

Evaluation of cell segmentation methods without reference segmentations

Haoran Chen and Robert F. Murphy*

Computational Biology Department, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213

ABSTRACT Cell segmentation is a cornerstone of many bioimage informatics studies, and inaccurate segmentation introduces error in downstream analysis. Evaluating segmentation results is thus a necessary step for developing segmentation methods as well as for choosing the most appropriate method for a particular type of sample. The evaluation process has typically involved comparison of segmentations with those generated by humans, which can be expensive and subject to unknown bias. We present here an approach to evaluating cell segmentation methods without relying upon comparison to results from humans. For this, we defined a number of segmentation quality metrics that can be applied to multichannel fluorescence images. We calculated these metrics for 14 previously described segmentation methods applied to datasets from four multiplexed microscope modalities covering five tissues. Using principal component analysis to combine the metrics, we defined an overall cell segmentation quality score and ranked the segmentation methods. We found that two deep learning-based methods performed the best overall, but that results for all methods could be significantly improved by postprocessing to ensure proper matching of cell and nuclear masks. Our evaluation tool is available as open source and all code and data are available in a Reproducible Research Archive.

Monitoring Editor

Diane Lidke
University of New Mexico

Received: Aug 25, 2022

Revised: Nov 28, 2022

Accepted: Dec 7, 2022

INTRODUCTION

Cell segmentation is the task of defining cell boundaries in images. It is a fundamental step for many image-based cellular studies, including analysis and modeling of subcellular patterns (Boland and Murphy, 2001), analysis of changes upon various perturbations (Carpenter *et al.*, 2006), cell tracking for investigating cell migration and proliferation (Garay *et al.*, 2013), and cell morphology analysis for discovering cell physiological states (Rittscher, 2010). Inaccurate cell segmentation introduces potential systematic error in all these downstream analyses. Because different methods may perform differently for different imaging modalities or tissues, it is important to evaluate potential cell segmentation methods to choose the most suitable for a specific application.

This article was published online ahead of print in MBoC in Press (<http://www.molbiolcell.org/cgi/doi/10.1091/mbc.E22-08-0364>) on December 14, 2022.

*Address correspondence to: Robert F. Murphy (murphy@cmu.edu).

Abbreviations used: CODEX, CO-Detection by indEXing; HuBMAP, Human Bio-Molecular Atlas Program; IMC, Imaging Mass Cytometry; MIBI, Multiplexed Ion Beam Imaging; TMC, Tissue Mapping Center.

© 2023 Chen and Murphy. This article is distributed by The American Society for Cell Biology under license from the author(s). It is available to the public under an Attribution 4.0 International Creative Commons CC-BY 4.0 License (<https://creativecommons.org/licenses/by/4.0/>).

“ASCB®,” “The American Society for Cell Biology®,” and “Molecular Biology of the Cell®” are registered trademarks of The American Society for Cell Biology.

Existing segmentation methods can be divided into two categories, geometry-based segmentation techniques (i.e., traditional computer vision techniques) and deep learning-based approaches. The former include but are not limited to threshold-based segmentation (Shen *et al.*, 2018), region-based segmentation (Panagiotakis and Argyros, 2018), the watershed algorithm and its variants (Ji *et al.*, 2015), active contours (Wu *et al.*, 2015), Chan–Vese segmentation (Fan *et al.*, 2013; Braiki *et al.*, 2020), and graph-cut based segmentation (Oyebode and Tapamo, 2016). In the deep learning category, conventional deep convolutional neural networks (CNNs) were initially applied to various cell segmentation tasks (Jung *et al.*, 2019; Sadanandan *et al.*, 2017). CNN models learn the feature mapping of an image and convert the feature map into a vector for pixel-wise classification (i.e. segmentation). Since its publication, the U-Net model and its variants have become a widely used alternative (Ronneberger *et al.*, 2015; Al-Kofahi *et al.*, 2018; Falk *et al.*, 2019; Long, 2020). Instance-segmentation methods such as Mask R-CNN, in addition, are prominently utilized for cell and nucleus segmentation (Johnson, 2018; Lv *et al.*, 2019; Fujita and Han, 2020). Ensemble methods with multiple deep learning frameworks were also developed specifically for nuclei segmentation (Vuola *et al.*, 2019; Kablan *et al.*, 2020).

The typical approach to evaluating cell segmentation methods is to compare the segmentation results with human-created

segmentation masks. For example, early work (Bamford, 2003) evaluated geometry-based algorithms on 20,000 cellular images annotated by three independent nonexpert observers. The author argued that evaluating cell segmentation does not necessarily require expert annotation. Caicedo *et al.* (2019) evaluated multiple deep learning strategies for cell/nuclei segmentation by calculating the accuracy and types of errors in contrast to the expert annotations. The authors created a prototype annotation tool to facilitate the annotation process. At the model development level, deep learning-based segmentation methods design loss functions to evaluate the similarity between the current segmentation mask and the human-created mask (Al-Kofahi *et al.*, 2018; Kromp *et al.*, 2020), which makes the model mimic the logic of human segmentation. Unfortunately, comparison with human-created segmentations assumes that people are good at this task; however, human segmentation can suffer from extensive intra- and interobserver variability (Wiesmann *et al.*, 2017; Vicar *et al.*, 2019). In addition, human segmentation can require extensive labor and cost. For example, the development of the latest DeepCell version (Greenwald *et al.*, 2022) required a very large number of human evaluations (and would have required far more if an active machine learning approach had not been used).

An alternative approach is creating simulated cell images (Wiesmann *et al.*, 2017) to assess segmentation performance. The simulated images, along with their segmentation results, are generated based on various information obtained from real fluorescent images including cell shape, cell texture, and cell arrangement. Because the image is simulated, the correct segmentation is known. Even though this method is efficient and reproducible, it reaches its limits while encountering images with high cell shape and texture variability.

For this study, we sought to define an objective evaluation approach that does not require a reference segmentation. We were motivated in large part by the desire to optimize the pipeline used for analysis of multichannel tissue images as part of the Human BioMolecular Atlas Program (HuBMAP; HuBMAP Consortium, 2019). We first defined a series of metrics based upon assumptions about the desired characteristics of good cell segmentation methods. We then identified currently available cell segmentation methods that had pretrained models and evaluated their performance on many images from multichannel imaging modalities. A principal component analysis (PCA) model was then trained using the metrics computed from the segmentation results of 14 methods on 637 multichannel tissue images across four imaging modalities and used to generate overall segmentation quality scores (Figure 1). We also evaluated the robustness of each method to various image degradations, such as adding noise. We found that as distributed, the Cellpose model (Stringer *et al.*, 2021) gave the best results, but that after postprocessing to ensure proper matching of cell and nuclear masks, two different DeepCell models (Bannon *et al.*, 2021) performed the best. We also found that our evaluation metrics are sensitive to undersegmentation error, which is very common in practice. Last, we found that our quality scores, which are obtained without the help of any human reference, not only capture the interobserver variance between two human experts, but also have a high correlation with three cell segmentation benchmarks using expert annotations.

To enable use by other investigators, we provide open source software that is able to evaluate the performance of any cell segmentation method on multichannel image inputs.

RESULTS

Generating masks and calculating metrics

We began by running all methods on images from four imaging modalities: CODEX, Cell DIVE, MIBI, and IMC (see *Methods* and Supplemental Table 1). There were 637 multichannel images in total. Each method generated whole-cell masks, and some methods also generated nuclear masks. For those methods that did not generate a nuclear mask, we provided a simple mask based on Otsu thresholding of the nuclear channel. For each method, an additional pair of masks was created with our “repair” procedure to eliminate unmatched nuclear and cell masks and remove nuclear regions outside the corresponding cell mask (see *Mask Processing* in *Methods*). This was done because unmatched cells and nuclei would be penalized by our evaluation metrics; the combination of the original method and the repair procedure was treated as a separate method to allow evaluation of each method either as originally provided or as most suitable for cell quantitation.

We then ran the evaluation pipeline to calculate the evaluation metrics (see Table 2 and the Supplemental Methods) for each method for each image. This process yielded two matrices of 637 images \times 14 methods \times 14 metrics, one for the original methods and one for the methods with repair.

To examine how the individual metrics varied across methods, we performed z-score standardization on each metric and averaged the metrics over all images for each method (Figure 2). We observed that the 14 methods have heterogeneous performance on the different metrics with both nonrepair and repair approaches. Methods after repair tend to have a significant increase in the average metric values (Figure 2A vs. Figure 2B), and the improvement of the FMCN metric directly reflects significantly better cell-nuclei matching after repair. We noticed that methods before repair have slightly higher cell uniformity (reflected by higher $1/[ACVC+1]$, FPCC, and AS metrics) due to a smaller number of cells with matching cell and nuclear masks. In both nonrepair and repair figures, the patterns of curves for all methods except the Voronoi are similar, with curves of deep learning-based methods similar to each other. We also noticed that different methods sometimes trade off a decrease in one or more metrics for an increase in others. For instance, DeepCell 0.6.0 and DeepCell 0.9.0 have opposite behavior on the NC and FFC metrics in Figure 2B. This comes from the fact that despite similar deep learning algorithms being applied, DeepCell 0.6.0 tends to segment fewer but larger cells than DeepCell 0.9.0. Our metrics accurately reflect this phenomenon (relatively lower NC but higher FFC). Our metrics also capture the improvement of DeepCell 0.12.3, which has a much larger training dataset than two previous versions, with both relatively high NC and FFC. We also observed that different metrics show different ranges of variance for the 14 methods. For instance, there is large variation among methods with respect to the NC metric (number of cells per 100 squared micrometers) and the FFC metric (fraction of image foreground occupied by cells) but less variation in $1/(ACVC_NUC+1)$ values (average of weighted average CV of cell type intensities over 1–10 clusters on nuclei).

Evaluating sensitivity of methods to added noise or downsampling

To test the robustness of the methods, we created perturbed images with various amounts of added zero-mean Gaussian noise or various extents of downsampling (see *Methods*) and evaluated the quality of the resulting segmentations. The Gaussian perturbations were done only to the images provided to the segmentation method, and not to the multichannel images used to calculate

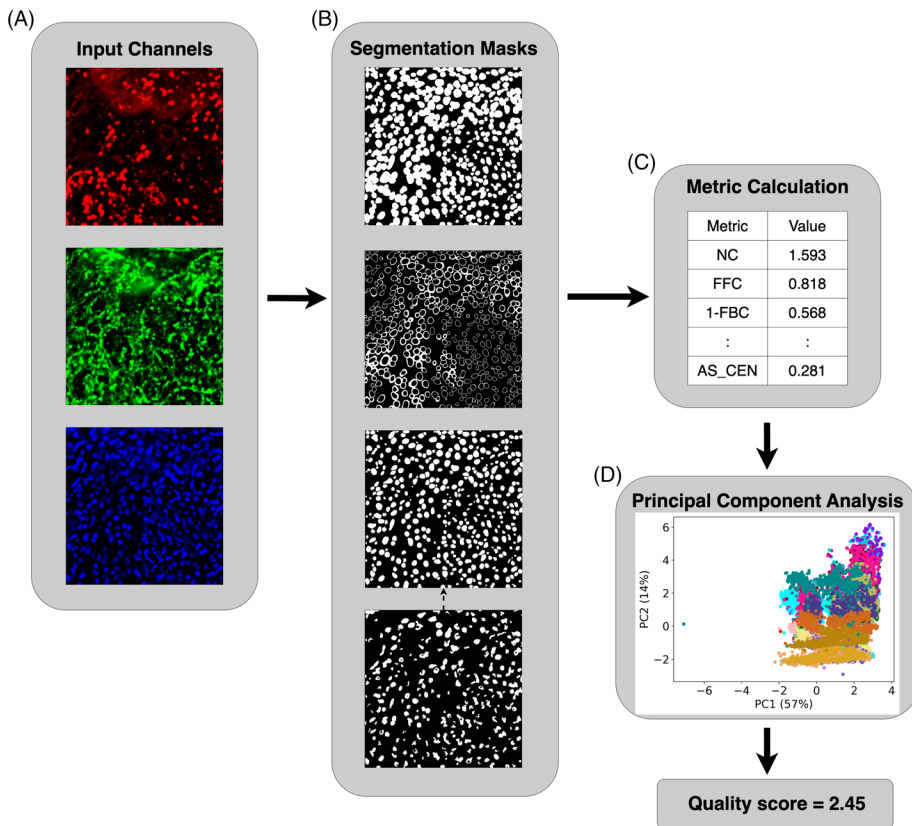


FIGURE 1: Our pipeline for cell segmentation evaluation. (A) Input channels for cell membrane (red), cytoplasm (green), and nucleus (blue) are available for each segmentation method, but methods only used one or two channels. (B) Methods generate segmentation masks: cell mask, cell outside nucleus mask, and nuclear mask (shown in the top three panels, respectively). For segmentation methods that do not output the nuclear mask, we used an Otsu-thresholded nuclear mask (bottom panel) as a substitute. We removed the pixels in the nuclear mask that were outside the cell membrane in the cell mask, as well as the cells and nuclei that did not have the corresponding nuclei and cells. (C) For each set of segmentation masks, we calculated 14 metrics to evaluate the quality of segmentation objectively under various assumptions. We then applied principal component analysis to the matrix of all metrics for all methods and all images. (D) The scatterplot shows a point for each segmentation for each image. Different colors represent different segmentation methods. Finally, a variance-weighted sum of PC1 and PC2 was used to generate an overall quality score for each combination of method and image.

the metrics. The downsampling was also done on the multichannel images, matching the same size as segmentation masks to calculate the metrics. This allowed us to assess method robustness: how sensitive the results for a given segmentation method were to the quality of the image. It also provided a parallel check of how well our metrics were performing based on the assumption that a perturbed image should yield a worse segmentation.

Segmentation quality score

One of our primary goals was to provide an overall segmentation quality score for each method on each image. To do this we created a principal component analysis (PCA) model using the metrics for all methods for all images with and without perturbation. The metrics were z-scored before PCA. The top two principal components (PCs) for each method across all modalities and images are shown in Figure 3 (and Supplemental Figure 3) with and without various amounts of perturbation. Because the number of datasets available for each modality varied, we averaged the top 2 PCs across all images within each modality and then averaged across

the modalities to balance the contributions from all modalities to the final model. With random Gaussian perturbation (Figure 3A and Supplemental Figure 3A), we observed that PC1 values for all methods tend to decrease as perturbation increases, confirming that the metrics perform as anticipated. We also observed that most of methods show decreased PC1 with increased downsampling (Figure 3B and Supplemental Figure 3B), with some exceptions. The exceptions occur because the downsampling process also removed noise in the images, causing some methods to get a slightly better PC1 score after downsampling. Thus, PC1 alone is not sufficient as an overall quality score.

The loadings of PC1 (Supplementary Figure 4A) are all positive, showing that all metrics have a synergistic effect on PC1 values. On the other hand, PC2 loadings (Supplementary Figure 4B) show that it is primarily an indicator of the overall coverage of a mask with high NC and FFC loadings. We therefore adopted the sum of PC1 and PC2 weighted by their explained variance as our final overall quality score. Rankings based on the quality score for all methods (with and without repair) averaged across all modalities and (unperturbed) images are shown in Figure 4, with Voronoi segmentation as a reference baseline. As expected, repair generally improves the overall metric for most methods. We observed that DeepCell methods have the highest overall performance with repair, and that they are also robust to both random Gaussian noise and downsampling (Figure 3). Cellpose methods perform the best among methods without repair. CellProfiler scores the best among methods with non-deep learning algorithms. Methods that are primarily used for cultured cell segmentation rather than tissue segmentation (e.g., CellX, AICS classic) tend to perform worse than the Voronoi baseline.

Model selection for particular tissue and data modality

One of the major potential uses of our evaluation pipeline is to select the most appropriate method for a specific imaging modality or tissue. Images from different modalities have a large variation in signal-to-noise and spatial resolution that may influence the performance of segmentation methods. We therefore separately evaluated images from different modalities (Supplemental Figures 5–8). While it performed the best across different modalities, the version of DeepCell that was optimal varied. Perhaps because they had the lowest resolution among the four modalities (1 micrometer pixel size), the relative ranking of the methods was quite different for IMC images than for other modalities. We also observed that CellProfiler ranked in the middle of DeepCell methods for CellDIVE and IMC modalities.

Segmentation method performance may also be expected to vary for different tissues due to potential differences in the shapes

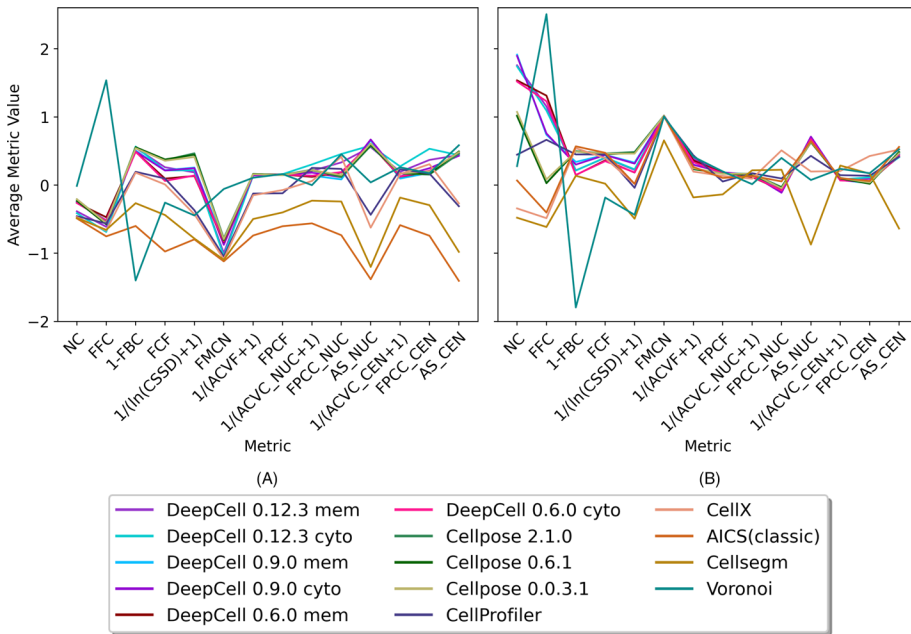


FIGURE 2: Heterogeneity of methods performance. Each metric was z-scored. For each segmentation method, each metric was averaged across all images. Panel (A) shows results from methods evaluated without mask repair (see text). Panel (B) shows the results from methods after the repair.

and spatial arrangement of cells. We therefore separately analyzed performance on five tissues: small intestine, large intestine, spleen, thymus, and lymph node (images for the first two tissues are solely from CODEX, and the latter three are available for both CODEX and IMC; Supplemental Figures 9-13). We observed that while DeepCell

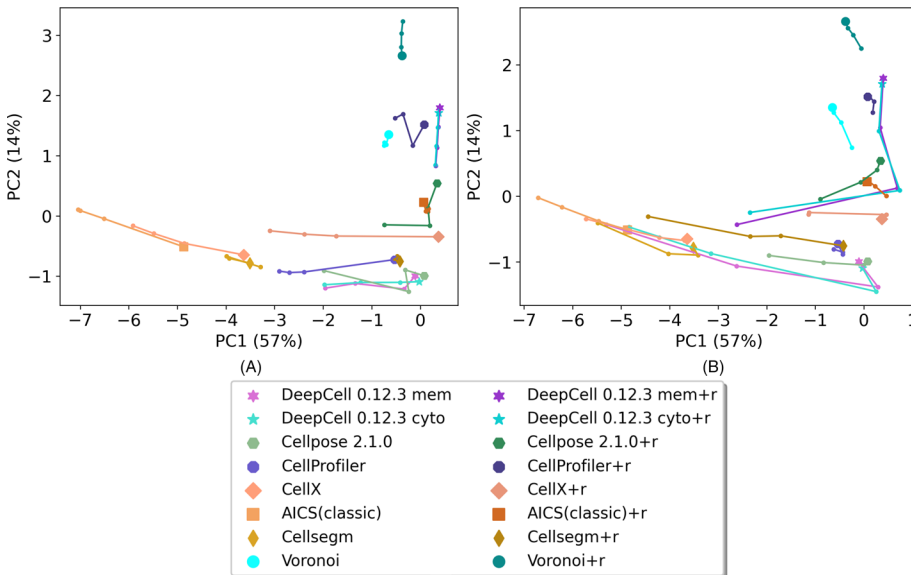


FIGURE 3: Top two principal components of each method on images from all modalities. (A) Results from images with Gaussian noise perturbation. (B) Result from images with downsampling. Each method is represented by a unique color. The points with a unique marker shape represent unperturbed images. The trajectories represent low, medium, and high Gaussian noise perturbation, A, or small, medium, and large degree of downsampling, B. Only the results from the latest DeepCell and Cellpose are shown, for better visualization. Supplemental Figure 3 shows results from the earlier versions.

models still performed the best among all methods (with the latest version being optimal), the rankings and overall quality scores of methods on different tissues vary. Interestingly, strong similar performance was observed for both DeepCell and Cellpose on large intestine.

Evaluating expert annotated segmentation

Because previous cell segmentation studies have focused on comparison with human experts, we used two lymph node CODEX multichannel images (tiles R001_X003_Y004 and R001_X004_Y003 in dataset HBM279. TQRS.775) and accompanying cell segmentation masks annotated by two experts (Dayao et al., 2022) to compare performance of various segmentation methods. We processed those images using all segmentation methods listed in Table 1 and evaluated the segmentation outputs along with expert-annotated masks. Because the expert annotations did not include nuclear masks, we slightly modified our evaluation metrics to calculate the cell uniformity metrics (i.e., $1/[ACVC + 1]$, FPCC, and AS, in Table 2) only on the cell masks and removed the FMCN metric. Accordingly, we retrained the PCA model (using all datasets) with the reduced number of metrics (10 instead of 14) for use only in this comparison with expert annotations.

The quality score ranking and top 2 PCs plot in Figure 5A reflect relatively high-accuracy expert annotation and also reveal the interobserver variability we describe in the Introduction. While Expert1 is among the best in the ranking, Expert2 has an overall score below the baseline. This emphasizes that while the cell segmentation annotated by experts may be helpful in many ways, it should be used with caution for the cell segmentation task.

We next sought to provide confidence that our quality scores obtained without a human expert would produce reliable results comparable to those for measures requiring human expert segmentation. To do this, we directly compared our quality scores with three benchmarks of cell segmentation quality (see *Methods*). These benchmarks were designed to be symmetric, meaning that the results are the same when either human annotation or computer segmentation is treated as the reference. The benchmarks were calculated for each method compared with each expert annotation in each image (Figure 6). We observed average Pearson correlation coefficients across the two images of 0.83, 0.81, and 0.72 between our quality score and the F1, Avg F1, and SEG' scores, respectively. Thus our quality scores are highly predictive of the benchmark scores that would have been obtained by

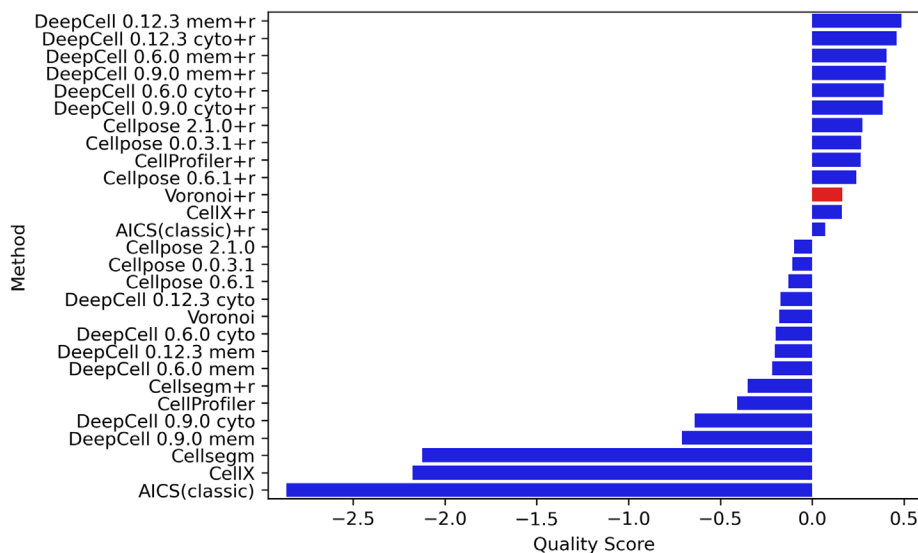


FIGURE 4: Overall quality score rankings of all modalities. All methods with and without repair are ranked by their quality scores. The methods that have a higher score than Voronoi are considered acceptable.

comparison to expert annotations. Specifically, the deep learning-based methods in the upper right corner of three benchmarks show higher similarity to both expert annotations in either image. The non-deep learning based methods, however, show lower agreement between quality score and either benchmark. This is because, while these benchmarks mainly measure the similarity between annotation and segmentation from a single perspective (i.e., overlapping), our quality score measures bad segmentation performance in various ways (e.g., missing cells, cells too small, unmatched cell, and nuclear masks). We also observed F1 scores around 0.7, Avg F1 scores around 0.37, and SEG' scores around 0.35 when the two expert annotations were compared, which are all lower than most deep

learning-based methods compared with either expert annotation. This interobserver variance in terms of segmentation quality again echoes the conclusion from Figure 5 and supports the value of our quality scores as a measure of cell segmentation accuracy.

Evaluating undersegmentation error

As shown in Figure 3, our quality scores appropriately reflect the degradation of cell segmentation quality that occurs after an input image is perturbed. As a final test of the sensitivity of our quality scores, we directly degraded cell segmentations by two approaches. The first was to simulate undersegmentation, in which adjacent or overlapping objects that should be distinct are merged. This often occurs in cell and nucleus segmentation, especially in dense tissues (Kromp et al., 2021). We simulated an undersegmentation scenario by merging cells in contact with each other. For this test, we started from masks generated by DeepCell 0.12.3, which is the best performer among all methods, on tiles from each of the CODEX datasets across all five tissue types, taking DAPI and E-cadherin as nuclear and cell membrane inputs (we selected one tile with enough cells to merge for each dataset). We sequentially merged pairs of contacting cells until the number of cells reached a target percentage (either 90% or 60%) of the original cell number (note that cells that were merged were not allowed to merge again with other cells, to better reflect typical undersegmentation performance, and that merging to 90% of the original cell number indicates that 20% of cells in the original mask have been merged). For each pair of nuclei belonging to the merged cells, we applied morphological opening on the inverted mask to connect the gap between the

learning-based methods compared with either expert annotation. This interobserver variance in terms of segmentation quality again echoes the conclusion from Figure 5 and supports the value of our quality scores as a measure of cell segmentation accuracy.

Method	Inputs			Requires parameters ^a	Output
	Cytoplasm	Cell membrane	Nucleus		Nuclear mask
DeepCell 0.12.3		X	X	Yes	Yes
DeepCell 0.12.3	X		X	Yes	Yes
DeepCell 0.9.0		X	X	No	Yes
DeepCell 0.9.0	X		X	No	Yes
DeepCell 0.6.0		X	X	No	Yes
DeepCell 0.6.0	X		X	No	Yes
Cellpose 2.1.0	X		X	No	Yes
Cellpose 0.6.1	X		X	No	Yes
Cellpose 0.0.3.1	X		X	No	Yes
CellSegm		X	X	Yes	No
CellX		X		Yes	No
CellProfiler		X	X	Yes	Yes
AICS (classic)		X	X	Yes	No
Voronoi			X	No	Yes

^aThe required parameters define the range of acceptable cell sizes.

TABLE 1: Segmentation methods evaluated.

Name	Metric	Mask(s) to calculate the metrics on
Number of cells per 100 squared micrometers	NC	Matched cell mask
Fraction of image foreground occupied by cells	FFC	Matched cell mask
Fraction of image background occupied by cells	1-FBC	Matched cell mask
Fraction of cell mask in foreground	FCF	Matched cell mask
Fraction of match between cells and nuclei	FMCN	Cell and nuclear masks prior to mask processing
Average CV of foreground pixels outside the cells	1/(ACVF+1)	Matched cell mask
Fraction of first PC of foreground pixels outside the cells	FPCF	Matched cell mask
Average of weighted average CV of cell type intensities over 1–10 clusters	1/(ACVC+1)	Matched nuclear mask (NUC) and cell excluding nucleus mask (CEN)
Average of weighted average fraction of the first PC of cell type intensities over 1–10 clusters	FPCF	Matched nuclear mask (NUC) and cell excluding nucleus mask (CEN)
Average of Silhouette score of clustering over 2–10 clusters	AS	Matched nuclear mask (NUC) and cell excluding nucleus mask (CEN)
Standard deviation of cell size	1/(ln(CSSD)+1)	Matched cell mask

TABLE 2: Summary of segmentation metrics.

nuclei to also simulate undersegmented nuclei (see Supplemental Figure 14).

As a second approach to degrading segmentation masks, we shifted the masks produced by DeepCell 0.12.3 relative to the original image. We shifted them by 0.1%, 1%, and 50% of the average of the two dimensions of the image (for a 1000 × 1000-pixel image, shifting 0.1% is 1 pixel right and 1 pixel down). Presumably, results for 0.1% shift should be lower than but similar to the original performance, while the scores of the ones shifted 50% should be among the lowest.

We then ran our evaluation pipeline using DeepCell 0.12.3 on these degraded masks and compared the results with quality scores of all methods on the original images. The results in Figure 7 show that the quality scores reflect the degree of degradation produced. The breakdown of metrics in Supplemental Figure 15 illustrates that our evaluation method is sensitive to undersegmentation, but the degree of sensitivity expected depends upon the extent of variation

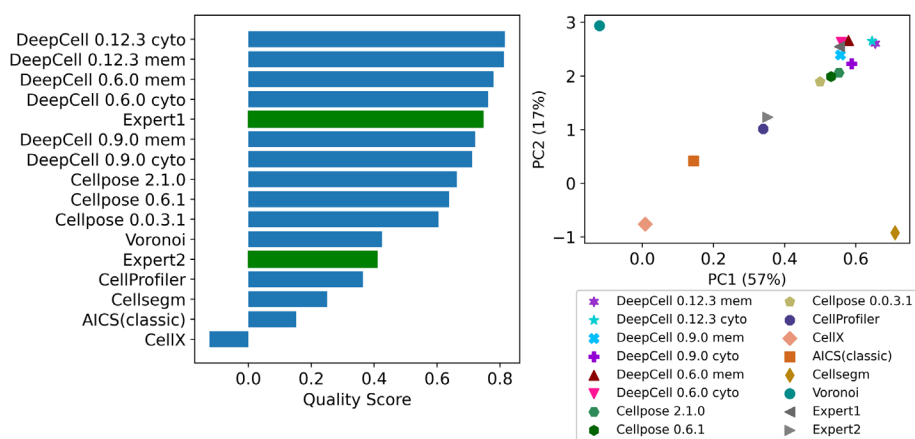


FIGURE 5: Quality score ranking and top 2 PCs of two expert-annotated segmentations comparing with results from all segmentation methods in Table 1. The accuracy of annotations from the Expert1 was reflected by our PCA model showing both high PC1 (good overall quality) and high PC2 (good tissue coverage) to the upper right corner among the group of deep learning-based models. The discrepancy between two experts was also captured, showing interobserver variability in the cell segmentation task.

in channel intensities among different cells (see *Discussion*). Notably, our evaluation method captured the degradation from merely shifting 0.1% with a slightly lower quality score than the original.

Evaluating similarity between results for different segmentation methods

Besides scoring each segmentation method separately, we also directly compared the segmentations produced by each pair of methods. We used two approaches: calculating a normalized distance between the single-method metric vectors of two methods, and calculating a set of metrics from direct comparison of the segmentation masks from two methods (see *Supplemental Methods*). In each approach, we concatenated the metric matrices across all images from multiple data modalities for all pairs of methods and applied PCA to obtain PC1 values as a difference score. These scores were normalized by subtracting the PC1 value of a method compared with itself (i.e., the lowest possible difference metric value). Figure 8 shows heat maps of the difference between the methods for the two sets of metrics. Higher values indicate more difference (lower similarity) between pairs of methods. Using the single-method metrics (Figure 8A), all the deep learning-based models (DeepCell and Cellpose) are relatively close to each other, consistent with the results shown above. The results for DeepCell with the membrane and cytoplasm markers are also close. However, pairwise comparisons (Figure 8B) of the masks show larger differences in the results for DeepCell and Cellpose (including between the three versions of DeepCell) and somewhat larger differences between the results for membrane and cytoplasm markers.

DISCUSSION

Many cell segmentation methods have been described, but approaches to evaluating them have been limited. We have described here our design of a set of reference-free

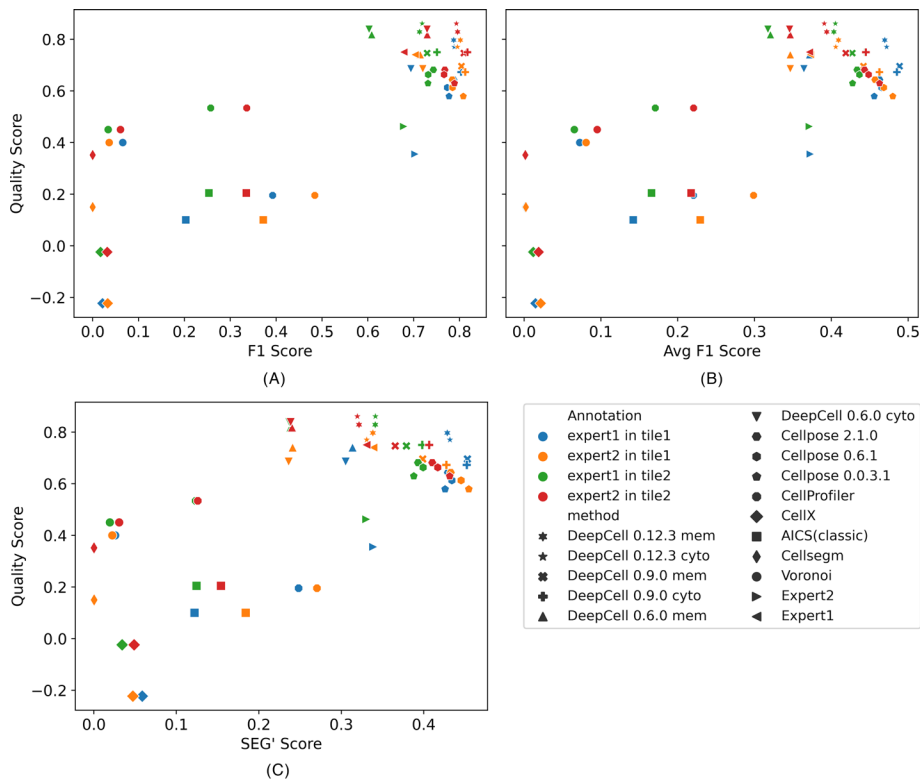


FIGURE 6: Comparing quality score with three benchmarks: F1 score (A), Avg F1 score (B) and SEG' score (C). Each marker represents a segmentation or an expert annotation. Each color represents the benchmark that was calculated with the marker's segmentation/annotation against an expert annotation in a CODEX image (tile). Pearson correlations between our quality score and three benchmarks are 0.83, 0.81, and 0.72, respectively.

metrics to comprehensively measure the performance of segmentation methods. The metrics were developed based on a series of assumptions about the desirable characteristics of cell segmentation, especially for multichannel images. We also trained a PCA model using these metrics to compare all pretrained segmentation models that we were able to identify. The results demonstrate the power of deep learning-based methods, with DeepCell and Cellpose performing the best. We also show that our metrics reflect the poorer performance expected for various image degradations, such as reducing pixel resolution, adding noise, and artificially increasing undersegmentation. Our evaluation approach also can be used to measure differences in the quality of different expert annotations. Our segmentation quality score highly correlates with three quality benchmarks that use expert annotations. Our open source tool has been incorporated into the image analysis pipeline for the HuBMAP project (HuBMAPConsortium, 2019).

The metrics are applicable across a range of image resolutions; the pixel size in the images we have used for evaluation range from 0.325 to 1 μm . Of course, the metrics (and the cell segmentation methods themselves) require an image resolution sufficient to adequately resolve single cells.

An important consideration to note is that our metrics are based in large part on the assumption that tissue images used for evaluation contain different cell types that vary in their expression of different markers. If a tissue image is largely uniform, then our measures of segmented cell homogeneity will yield similar results regardless of the accuracy of cell segmentation masks. The greater the number of channels measured and the larger the differences between cell

types, the more discriminating our metrics will be. The metrics also do not make any assumptions about allowable cell shapes. Measuring cell shape is an inherent problem for 2D tissue images where 3D cells are only seen in a thin layer in the z-axis. Cells belonging to the same cell type might be captured at different angles or heights and therefore display different shapes. Although our metrics of homogeneity at the cell level solve this issue indirectly (since cells from the same cell type but at different views should have similar channel compositions), segmentation methods that produce cells and nuclei with unexpected shapes are not directly penalized. Therefore, we plan in future work to incorporate metrics using spherical harmonic transform-based shape descriptors that have been shown to provide the best representation of cell and nuclear shapes (Ruan and Murphy, 2019).

We also note that while this study was inspired by the multichannel tissue images of the HuBMAP project, our metrics are also applicable to evaluate segmentation methods on cultured cells. However, as with tissue images, the sensitivity of the metrics depends on the number of channels and the diversity among individual cells.

To make our approach widely available, we provide an open source pipeline for calculating the segmentation metrics for a given set of images and a given segmentation method. Our evaluation pipeline provides a platform for users to choose segmentation methods for an individual image, tissue, and/or imaging modalities. While prior algorithmic approaches may have claimed the highest accuracy against different manually annotated training datasets, our method directly benchmarks them under the same set of measures.

The use of three-dimensional multichannel tissue imaging is gradually growing. Some of the segmentation methods we have evaluated are capable of both 2D and 3D segmentation (see *Methods*), whereas algorithms such as the deep learning version of the Allen cell structure segmenter (Chen et al., 2020) and nnUNet (Isensee et al., 2021) only focus on segmenting 3D cell images. We are in the process of extending our evaluation work to 3D segmentation.

We note that in the future our quality scores could be included in the loss functions for training cell segmentation models in order to further improve model performance. They also provide a useful measure of image quality, since, within a given tissue and modality, better-quality images can be expected to provide better segmentation results. This could potentially allow tissue images to be screened before acceptance by large projects such as HuBMAP.

METHODS

Images

We obtained image data from the HuBMAP project for four multiplexed modalities. For each modality, the nuclear, cytoplasmic, and cell membrane channels chosen for segmentation (Supplemental Table 1) were either those recommended by the Tissue Mapping Center (TMC) that produced the datasets, or selected based on

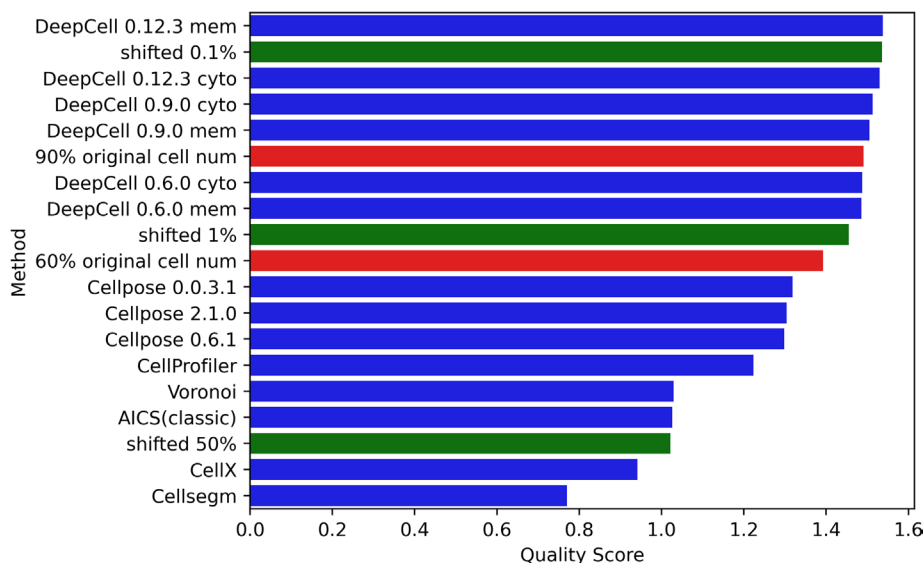


FIGURE 7: Quality score ranking of simulated masks with a range of undersegmentation error and shifting. Undersegmented masks (red) and shifted masks (green) were created from masks produced by DeepCell 0.12.3 using a cell membrane marker as described in the *Methods*.

peer-reviewed literature and subcellular localization of antibody targeted protein annotated by UniProt or Gene Ontology. Details about input markers are summarized in Supplemental Table 2.

CODEX

CO-Detection by indEXing (CODEX) is an imaging technique generating highly multiplexed images of fluorescently tagged antigens (Goltsev *et al.*, 2018). The HuBMAP portal (<https://portal.hubmap-consortium.org/>) contains 26 CODEX datasets on five tissues (large intestine, small intestine, thymus, spleen, and lymph node) produced by two TMCs: Stanford University and the University of Florida. The datasets from the Stanford TMC all contain 47 channels. Each dataset consists of four tissue regions, which are individually split into a grid of tiles with size $\sim 1440 \times 1920$. The University of

Florida TMC generated 11-channel CODEX datasets. Each dataset has only one region, which is divided into a grid in the same manner. To ensure both tissue coverage and evaluation efficiency, we created evaluation datasets with a subset of the tiles chosen to be non-neighboring (to avoid evaluating cells in the overlap region twice) and tiles that were not on the edge of the grid (to prevent large fractions of the image consisting of background on the edge of tissue). Three hundred ninety-three tiles were selected in total. The pixel size of CODEX images is $0.37745 \mu\text{m}$.

The different antigens used by the two TMCs required us to select different channels as segmentation input. For Stanford TMC datasets, we used the Hoechst, Cytokeratin, and CD45 channels as nuclear, cytoplasmic, and cell membrane channels, respectively. For the University of Florida TMC datasets, we used DAPI, CD107a, and E-cadherin. To ensure consistency of evaluation results from CODEX data, only the five channels that are common between the datasets of the two TMCs were used to calculate channel homogeneity metrics (see the Supplemental Methods).

peer-reviewed literature and subcellular localization of antibody targeted protein annotated by UniProt or Gene Ontology. Details about input markers are summarized in Supplemental Table 2.

Cell DIVE

Cell DIVE is another antibody-based multiplexing technique (Gerdes *et al.*, 2013). We had access to 12 regions of Cell DIVE data. Each region contains 26 images with 19 channels. Cell DIVE datasets are much larger than those from CODEX, consisting of $\sim 10,000 \times \sim 15,000$ pixels. To boost the efficiency of the pipeline while ensuring the coverage of datasets, we took the first image of each region and split it into a grid of tiles similar to those of CODEX datasets. Each tile has roughly 1000×1000 pixels, which is efficiently applicable to all segmentation methods. To exclude tiles with few cells, we calculated three metrics on each channel of each

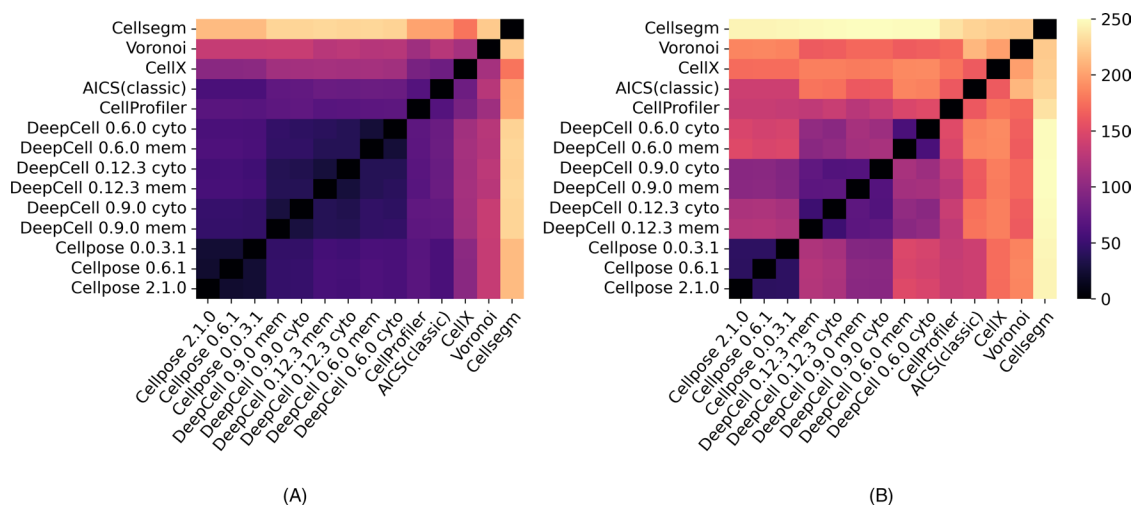


FIGURE 8: Heat maps of difference scores between segmentation methods. A score of zero represents no difference (i.e., identical segmentations). Larger (lighter) values indicate greater differences between methods. (A) Heat map generated by the difference between metrics of two single methods. (B) Heat map generated by the pairwise metrics that directly compare two methods' segmentations. The methods are ordered by clustering.

tile and applied K-means clustering to identify a cluster consisting of images with dense cells (see Supplemental Figure 1). We selected DAPI, cytokeratin, and P-cadherin as nuclear, cytoplasmic, and cell membrane inputs for segmentation. The pixel size of Cell DIVE images is 0.325 μm .

MIBI

MIBI (Multiplexed Ion Beam Imaging) uses a secondary-ion mass spectrometer (SIMS) to image antibodies tagged with monoisotopic metallic reporters (Angelo *et al.*, 2014). We obtained ten 29-channel MIBI datasets. Each image consists of 1024×1024 pixels. HH3, Pan-Keratin, and E-cadherin were utilized as the nuclear, cytoplasmic, and cell membrane channels for segmentation. The spot size (equivalent to pixel size) of MIBI images is 0.391 μm .

IMC

As an expansion of mass cytometry, imaging mass cytometry (IMC) uses laser ablation to generate plumes of particles that are carried to the mass cytometer by a stream of inert gas (Chang *et al.*, 2017). Thirteen out of a total of 2D 39-channel IMC datasets of spleen, thymus, and lymph node were available from the University of Florida TMC. In this modality, Ir191, SMA, and HLA-ABC were chosen as input channels for segmentation. For images without Ir191 channel, Histone channel was substituted as nuclear input. The pixel size of IMC images is 1 μm .

SEGMENTATION METHODS

We created a python wrapper for each method to adapt all methods to a common pipeline. The input channels required by each method are shown in Supplemental Table 1. All methods generate a segmentation mask for cell boundaries in an indexed image format with some methods also generating nuclear boundaries.

DeepCell. DeepCell is designed for robustly segmenting mammalian cells (Van Valen *et al.*, 2016; Moen *et al.*, 2019; Bannon *et al.*, 2021; Greenwald *et al.*, 2022). A feature pyramid network (Lin *et al.*, 2017) with PanopticNets architecture trained by more than 1 million paired whole-cell and nuclear annotations is embedded in the latest deepcell-tf package on the Python TensorFlow platform. We tested three versions of DeepCell software with different trained deep learning models. For each model, a nuclear intensity channel is a mandatory input and we chose cytoplasm or cell membrane as the secondary input image to generate two different whole-cell segmentation masks (these were treated as different methods). The latest version of DeepCell (v0.12.3) requires specification of pixel size in micrometers.

Cellpose. Cellpose is a generalist, U-Net-based algorithm that segments cells in various image modalities (Stringer *et al.*, 2021). The key idea of Cellpose is to track the gradient flow in the labels to predict cell boundaries. A PyTorch model "cyto" was trained with multiple datasets mainly with cytoplasmic and nuclear markers. We applied three versions of Cellpose with different pretrained models to segment whole cells on 2D cytoplasmic and nuclear intensity images. Unlike DeepCell, which uses the same model to generate whole-cell and nuclear masks, Cellpose contains a separate U-Net "nuclei" model for the segmentation of nuclear channels.

CellProfiler. CellProfiler 4.0 (Carpenter *et al.*, 2006; Kametsky *et al.*, 2011) contains a module for cell segmentation that applies traditional thresholding and propagation algorithms to segment nuclei and then cells. Required parameters that define the range of

acceptable nuclear sizes and cytoplasmic thickness were chosen for each imaging dataset. To facilitate batch computation and avoid CellProfiler's eight-bit cell index limitation in saving output images, we modified the CellProfiler module to store the indexed segmentation masks directly as NumPy arrays.

CellX. CellX is a MATLAB package that using traditional image processing operations and requires only a membrane marker image as input (Mayer *et al.*, 2013; Dimopoulos *et al.*, 2014). Required parameters that define the range of acceptable cell radii and maximum cell length (since cells are not perfect circles) were chosen for each dataset. Because the algorithm does not generate a nuclear mask, we paired the cell mask with a default Otsu-thresholded nuclear mask (see below under *Voronoi Segmentation*).

Cellsegm. Cellsegm is a MATLAB-based toolbox providing automated whole-cell segmentation (Hodneland *et al.*, 2013). Required parameters that define the range of acceptable cell sizes were chosen for each dataset. Because the algorithm does not generate a nuclear mask, we paired the cell mask with a default Otsu-thresholded nuclear mask.

The Allen Cell Structure Segmenter. There are two branches of the Allen Cell Structure Segmenter (AICS): a deep learning version and a classic image processing version (Chen *et al.*, 2020). The deep learning version of AICS only provides nuclear segmentation, so it was not evaluated. The classic version consists of thresholding and filtering operations and requires estimated minimum cell areas (which were chosen for each dataset based on pixel size). Because the algorithm does not generate a nuclear mask, we paired the cell mask with a default Otsu-thresholded nuclear mask.

Voronoi segmentation. To provide a baseline method, we used a simple method based on the Voronoi diagram. Otsu thresholding was applied on the nucleus channel as a first step to obtaining a nuclear mask. Then a Voronoi diagram was created to partition the image along lines equidistant from the centroid of each nucleus.

Degrading images for robustness analysis. To evaluate the robustness of each method, each image was degraded in two ways. In the first, zero-mean Gaussian noise was added to each pixel. The SD was set to various levels based on the typical channel intensity of a given modality: 500, 1000, and 1500 for CODEX and Cell DIVE images, and 5, 10, and 15 for MIBI and IMC data. For the second perturbation, we downsampled the images to 70%, 50%, and 30%, respectively, on both dimensions. Note that the Gaussian noise was only added to the images used for segmentation; the evaluation of the resulting masks was done using the original images.

Mask processing. For some segmentation methods, finding cell boundaries is done independent of finding nuclear boundaries. This may mean that the final segmentation masks include nuclei that do not have a corresponding cell boundary and vice versa. We assumed that a good segmentation method would minimize this and therefore defined a metric to capture this aspect of a segmentation method. To calculate it, each cell was matched to any nuclei contained within it. All cells that did not have corresponding nuclei were counted as mismatched, and vice versa. For cells that had multiple corresponding nuclei, the one with the smallest fraction of mismatched pixels was kept. All mismatched cells and nuclei were removed from the calculation of other metrics (see Supplemental Figure 2).

Segmentation methods may also generate misshaped nuclei that have pixels outside their corresponding cells—that is, nuclei that protrude through the cell membrane. Across all methods in all images, this occurred for an average of 60.6% of segmented nuclei. To solve this issue, we applied two posterior approaches. The first approach considered all misshaped nuclei and their corresponding cells to be mismatched even if they had a one-to-one relationship. Alternatively, we developed a “repair” pipeline that trimmed the mask of misshaped nuclei. The combination of a segmentation method followed by repair was evaluated as a distinct segmentation method.

To evaluate the segmentation performance using the different channel intensities, we expected that nuclear protein composition would be considerably different from the cytoplasm and the cell membrane. We therefore calculated our metrics using two masks: the (repaired) nuclear mask and a “Cell Excluding Nucleus” mask calculated by removing the nuclear mask from the cell mask.

After this mask processing step, each cell has a one-on-one matching relationship among its cell, nucleus, and cell excluding nucleus masks.

Evaluation metrics not requiring reference segmentation. We defined 14 metrics to evaluate the performance of a single segmentation method without requiring a reference segmentation. These are of two types: metrics that assess the coverage of a segmentation mask on the image, and metrics that measure various types of uniformity at the pixel and cell levels on multiplexed images. Each metric is derived under an assumption based on general concepts from cell biology. They are described in the Supplemental Methods and summarized in Table 2.

We also defined 10 metrics for comparing two segmentation methods (or one method with a human-generated segmentation). These are also described in the Supplemental Methods.

Benchmarks for comparing segmentation with annotations. We adapted three benchmarks to quantify the segmentation quality comparison with expert annotations. The F1 score has been widely used for benchmarking cell segmentation performance (Caicedo *et al.*, 2019; Greenwald *et al.*, 2022). The first step in calculating the F1 score is to determine matched pairs of cells from two masks. For each cell on the reference mask (which could be either mask), we calculated the Jaccard index (as the equation below) with overlapping cells in the other mask:

$$JI(R, S) = \frac{|R \cap S|}{|R \cup S|}$$

We set $JI = 0.3$ as the threshold that two cells must satisfy to be considered matched. If multiple cells have JI above the threshold, the one with the highest JI is selected. Based on this, we counted the number of true positive cells (TP, if two cells have JI above threshold) as well as the false positive and false negative cells (FP and FN, if two cells have JI below threshold). We calculated the F1 score by the following equations. Note that the F1 score remains the same regardless of which mask is the reference (swapping FP and FN will not change F1 score). Therefore, it is a symmetric measurement:

$$F_1 = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

where

$$\text{precision} = \frac{TP}{TP + FP}, \text{ recall} = \frac{TP}{TP + FN}$$

One issue with the F1 score is that it demands manual selection of the JI threshold. Different threshold choices may lead to different conclusions. To improve this, we averaged F1 scores with JI thresholds from 0 to 1 with 0.01 step sizes as our second benchmark (Avg F1 score).

The third benchmark, the SEG score, is widely applied for cell tracking (Maska *et al.*, 2014). The original SEG score scans each cell in the reference mask and finds the best matched cell in the query mask by the following condition:

$$|R \cap S| > 0.5 * |R|$$

where R is a reference cell and S is the segmented cell from a segmentation method. If the intersection (overlap) is greater than half the area of the reference cell, the two cells are considered matched. (Note that the constant 0.5 is analogous but not identical to setting a JI threshold.) If there are multiple segmented cells satisfying this condition, the one with the greatest overlap with R remains. For each reference cell, JI is used to qualify the similarity with its matched cell. For the ones with no matched cell, $JI = 0$. The final SEG score is the average of JI over all reference cells.

One major issue with the original SEG score definition comes from its asymmetry. If the query mask contains many extra cells that do not match any cell in the reference mask, the SEG score will be the same as if the query mask did not have any extra cells. However, segmenting extra unmatched cells is clearly worse. To solve this issue, we calculated two SEG scores, one using the human-labeled mask as reference and another using the computer-generated mask as reference, and averaged them as the final benchmark. We refer to this as the SEG' (SEG prime) score which is symmetric and accounts properly for unequal numbers of cells between the masks (and does not require an assumption of which mask is correct).

Availability. All code needed to calculate the evaluation metrics, as well as scripts to download test images, is available at <https://github.com/murphygroup/CellSegmentationEvaluator>. A reproducible research archive containing all code and intermediate results (which enable recreation of all figures and tables) is available at <https://github.com/murphygroup/ChenMurphy2DSegEvalRRA>.

ACKNOWLEDGMENTS

We thank Matthew Ruffalo for helpful discussions. This work was supported in part by a grant from the National Institutes of Health Common Fund, OT2 OD026682, and by a traineeship to H.C. under Training Grant T32 EB009403.

REFERENCES

- Al-Kofahi Y, Zaltsman A, Graves R, Marshall W, Rusu M (2018). A deep learning-based algorithm for 2-D cell segmentation in microscopy images. *BMC Bioinformatics* 19, 365.
- Angelo M, Bendall SC, Finck R, Hale MB, Hitzman C, Borowsky AD, Levenson RM, Lowe JB, Liu SD, Zhao S, *et al.* (2014). Multiplexed ion beam imaging of human breast tumors. *Nat Med* 20, 436–442.
- Bamford P (2003). Empirical comparison of cell segmentation algorithms using an annotated dataset. In: *Proceedings 2003 International Conference on Image Processing* (Cat. No. 03CH37429).
- Bannon D, Moen E, Schwartz M, Borba E, Kudo T, Greenwald N, Vijayakumar V, Chang B, Pao E, Osterman E, *et al.* (2021). DeepCell Kiosk: scaling deep learning-enabled cellular image analysis with Kubernetes. *Nat Methods* 18, 43–45.
- Boland MV, Murphy RF (2001). A neural network classifier capable of recognizing the patterns of all major subcellular structures in fluorescence microscope images of HeLa cells. *Bioinformatics* 17, 1213–1223.
- Braiki M, Benzinou A, Nasreddine K, Hymery N (2020). Automatic human dendritic cells segmentation using K-means clustering and Chan–Vese active contour model. *Comput Methods Programs Biomed* 195, 105520.

- Caicedo JC, Roth J, Goodman A, Becker T, Karhohs KW, Broisin M, Molnar C, McQuin C, Singh S, Theis FJ, Carpenter AE (2019). Evaluation of deep learning strategies for nucleus segmentation in fluorescence images. *Cytometry A* 95, 952–965.
- Carpenter AE, Jones TR, Lamprecht MR, Clarke C, Kang IH, Friman O, Guertin DA, Chang JH, Lindquist RA, Moffat J, et al. (2006). CellProfiler: image analysis software for identifying and quantifying cell phenotypes. *Genome Biol* 7, R100.
- Chang Q, Ornatsky OI, Siddiqui I, Loboda A, Baranov VI, Hedley DW (2017). Imaging mass cytometry. *Cytometry A* 91, 160–169.
- Chen J, Ding L, Viana MP, Lee H, Sluezwski MF, Morris B, Hendershott MC, Yang R, Mueller IA, Rafelski SM (2020). The Allen Cell and Structure Segmenter: a new open source toolkit for segmenting 3D intracellular structures in fluorescence microscopy images. *BioRxiv* 491035.
- Dayao MT, Brusko M, Wasserfall C, Bar-Joseph Z (2022). Membrane marker selection for segmenting single cell spatial proteomics data. *Nat Commun* 13, 1999.
- Dimopoulos S, Mayer CE, Rudolf F, Stelling J (2014). Accurate cell segmentation in microscopy images using membrane patterns. *Bioinformatics* 30, 2644–2651.
- Falk T, Mai D, Bensch R, Cicek O, Abdulkadir A, Marrakchi Y, Bohm A, Deubner J, Jackel Z, Seiwald K, et al. (2019). U-Net: deep learning for cell counting, detection, and morphometry. *Nat Methods* 16, 67–70.
- Fan J, Li S, Fan CZ, Zhang C (2013). Color cell image segmentation based on Chan–Vese model for vector-valued images. *J Software Eng Appl* 6, 554.
- Fujita S, Han X-H (2020). Cell detection and segmentation in microscopy images with improved mask R-CNN. In: *Proceedings of the Asian Conference on Computer Vision*.
- Garay T, Juhasz E, Molnar E, Eisenbauer M, Czirik A, Dekan B, Laszlo V, Hoda MA, Dome B, Timar J, et al. (2013). Cell migration or cytokinesis and proliferation?—revisiting the “go or grow” hypothesis in cancer cells in vitro. *Exp Cell Res* 319, 3094–3103.
- Gerdes MJ, Sevinsky CJ, Sood A, Adak S, Bello MO, Bordwell A, Can A, Corwin A, Dinn S, Filkins RJ, et al. (2013). Highly multiplexed single-cell analysis of formalin-fixed, paraffin-embedded cancer tissue. *Proc Natl Acad Sci USA* 110, 11982–11987.
- Goltsev Y, Samusik N, Kennedy-Darling J, Bhate S, Hale M, Vazquez G, Black S, Nolan GP (2018). Deep profiling of mouse splenic architecture with CODEX multiplexed imaging. *Cell* 174, 968–981.e915.
- Greenwald NF, Miller G, Moen E, Kong A, Kagel A, Dougherty T, Fullaway CC, McIntosh BJ, Leow KX, Schwartz MS, et al. (2022). Whole-cell segmentation of tissue images with human-level performance using large-scale data annotation and deep learning. *Nat Biotechnol* 40, 555–565.
- Hodneland E, Kogel T, Frei DM, Gerdes HH, Lundervold A (2013). CellSegm—a MATLAB toolbox for high-throughput 3D cell segmentation. *Source Code Biol Med* 8, 16.
- HuBMAP Consortium (2019). The human body at cellular resolution: the NIH Human Biomolecular Atlas Program. *Nature* 574, 187–192.
- Isensee F, Jaeger PF, Kohl SAA, Petersen J, Maier-Hein KH (2021). nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation. *Nat Methods* 18, 203–211.
- Ji X, Li Y, Cheng J, Yu Y, Wang M (2015). Cell image segmentation based on an improved watershed algorithm. 2015 8th International Congress on Image and Signal Processing (CISP).
- Johnson JW (2018). Adapting mask-rcnn for automatic nucleus segmentation. *arXiv preprint* <https://doi.org/10.48550/arXiv.1805.00500>.
- Jung H, Lodhi B, Kang J (2019). An automatic nuclei segmentation method based on deep convolutional neural networks for histopathology images. *BMC Biomed Eng* 1, 24.
- Kablan EB, Dogan H, Ercin ME, Ersoz S, Ekinci M (2020). An ensemble of fine-tuned fully convolutional neural networks for pleural effusion cell nuclei segmentation. *Comput Electrical Eng* 81, 106533.
- Kamentsky L, Jones TR, Fraser A, Bray MA, Logan DJ, Madden KL, Ljosa V, Rueden C, Eliceiri KW, Carpenter AE (2011). Improved structure, function and compatibility for CellProfiler: modular high-throughput image analysis software. *Bioinformatics* 27, 1179–1180.
- Kromp F, Bozsaky E, Rifatbegovic F, Fischer L, Ambros M, Berneder M, Weiss T, Ladic D, Dorr W, Hanbury A, et al. (2020). An annotated fluorescence image dataset for training nuclear segmentation methods. *Sci Data* 7, 262.
- Kromp F, Fischer L, Bozsaky E, Ambros IM, Dorr W, Beiske K, Ambros PF, Hanbury A, Taschner-Mandl S (2021). Evaluation of deep learning architectures for complex immunofluorescence nuclear image segmentation. *IEEE Trans Med Imaging* 40, 1934–1949.
- Lin T-Y, Dollár P, Girshick R, He K, Hariharan B, Belongie S (2017). Feature pyramid networks for object detection. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Long F (2020). Microscopy cell nuclei segmentation with enhanced U-Net. *BMC Bioinformatics* 21, 8.
- Lv G, Wen K, Wu Z, Jin X, An H, He J (2019). Nuclei R-CNN: improve mask R-CNN for nuclei segmentation. 2019 IEEE 2nd International Conference on Information Communication and Signal Processing (ICICSP).
- Maska M, Ulman V, Svoboda D, Matula P, Matula P, Ederra C, Urbiola A, Espana T, Venkatesan S, Balak DM, et al. (2014). A benchmark for comparison of cell tracking algorithms. *Bioinformatics* 30, 1609–1617.
- Mayer C, Dimopoulos S, Rudolf F, Stelling J (2013). Using CellX to quantify intracellular events. *Curr Protocols Mol Biol* 101, 14.22. 11-14.22. 20.
- Moen E, Borba E, Miller G, Schwartz M, Bannon D, Koe N, Camplisson I, Kyme D, Pavelchek C, Price T (2019). Accurate cell tracking and lineage construction in live-cell imaging experiments with deep learning. *Biorxiv* 803205.
- Oyebode KO, Tapamo JR (2016). Adaptive parameter selection for graph cut-based segmentation on cell images. *Image Anal Stereol* 35, 29–37.
- Panagiotakis C, Argyros AA (2018). Cell segmentation via region-based ellipse fitting. 2018 25th IEEE International Conference on Image Processing (ICIP).
- Rittscher J (2010). Characterization of biological processes through automated image analysis. *Annu Rev Biomed Eng* 12, 315–344.
- Ronneberger O, Fischer P, Brox T (2015). U-net: Convolutional networks for biomedical image segmentation. *International Conference on Medical Image Computing and Computer-Assisted Intervention*.
- Ruan X, Murphy RF (2019). Evaluation of methods for generative modeling of cell and nuclear shape. *Bioinformatics* 35, 2475–2485.
- Sadanandan SK, Ranefall P, Le Guyader S, Wahlby C (2017). Automated training of deep convolutional neural networks for cell segmentation. *Sci Rep* 7, 7860.
- Shen SP, Tseng HA, Hansen KR, Wu R, Gritton HJ, Si J, Han X (2018). Automatic cell segmentation by adaptive thresholding (ACSAT) for large-scale calcium imaging datasets. *eNeuro* 5. <https://doi.org/10.1523/ENEURO.0056-18.2018>
- Stringer C, Wang T, Michaelos M, Pachitariu M (2021). Cellpose: a generalist algorithm for cellular segmentation. *Nat Methods* 18, 100–106.
- Van Valen DA, Kudo T, Lane KM, Macklin DN, Quach NT, DeFelice MM, Maayan I, Tanouchi Y, Ashley EA, Covert MW (2016). Deep learning automates the quantitative analysis of individual cells in live-cell imaging experiments. *PLoS Comput Biol* 12, e1005177.
- Vicar T, Balvan J, Jaros J, Jug F, Kolar R, Masarik M, Gumulec J (2019). Cell segmentation methods for label-free contrast microscopy: review and comprehensive comparison. *BMC Bioinformatics* 20, 360.
- Vuola AO, Akram SU, Kannala J (2019). Mask-RCNN and U-net ensemble for nuclei segmentation. 2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019).
- Wiesmann V, Bergler M, Palmisano R, Prinzen M, Franz D, Wittenberg T (2017). Using simulated fluorescence cell micrographs for the evaluation of cell image segmentation algorithms. *BMC Bioinformatics* 18, 1–12.
- Wu P, Yi J, Zhao G, Huang Z, Qiu B, Gao D (2015). Active contour-based cell segmentation during freezing and its application in cryopreservation. *IEEE Trans Biomed Eng* 62, 284–295.