performance using the AP metric, we used the Common Objects in  
522 Context (COCO) evaluation package, a widely used, standardized benchmarking tool in the field  
523 of instance segmentation. Segmentation results were first converted into the COCO format  
524 before the AP metric was computed using the package. To eliminate redundant detections and  
525 ensure that each object is uniquely identified, the package implements the Non-Maximum  
526 Suppression (NMS) procedure. NMS selectively filters out overlapping bounding boxes,  
527 retaining only the box with the highest confidence score while discarding others with substantial  
528 overlap, as determined by the IoU threshold. Since methods such as Mesmer, Cellpose2, and  
529 CelloType do not generate confidence score for the predicted segmentation masks, we arbitrarily  
530 assigned the confidence score to be 1. For the CelloType variant that outputs the confidence  
531 score (CelloType\_C), we used the actual confidence scores computed by the method when  
532 applying the NMS procedure.

533

534

## 535 **Datasets**

536

### 537 **TissueNet dataset**

538

539 The TissueNet dataset<sup>4</sup> consists of 2,601 training and 1,249 test multiplexed images collected  
540 using multiple imaging platforms and tissue types. Imaging platforms include CODEX, CycIF,  
541 IMC, MIBI, MxIF and Vectra. Tissue types include breast, gastrointestinal, immune cells, lung,  
542 pancreas, and skin. Although many images have dozens of protein markers, all images contain at  
543 least two channels necessary for cell/nucleus segmentation: a cell membrane channel and a  
544 nuclear channel. Each image contains a manual segmentation of cells and/or nuclei. Each  
545 training and test image has a dimension of  $512 \times 512$  pixels and  $256 \times 256$  pixels, respectively.

546

### 547 **Cellpose Cyto dataset**

548

549 The Cyto dataset<sup>6</sup> consists of images from a variety of sources, including: 1) Cells (Cell Image  
550 Library) set: 100 fluorescent images of cultured neurons with both cytoplasmic and nuclear  
551 stains obtained from the Cell Image Library database (<http://www.cellimagelibrary.org>); 2) Cells  
552 (Fluorescent) set: 216 fluorescent images of cells visualized with cytoplasmic markers. This set  
553 contains images from BBBC020, BBBC007v1, mouse cortical and hippocampal cells expressing

554 GCaMP6 imaged using a two-photon microscope, confocal images of mouse cortical neurons,  
555 and the rest were obtained through Google image search; 3) Cells (non-fluorescent) set: 50  
556 brightfield microscopy images from OMERO and Google image search; 4) Cells (Membrane)  
557 set: 58 fluorescent images of cells with membrane maker, 40 of which were from the Micro-Net  
558 image set and the rest were obtained through Google image search; 5) Other microscopy set: 86  
559 images of other types of microscopy that contain either non-cells or cells with atypical  
560 appearances. These images were obtained through Google image search; 6) Non-microscopy set:  
561 98 images of non-microscopy images obtained through Google search of repeating objects  
562 including images of fruits, vegetables, artificial materials, fish and reptile scales, starfish,  
563 jellyfish, sea urchins, rocks, seashells, etc. All images in the dataset were manually segmented by  
564 a human operator.

565

### 566 **Colorectal cancer CODEX dataset**

567

568 This dataset contains CODEX images of 140 human colorectal samples stained with a 56  
569 fluorescent antibodies and 2 nuclear stains<sup>26</sup>. Cells were segmented using Mesmer. Cell types  
570 were annotated by the authors using a combination of iterative clustering and manual  
571 examination of marker expression profiles and morphology. For each tissue image in the dataset,  
572 image patches of 512 x 512 pixels were generated.

573

### 574 **Bone marrow CODEX dataset**

575

576 This dataset<sup>27</sup> contains CODEX images of 12 human bone marrow samples stained with 54  
577 fluorescent antibodies and one nuclear stain. Hematopoietic cell types were annotated by the  
578 authors using a combination of iterative clustering and manual examination of marker expression  
579 profiles and morphology. Adipocytes and trabecular bone fragments were manually annotated by  
580 the authors.

581

### 582 **Code availability**

583

584 CelloType is available at: <https://github.com/tanlabcode/CelloType>.

585

586

### 587 **Figure Legends**

588

#### 589 **Figure 1 – Overview of CelloType.**

590

591 **a)** Overall architecture, input, and output of CelloType. First, a Transformer-based feature  
592 extractor is employed to derive multi-scale features ( $C_b$ ) from the image. Second, using a  
593 Transformer-based architecture, the DINO object detection module extracts latent features ( $C_e$ )  
594 and query embeddings ( $q_c$ ) that are combined to generate object detection boxes with cell type  
595 labels. Subsequently, the MaskDINO module integrates the extracted image features with  
596 DINO's outputs, resulting in detailed instance segmentation and cell type classification. During  
597 training, the model is optimized based on an overall loss function ( $Loss$ ) that considers losses  
598 based on cell segmentation mask ( $\lambda_{mask}L_{mask}$ ), bounding box ( $\lambda_{box}L_{box}$ ), and cell type label  
599 ( $\lambda_{cls}L_{cls}$ ). **b)** Input, output, and architecture of the DINO module. The DINO module consists of

600 a multi-layer Transformer and multiple prediction heads. DINO starts by flattening the multi-  
601 scale features from the Transformer-based feature extractor. These features are merged with  
602 positional embeddings to preserve spatial context (step 1 in the figure). DINO then employs a  
603 mixed query selection strategy, initializing positional queries ( $Q_{pos}$ ) as anchor detection boxes  
604 and maintaining content queries ( $Q_{content}$ ) as learnable features, thus adapting to the diverse  
605 characteristics of cells (step 2). The model refines these anchor boxes through decoder layers  
606 using deformable attention mechanism and employs contrastive denoising training by  
607 introducing noise to ground truth (GT) labels and boxes to improve robustness and accuracy.  
608 Then a linear projection acts as the classification branch to produce the classification results for  
609 each box (step 3). **c) Multi-scale ability of CelloType.** CelloType is versatile and can perform a  
610 range of end-to-end tasks at different scales, including cell segmentation, nuclear segmentation,  
611 microanatomical structure segmentation, and full instance segmentation with corresponding class  
612 annotations.

613

### 614 **Figure 2 – Evaluation of segmentation accuracy using TissueNet datasets**

615

616 **a)** Average Precision (AP) across Intersection over Union (IoU) thresholds for cell segmentation  
617 by Mesmer, Cellpose2, CelloType and CelloType\_C (CelloType with confidence score). Mean  
618 AP value across IoU thresholds of 0.5-0.9 (mAP) for each method is indicated in parenthesis. **b)**  
619 AP across IoU thresholds for nuclear segmentation. **c)** Performance of methods stratified by  
620 imaging platform and tissue type. The top left heatmap shows the mAP scores for cell  
621 segmentation stratified by imaging platform, including CODEX, CyCIF, IMC, MIBI, MxIF and  
622 Vertra. The top right heatmap shows the mAP scores for cell segmentation stratified by tissue  
623 type, including breast, gastrointestinal, immune, pancreas and skin. The second row of heatmaps  
624 shows the mAP values for nuclear segmentation. **d)** Representative examples of cell  
625 segmentation of immune tissue imaged using Vectra platform. Blue, nuclear channel;  
626 green, membrane channel; white, cell boundary. The red box highlights a representative region  
627 that the methods perform differently. The AP75 score (Average precision at IoU threshold of  
628 0.75) is displayed on the images. **e)** Representative examples of nuclear segmentation of  
629 gastrointestinal tissue using the IMC platform. The AP50 scores are shown on the images.

630

### 631 **Figure 3 – Evaluation of segmentation accuracy using Cellpose Cyto dataset**

632

633 **a)** Average precision (AP) across Intersection over Union (IoU) thresholds for Cellpose2,  
634 CelloType and CelloType\_C (CelloType with confidence score). Mean AP value across IoU  
635 thresholds of 0.5-0.9 (mAP) for each method is indicated in parenthesis. **b)** Mean AP values of  
636 Cellpose2, CelloType, and CelloType\_C stratified by imaging modalities and cell types. The test  
637 dataset comprises microscopy and non-microscopy images from the Cellpose Cyto dataset that  
638 comprises 6 subsets, including Cells (Cell Image Library), Cells (Fluorecent), Cells (Non-  
639 fluorecent), Cells (Membrane), Other microscopy, and Non-microscopy. **c)** Representative  
640 examples of cell segmentation of a microscopy image by the compared methods. The red boxes  
641 highlight a representative region that the methods perform differently. The AP75 score is  
642 displayed on the images. **d)** Representative examples of cell segmentation of a non-fluorescent  
643 image by the compared methods.

644

### 645 **Figure 4 – CelloType performs joint segmentation and cell type classification.**

646  
647 **a)** Barplot showing AP50 values for cell type annotation by the two compared methods.  
648 **b)** Line plot showing the relationship between classification accuracy and confidence score  
649 threshold by the two methods. **c)** Representative examples of cell segmentation and classification  
650 results using the colorectal cancer CODEX dataset. Each row represents a 200 by 200 pixels  
651 field of view (FOV) of a CODEX image. Each FOV shows predicted cell segmentation masks  
652 (boxes) and cell types (colors). Ground Truth, manually annotated cell types; CelloType, end-to-  
653 end cell segmentation and cell type classification; Cellpose2+CellSighter, cell segmentation by  
654 Cellpose 2 followed by cell type classification by CellSighter. Randomly selected confidence  
655 scores for cell classification computed by the two methods were displayed next to the predicted  
656 instances.

657  
658 **Figure 5 – Performance benchmarking of Cellpose2 and CellSighter.**

659 Each method was evaluated for its originally intended task, namely Cellpose2 for segmentation  
660 and CellSighter for cell classification. Colorectal cancer CODEX dataset was used for  
661 benchmarking purpose. **a)** AP value of segmentation across a range of IoU thresholds. Mean AP  
662 value (mAP) is shown in parenthesis. **b)** Heatmap showing the confusion matrix of CellSighter  
663 cell type classification results. Ground truth cell segmentation masks were used as input to  
664 CellSighter. Each grid in the heatmap includes an accuracy score and the count of cells. **c)**  
665 Barplot showing the precision scores for each class identified by the CellSighter model based on  
666 the ground truth cell segmentation mask, with an overall mean precision of 0.53.

667  
668 **Figure 6 – CelloType supports joint multi-scale segmentation and classification.**

669 **a)** Performance evaluation of CelloType stratified by cell and microanatomic structure types.  
670 The bar plot shows the mean and 95% confidence interval of AP50 values in 5-fold cross-  
671 validation experiments. **b)** Line plot showing the relationship between classification accuracy  
672 and confidence score threshold. **c)** Representative examples of multi-scale segmentation and  
673 classification using human bone marrow CODEX data. The first row of images shows an  
674 example of bone marrow area consisting of various types of smaller hematopoietic cells and  
675 much larger adipocytes. The second row of images shows an example of bone marrow area  
676 consisting of various hematopoietic cell types and microanatomic structure such as trabecula  
677 bone fragments. Randomly selected confidence scores for cell classification were displayed next  
678 to the predicted instances.

679  
680  
681 **Supplemental Figure 1 – Distribution of cell types in the colorectal cancer CODEX dataset.**

682  
683 **Supplemental Figure 2 – Distribution of cell types in the human bone marrow CODEX**  
684 **dataset.**

685  
686 **Acknowledgments**  
687

688 The authors thank the Children’s Hospital of Philadelphia Research Information Services for  
689 providing computing support. This work was supported by the National Cancer Institute (NCI)  
690 Human Tumor Atlas Network grant under award #U2C CA233285 (K.T.) and the National  
691 Institutes of Health (NIH) Human Biomolecular Atlas Program grant under award #U54



692 HL165442 (K.T.).

693

694

### 695 **Author Contributions**

696

697 M.P. and K.T. conceived and designed the study. M.P. implemented the CelloType algorithm  
698 and wrote the software, M.P. and T.R. performed data analysis, X.W. and K.T. supervised the  
699 overall study. M.P. and K.T. wrote the manuscript with input from all authors.

700

### 701 **Competing interests**

702 The authors declare no competing interests.

703

704

### 705 **References**

706

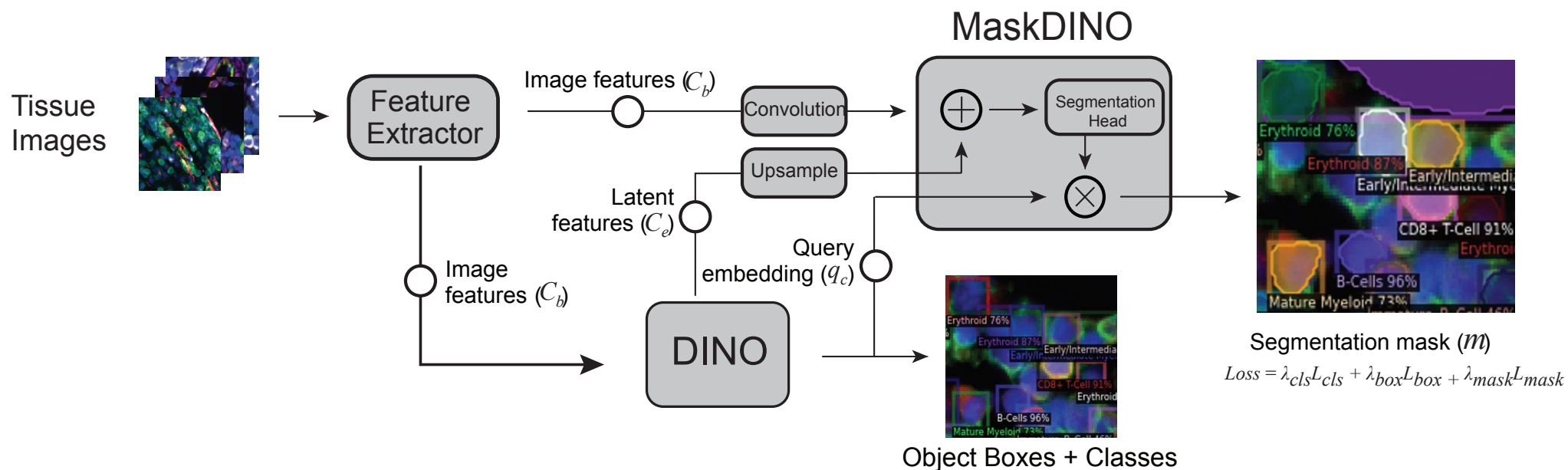
- 707 1. Bressan, D., Battistoni, G., and Hannon, G.J. (2023). The dawn of spatial omics. *Science*  
708 *381*, eabq4964. [10.1126/science.abq4964](https://doi.org/10.1126/science.abq4964).
- 709 2. Rozenblatt-Rosen, O., Regev, A., Oberdoerffer, P., Nawy, T., Hupalowska, A., Rood,  
710 J.E., Ashenberg, O., Cerami, E., Coffey, R.J., Demir, E., et al. (2020). The Human Tumor  
711 Atlas Network: Charting Tumor Transitions across Space and Time at Single-Cell  
712 Resolution. *Cell* *181*, 236-249. [10.1016/j.cell.2020.03.053](https://doi.org/10.1016/j.cell.2020.03.053).
- 713 3. Hu, B.C., Writing, G., Snyder, M.P., Lin, S., Posgai, A., Atkinson, M., Regev, A., Rood,  
714 J., Rozenblatt-Rosen, O., Gaffney, L., et al. (2019). The human body at cellular  
715 resolution: the NIH Human Biomolecular Atlas Program. *Nature* *574*, 187-192.  
716 [10.1038/s41586-019-1629-x](https://doi.org/10.1038/s41586-019-1629-x).
- 717 4. Greenwald, N.F., Miller, G., Moen, E., Kong, A., Kagel, A., Dougherty, T., Fullaway,  
718 C.C., McIntosh, B.J., Leow, K.X., Schwartz, M.S., et al. (2022). Whole-cell segmentation  
719 of tissue images with human-level performance using large-scale data annotation and  
720 deep learning. *Nat Biotechnol* *40*, 555-565. [10.1038/s41587-021-01094-0](https://doi.org/10.1038/s41587-021-01094-0).
- 721 5. He, K., Zhang, X., Ren, S., Sun, Ji. (2016). Deep Residual Learning for Image  
722 Recognition. 2016 IEEE Conference on Computer Vision and Pattern Recognition  
723 (CVPR), (IEEE), pp. 770-778.
- 724 6. Stringer, C., Wang, T., Michaelos, M., and Pachitariu, M. (2021). Cellpose: a generalist  
725 algorithm for cellular segmentation. *Nature Methods* *18*, 100-106. [10.1038/s41592-020-](https://doi.org/10.1038/s41592-020-01018-x)  
726 [01018-x](https://doi.org/10.1038/s41592-020-01018-x).
- 727 7. Pachitariu, M., and Stringer, C. (2022). Cellpose 2.0: how to train your own model.  
728 *Nature Methods* *19*, 1634-1641. [10.1038/s41592-022-01663-4](https://doi.org/10.1038/s41592-022-01663-4).
- 729 8. Ronneberger, O., Fischer, P., and Brox, T. (2015). U-Net: Convolutional Networks for  
730 Biomedical Image Segmentation. In *Medical Image Computing and Computer-Assisted*  
731 *Intervention – MICCAI 2015*, N. Navab, J. Hornegger, W.M. Wells, and A.F. Frangi,  
732 eds. (Springer International Publishing), pp. 234-241.
- 733 9. Amitay, Y., Bussi, Y., Feinstein, B., Bagon, S., Milo, I., and Keren, L. (2023).  
734 CellSighter: a neural network to classify cells in highly multiplexed images. *Nature*  
735 *Communications* *14*, 4302. [10.1038/s41467-023-40066-7](https://doi.org/10.1038/s41467-023-40066-7).

- 736 10. Zhang, W., Li, I., Reticker-Flynn, N.E., Good, Z., Chang, S., Samusik, N., Saumyaa, S.,  
737 Li, Y., Zhou, X., Liang, R., et al. (2022). Identification of cell types in multiplexed in situ  
738 images by combining protein expression and spatial information using CELESTA.  
739 *Nature Methods* *19*, 759-769. 10.1038/s41592-022-01498-z.
- 740 11. Isensee, F., Jaeger, P.F., Kohl, S.A.A., Petersen, J., and Maier-Hein, K.H. (2021). nnU-  
741 Net: a self-configuring method for deep learning-based biomedical image segmentation.  
742 *Nature Methods* *18*, 203-211. 10.1038/s41592-020-01008-z.
- 743 12. Vaswani, A., Shazeer, N.M., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.,  
744 Polosukhin, I. (2017). Attention is all you need. *Neural Information Processing Systems*,  
745 6000-6010.
- 746 13. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., and Zagoruyko, S. (2020).  
747 End-to-End Object Detection with Transformers. In *Computer Vision – ECCV 2020*, A.  
748 Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, eds. (Springer International Publishing),  
749 pp. 213-229.
- 750 14. Zhang, H., Li, F., Liu, S., Zhang, L., Su, H., Zhu, J., Ni, L.M., and Shum, H.-Y. (2022).  
751 DINO: DETR with Improved DeNoising Anchor Boxes for End-to-End Object  
752 Detection. *arXiv*.
- 753 15. Hörst, F., Rempe, M., Heine, L., Seibold, C., Keyl, J., Baldini, G., Ugurel, S., Siveke, J.,  
754 Grünwald, B., Egger, J., and Kleesiek, J. (2023). CellViT: Vision Transformers for  
755 Precise Cell Segmentation and Classification. *arXiv*.
- 756 16. Li, F., Zhang, H., Xu, H., Liu, S., Zhang, L., Ni, L.M., and Shum, H.-Y. (2023). Mask  
757 DINO: Towards A Unified Transformer-based Framework for Object Detection and  
758 Segmentation. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
759 (CVPR), (IEEE), pp. 3041-3050.
- 760 17. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. (2021). Swin  
761 Transformer: Hierarchical Vision Transformer using Shifted Windows. *2021 IEEE/CVF*  
762 *International Conference on Computer Vision (ICCV)*, (IEEE), pp. 9992-10002.
- 763 18. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T.,  
764 Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. (2021). An Image is Worth  
765 16x16 Words: Transformers for Image Recognition at Scale. *arXiv*.
- 766 19. Goltsev, Y., Samusik, N., Kennedy-Darling, J., Bhate, S., Hale, M., Vazquez, G., Black,  
767 S., and Nolan, G.P. (2018). Deep Profiling of Mouse Splenic Architecture with CODEX  
768 Multiplexed Imaging. *Cell* *174*, 968-981.e915. 10.1016/j.cell.2018.07.010.
- 769 20. Lin, J.-R., Fallahi-Sichani, M., and Sorger, P.K. (2015). Highly multiplexed imaging of  
770 single cells using a high-throughput cyclic immunofluorescence method. *Nature*  
771 *Communications* *6*, 8390. 10.1038/ncomms9390.
- 772 21. Ali, H.R., Jackson, H.W., Zanutelli, V.R.T., Danenberg, E., Fischer, J.R., Bardwell, H.,  
773 Provenzano, E., Team, C.I.G.C., Ali, H.R., Al Sa'd, M., et al. (2020). Imaging mass  
774 cytometry and multiplatform genomics define the phenogenomic landscape of breast  
775 cancer. *Nature Cancer* *1*, 163-175. 10.1038/s43018-020-0026-6.
- 776 22. Keren, L., Bosse, M., Thompson, S., Risom, T., Vijayaragavan, K., McCaffrey, E.,  
777 Marquez, D., Angoshtari, R., Greenwald, N.F., Fienberg, H., et al. (2019). MIBI-TOF: A  
778 multiplexed imaging platform relates cellular phenotypes and tissue structure. *Science*  
779 *Advances* *5*, eaax5851. 10.1126/sciadv.aax5851.
- 780 23. Gerdes, M.J., Sevinsky, C.J., Sood, A., Adak, S., Bello, M.O., Bordwell, A., Can, A.,  
781 Corwin, A., Dinn, S., Filkins, R.J., et al. (2013). Highly multiplexed single-cell analysis

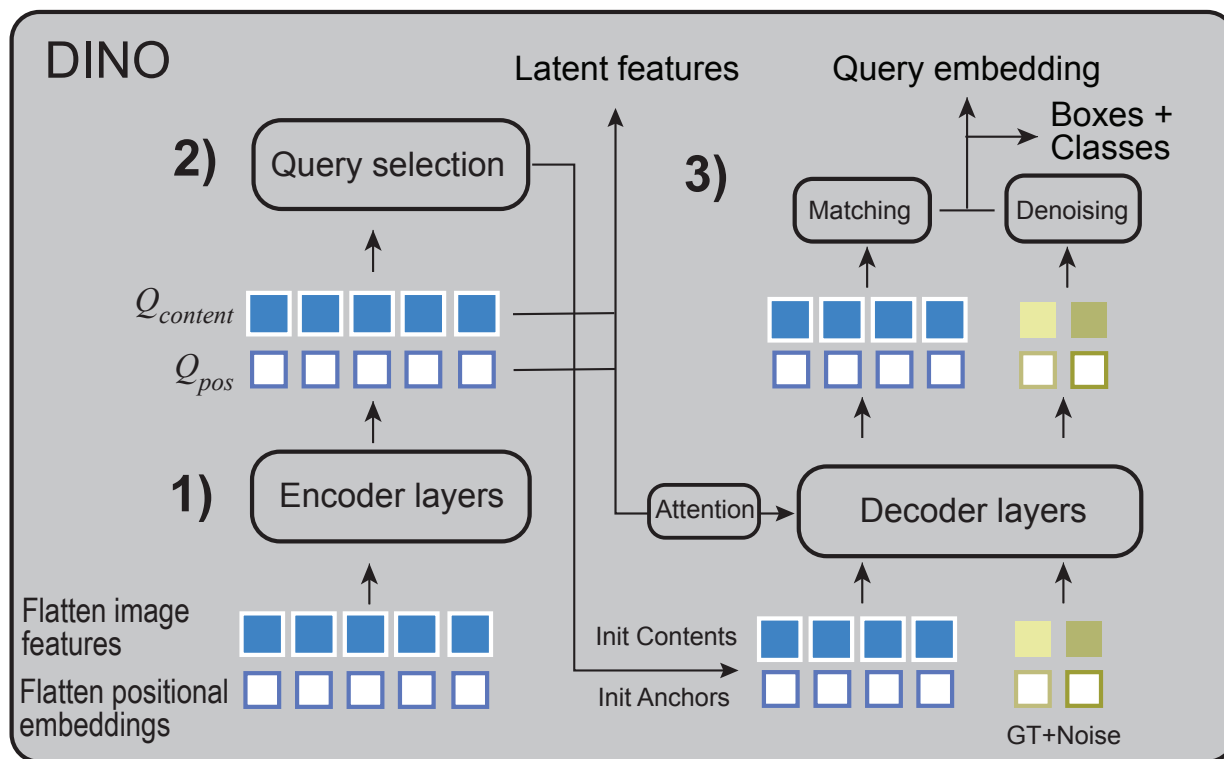
- 782 of formalin-fixed, paraffin-embedded cancer tissue. *Proceedings of the National*  
783 *Academy of Sciences* *110*, 11982-11987. 10.1073/pnas.1300136110.
- 784 24. Fiore, C., Bailey, D., Conlon, N., Wu, X., Martin, N., Fiorentino, M., Finn, S., Fall, K.,  
785 Andersson, S.-O., Andren, O., et al. (2012). Utility of multispectral imaging in automated  
786 quantitative scoring of immunohistochemistry. *J Clin Pathol* *65*, 496-502.  
787 10.1136/jclinpath-2012-200734.
- 788 25. Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and  
789 Zitnick, C.L. (2014). Microsoft COCO: Common Objects in Context. In *Computer*  
790 *Vision – ECCV 2014*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, eds. (Springer  
791 International Publishing), pp. 740-755.
- 792 26. Schürch, C.M., Bhate, S.S., Barlow, G.L., Phillips, D.J., Noti, L., Zlobec, I., Chu, P.,  
793 Black, S., Demeter, J., McIlwain, D.R., et al. (2020). Coordinated Cellular  
794 Neighborhoods Orchestrate Antitumoral Immunity at the Colorectal Cancer Invasive  
795 Front. *Cell* *182*, 1341-1359.e1319. 10.1016/j.cell.2020.07.005.
- 796 27. Bandyopadhyay, S., Duffy, M., Ahn, K., Pang, M., Smith, D., Duncan, G., Sussman, J,  
797 Zhang, I, Huang, J, Lin, Y, Xiong, B, Imtiaz, T, Chen, C, Thadi, A, Chen, C, Xu, J,  
798 Reichart, M, Pillai, V, Snaith, O, Oldridge, D, Bhattacharyya, S, Maillard, I, Carroll, M,  
799 Nelson, C, Qin, L, Tan, K (2024). Mapping the cellular biogeography of human bone  
800 marrow niches using single-cell transcriptomics and proteomic imaging. bioRxiv.  
801 10.1101/2024.03.14.585083.
- 802 28. Feng, Y., Wang, Y., Li, H., Qu, M., and Yang, J. (2023). Learning what and where to  
803 segment: A new perspective on medical image few-shot segmentation. *Medical Image*  
804 *Analysis* *87*, 102834. 10.1016/j.media.2023.102834.
- 805 29. Xie, E., Ding, J., Wang, W., Zhan, X., Xu, H., Sun, P., Li, Z., and Luo, P. (2021). DetCo:  
806 Unsupervised Contrastive Learning for Object Detection. 2021 IEEE/CVF International  
807 Conference on Computer Vision (ICCV), (IEEE), pp. 8372-8381.
- 808 30. Shang, L., and Zhou, X. (2022). Spatially aware dimension reduction for spatial  
809 transcriptomics. *Nature Communications* *13*, 7203. 10.1038/s41467-022-34879-1.
- 810 31. Zhu, X., Su, W, Lu, L, Li, B, Wang, X, Dai, J (2021). Deformable DETR: Deformable  
811 Transformer for End-to-End Object Detection. 2021 International Conference on  
812 Learning Representations (ICLR), (IEEE).
- 813 32. Li, F., Zhang, H., Liu, S., Guo, J., Ni, L.M., and Zhang, L. (2022). DN-DETR: Accelerate  
814 DETR Training by Introducing Query DeNoising. 2022 Computer Vision and Pattern  
815 Recognition Conference (CVPR), (IEEE).
- 816 33. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z.,  
817 Gimelshein, N., Antiga, L., et al. PyTorch: An Imperative Style, High-Performance Deep  
818 Learning Library.
- 819 34. Padilla, R., Netto, S.L., and Da Silva, E.A.B. (2020). A Survey on Performance Metrics  
820 for Object-Detection Algorithms. 2020 International Conference on Systems, Signals and  
821 Image Processing (IWSSIP), (IEEE), pp. 237-242.
- 822 35. Hirling, D., Tasnadi, E., Caicedo, J., Caroprese, M.V., Sjögren, R., Aubreville, M., Koos,  
823 K., and Horvath, P. (2024). Segmentation metric misinterpretations in bioimage analysis.  
824 *Nature Methods* *21*, 213-216. 10.1038/s41592-023-01942-8.
- 825

# Figure 1

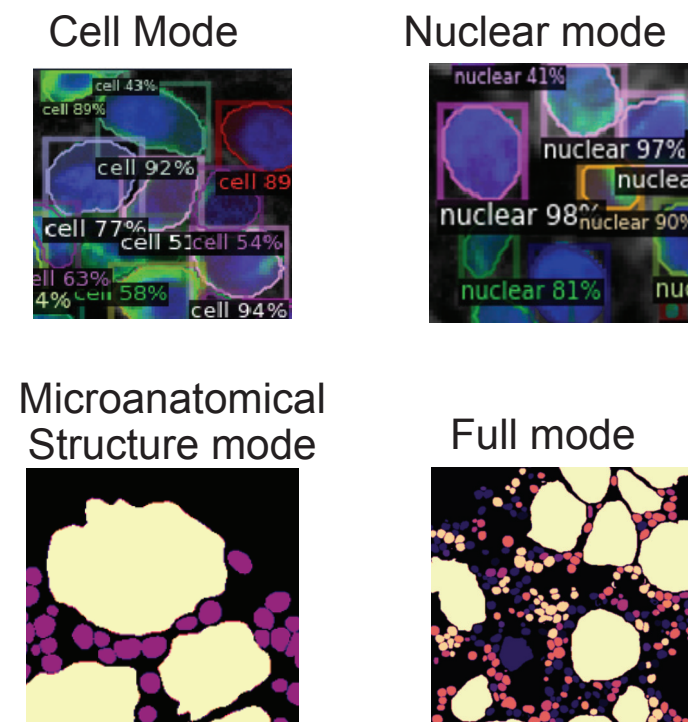
**a**

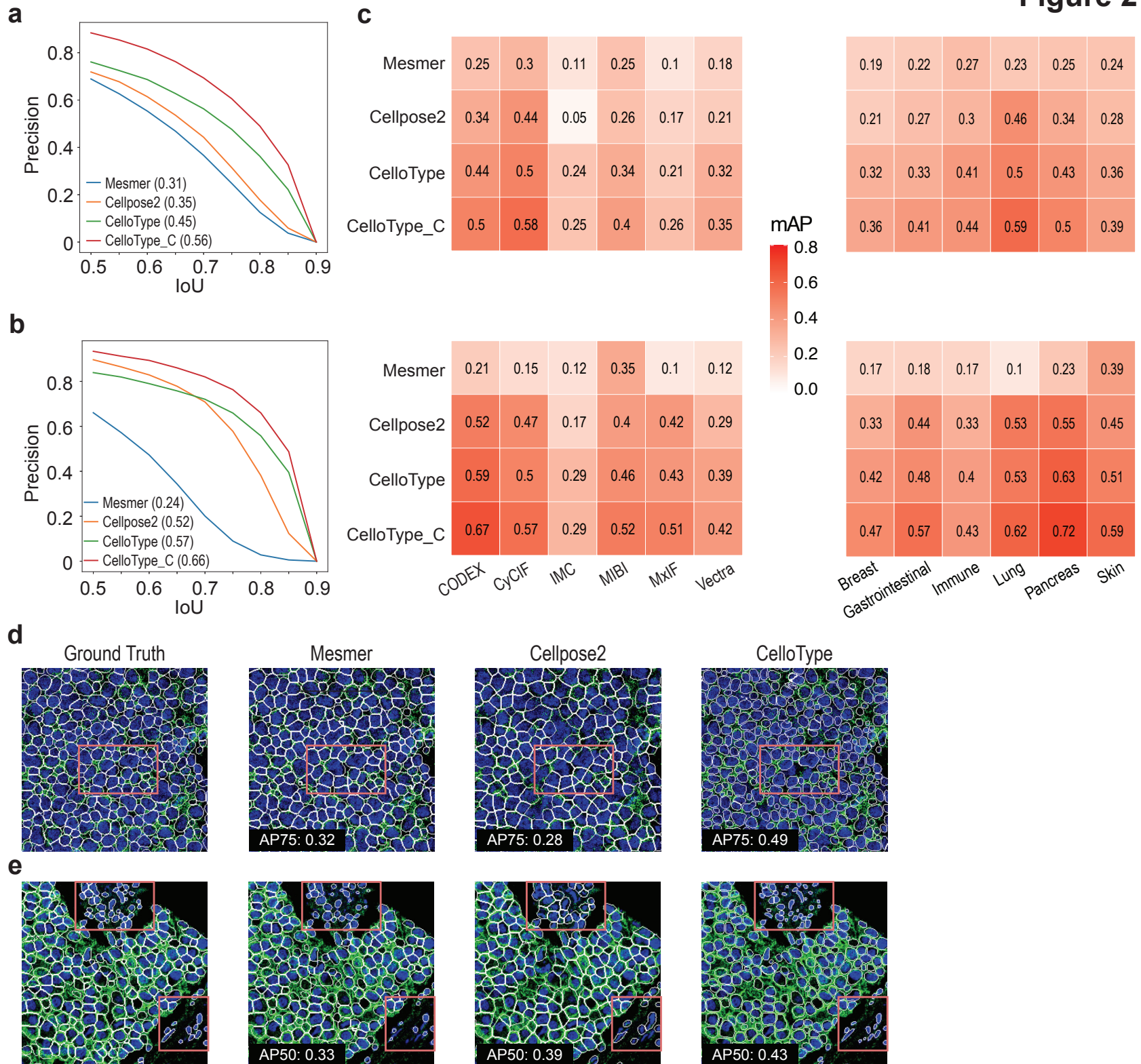


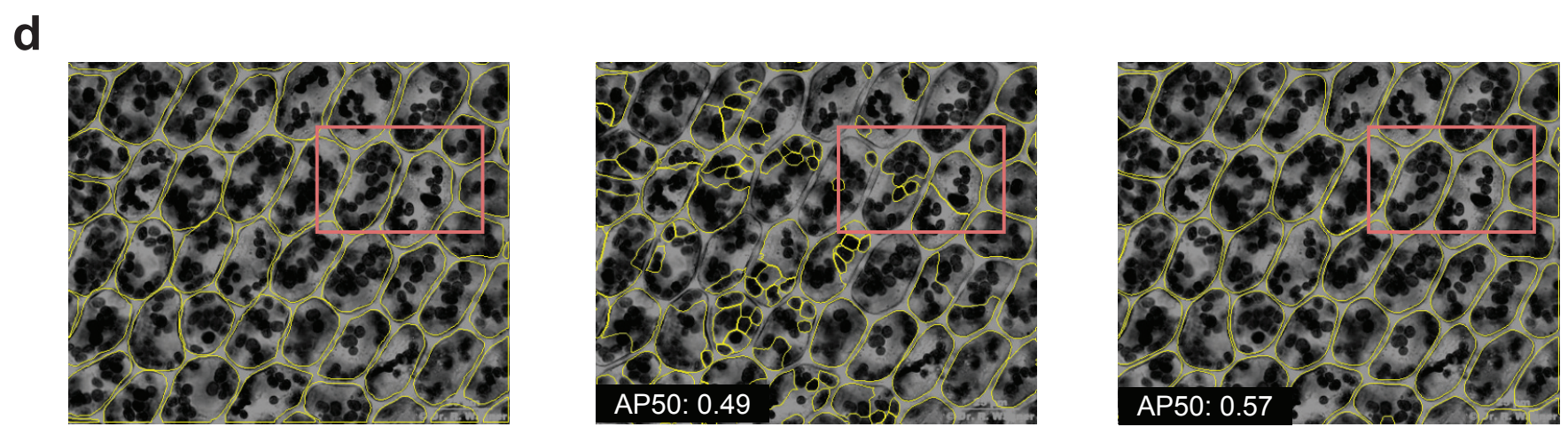
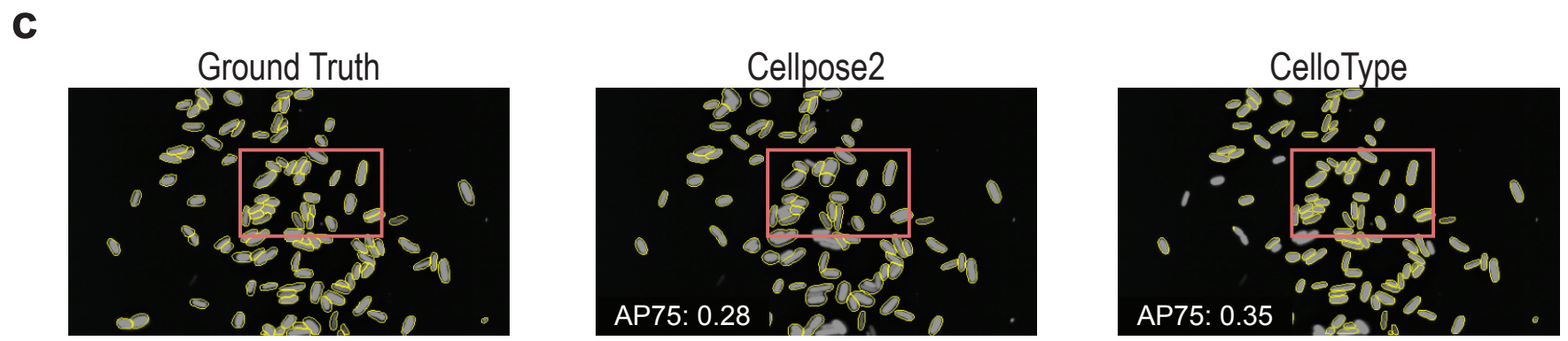
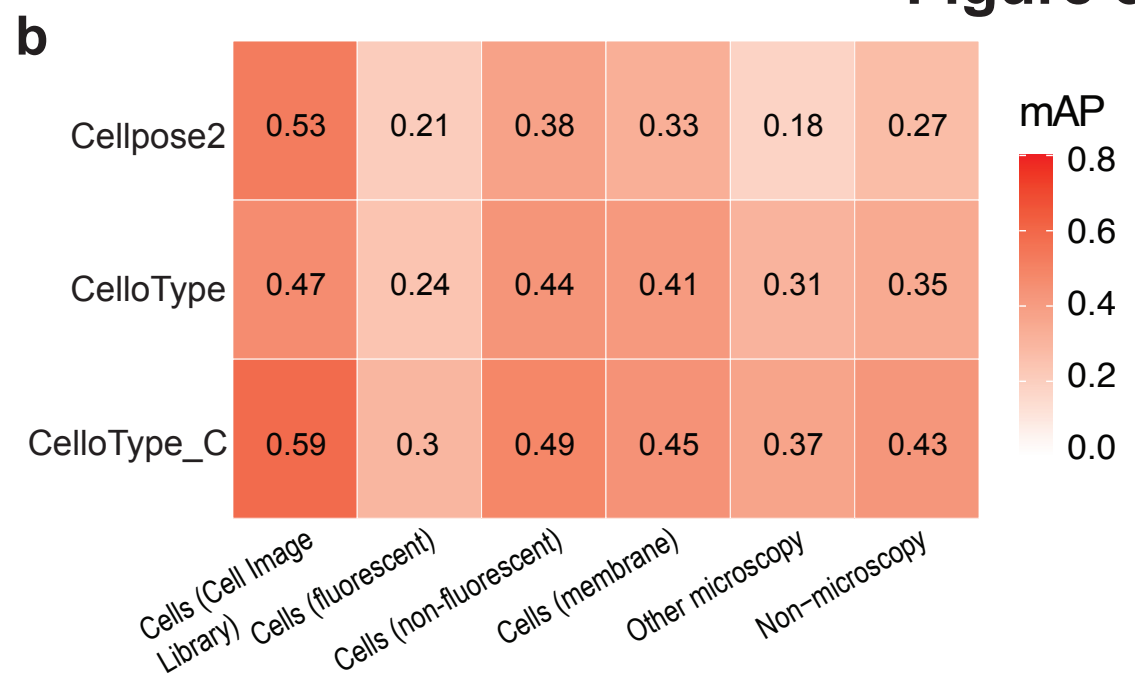
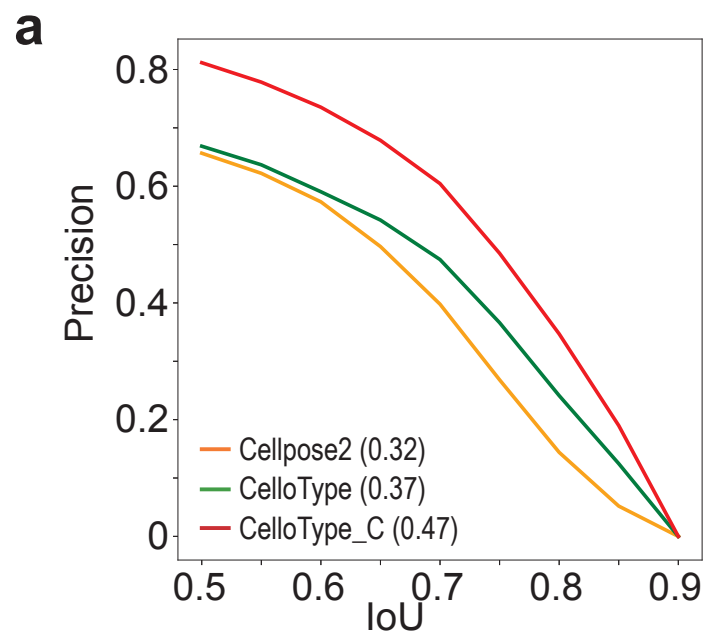
**b**



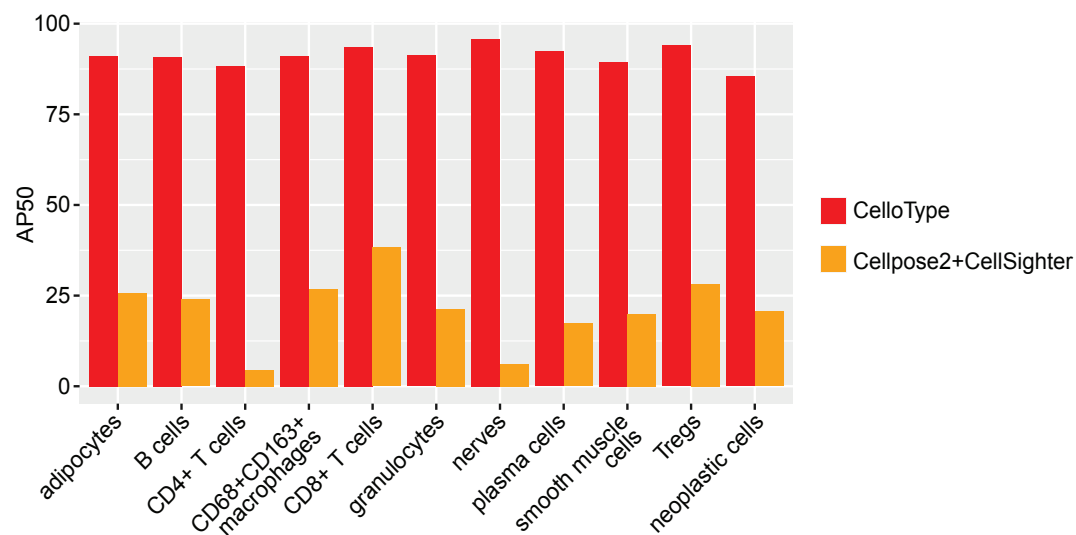
**c**



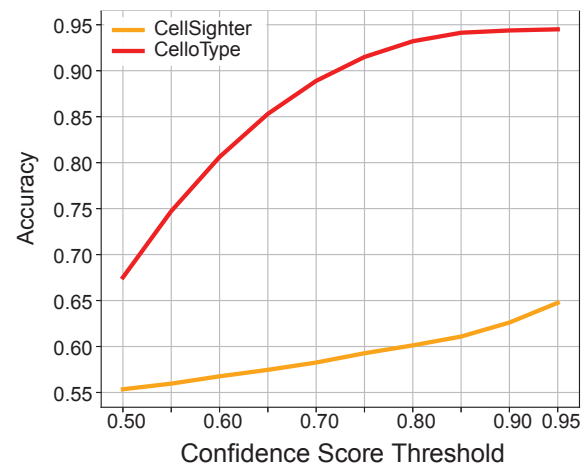
**Figure 2**

**Figure 3**

**a**



**b**

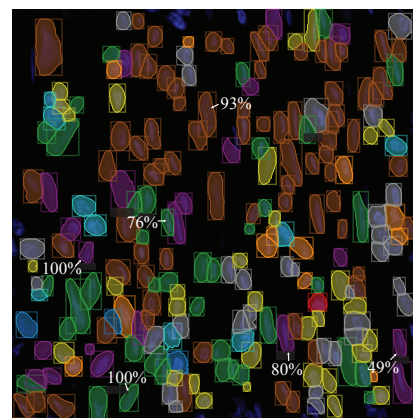
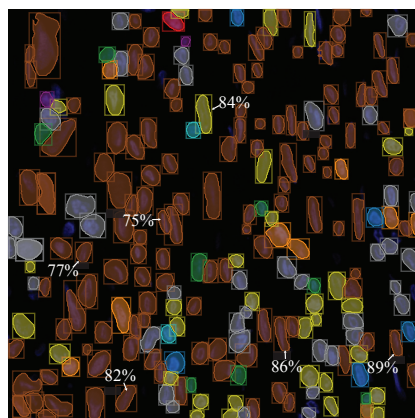
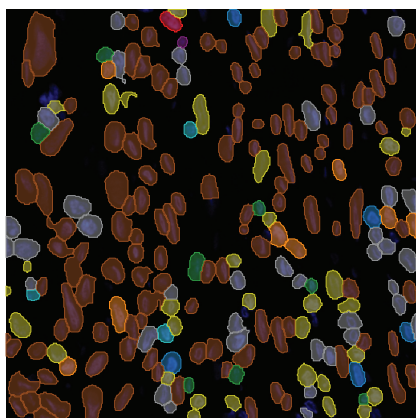
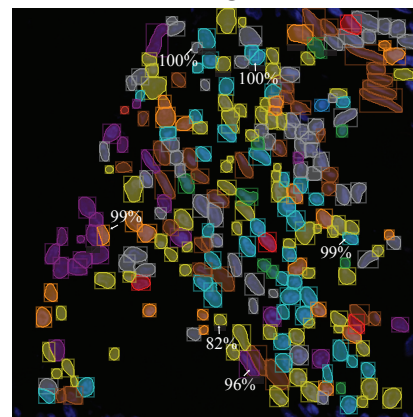
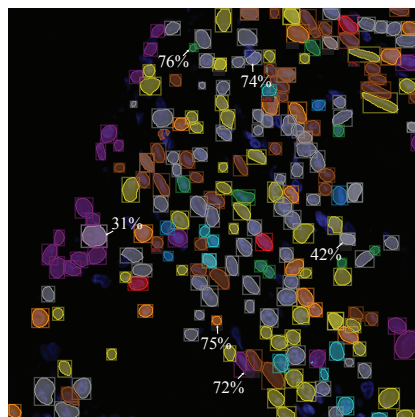
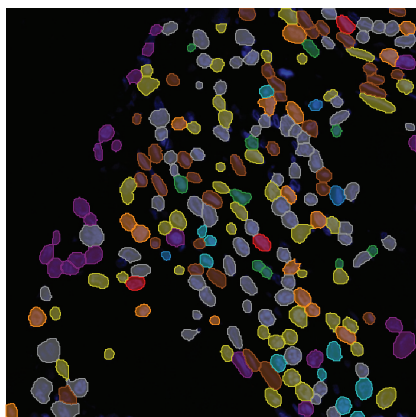


**c**

Ground Truth

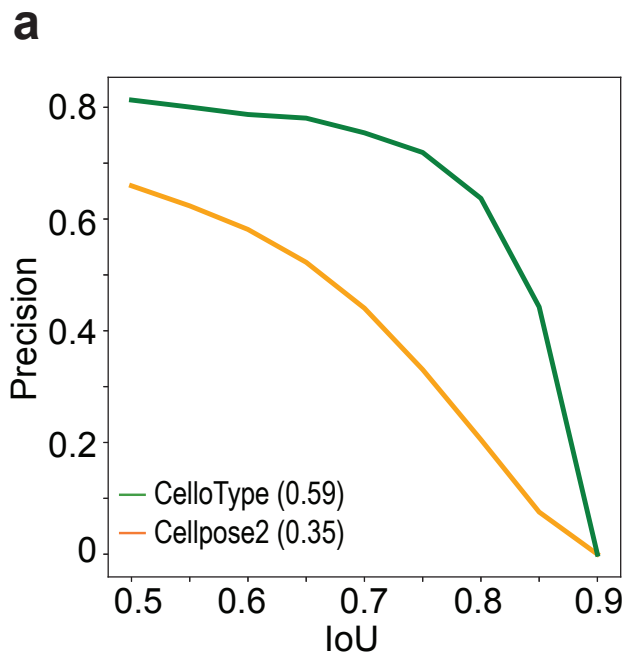
CelloType

Cellpose2 +  
CellSighter



- granulocytes
- CD68+ CD163+ macrophages
- adipocytes
- plasma cells
- CD8+ T cells
- Tregs
- CD4+ T cells
- B cells
- smooth muscle cells
- nerves
- neoplastic cells
- others

# Figure 5

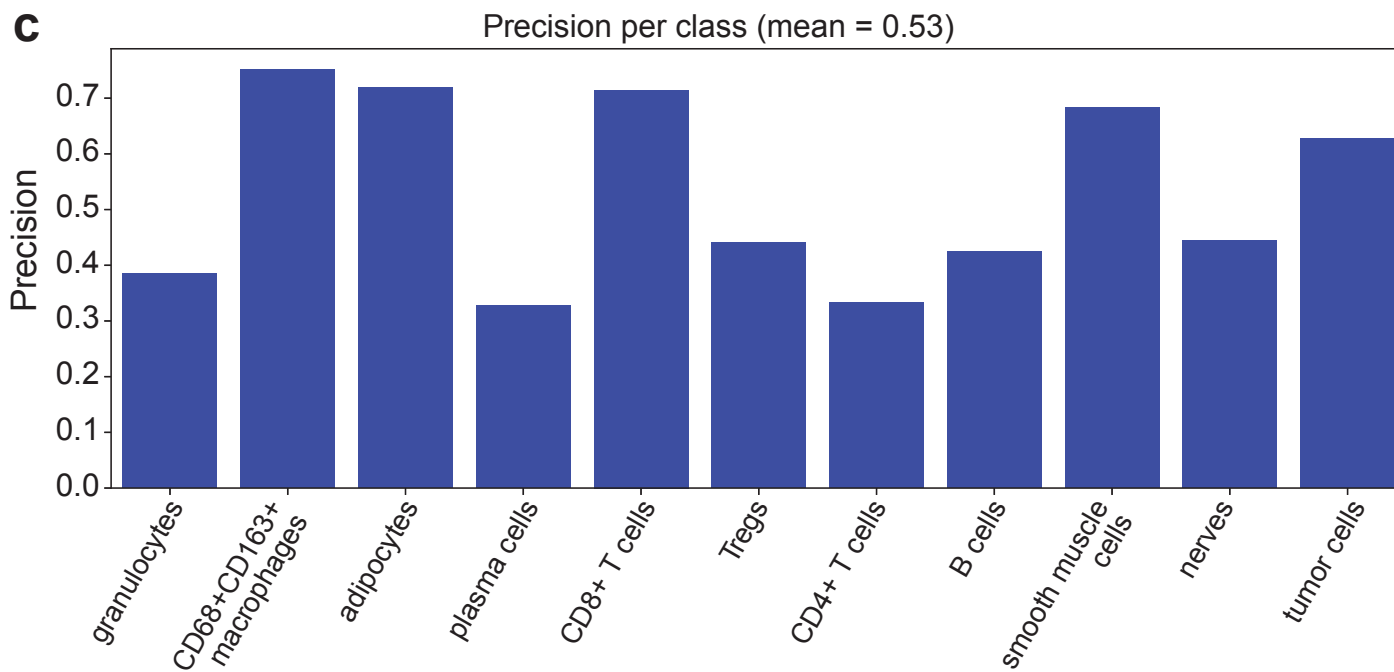


**b**

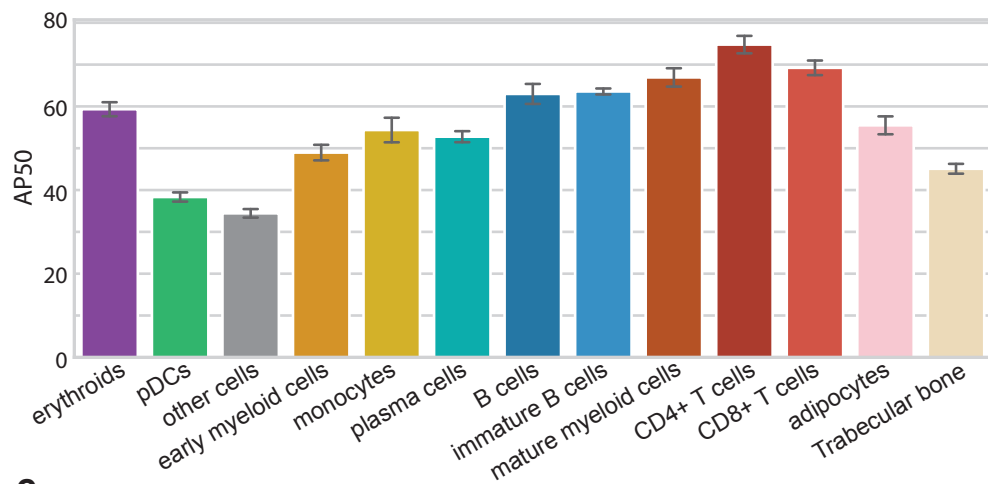
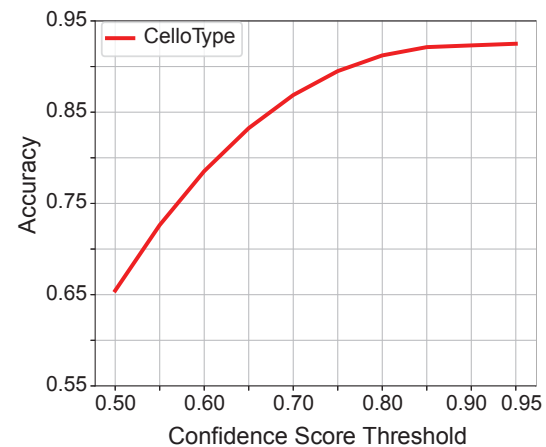
	<i>granulocytes</i>	<i>tumor cells</i>	<i>CD68+CD163+ macrophages</i>	<i>adipocytes</i>	<i>plasma cells</i>	<i>CD8+ T cells</i>	<i>Tregs</i>	<i>CD4+ T cells</i>	<i>B cells</i>	<i>smooth muscle cells</i>	<i>nerves</i>
<i>granulocytes</i>	55.0 (266)	22.5 (109)	8.7 (42)	0.4 (2)	1.9 (9)	1.9 (9)	0.4 (2)	1.2 (6)	0.6 (3)	7.4 (36)	
<i>tumor cells</i>	1.8 (17)	84.1 (791)	3.9 (37)	0.6 (6)	1.4 (13)	1.5 (14)			0.3 (3)	6.2 (58)	
<i>CD68+CD163+ macrophages</i>	2.0 (31)	2.2 (33)	87.4 (1334)	0.4 (6)	0.9 (13)	2.3 (35)		0.3 (4)	0.5 (7)	3.9 (60)	
<i>adipocytes</i>		11.6 (11)	4.2 (4)	77.9 (74)		2.1 (2)	2.1 (2)		2.1 (2)		
<i>plasma cells</i>	3.3 (14)	14.8 (62)	15.0 (63)		61.0 (256)	2.4 (10)			1.0 (4)	2.6 (11)	
<i>CD8+ T cells</i>	0.9 (6)	2.2 (14)	7.6 (49)	0.2 (1)	0.9 (6)	84.3 (541)	0.2 (1)	0.9 (6)	1.1 (7)	1.7 (11)	
<i>Tregs</i>	0.9 (1)	1.8 (2)	12.8 (14)			5.5 (6)	61.5 (67)	12.8 (14)	1.8 (2)	2.8 (3)	
<i>CD4+ T cells</i>		5.6 (1)	16.7 (3)				5.6 (1)	61.1 (11)	11.1 (2)		
<i>B cells</i>	0.7 (2)	8.4 (25)	12.8 (38)	1.7 (5)	1.0 (3)	4.7 (14)	0.3 (1)	2.7 (8)	66.4 (198)	1.0 (3)	0.3 (1)
<i>smooth muscle cells</i>	1.7 (13)	3.2 (24)	5.1 (38)		1.7 (13)	1.3 (10)			1.2 (9)	85.2 (640)	0.4 (3)
<i>nerves</i>		22.2 (2)			11.1 (1)				11.1 (1)	11.1 (1)	44.4 (4)

Prediction

Ground Truth





**Figure 6****a****b****c**