

¹ **IMCell^{XMBD}: A statistical approach for robust cell identification and quantification from imaging mass cytometry images**

⁴ Xu Xiao^{1,2,#}, Naifei Su^{3,#}, Yan Kong⁴, Lei Zhang⁵, Xin Ding⁶, Wenxian Yang^{3,*}, Rongshan
⁵ Yu^{1,2,3,*}

⁶ ¹ Department of Computer Science, School of Informatics, Xiamen University, Xiamen, China

⁷ ² National Institute for Data Science in Health and Medicine, Xiamen University, Xiamen, China

⁸ ³ Aginome Scientific, Xiamen, China

⁹ ⁴ Peking University Cancer Hospital and Institute, Beijing, China

¹⁰ ⁵ School of Life Science, Xiamen University, Xiamen, China

¹¹ ⁶ Xiamen Zhongshan Hospital, Xiamen University, Xiamen, China

¹² * Corresponding author: rsyu@xmu.edu.cn

¹³

¹⁴ **Imaging Mass Cytometry (IMC) has become a useful tool in biomedical research due to its capability to measure over 100 markers simultaneously. Unfortunately, some protein channels can be very noisy, which may significantly affect the phenotyping results without proper processing.** We developed IMCell^{XMBD}¹, a highly effective and generalizable cell identification and quantification method for imaging mass cytometry images. IMCell performs cell-level denoising by subtracting an estimated background noise value from cell protein expressions,

¹XMBD: Xiamen Big Data, a biomedical open software initiative in the National Institute for Data Science in Health and Medicine, Xiamen University, China.

20 identifies positive cells from negative cells by comparing the distribution between segmented
21 cells and decoy cells, and normalize the protein expression levels of the identified positive
22 cells for downstream data analysis. Experimental results demonstrate that our method sig-
23 nificantly improves the reliability of cell phenotyping which is essential for using IMC in
24 biomedical studies.

25 **1 Introduction**

26 Analysis of the heterogeneity of cells is critical to discover the complexity and factuality of the life
27 system. Recently, single-cell sequencing technologies have been increasingly used in the research
28 of developmental physiology and disease ^{1–4}, but the spatial context of individual cells in the tissue
29 is lost due to tissue dissociation in these technologies. On the other hand, traditional immunohisto-
30 chemistry (IHC) and immunofluorescence (IF) preserve spatial context but the number of markers
31 is limited. The development of multiplex IHC/IF (mIHC/mIF) technologies has enabled the detec-
32 tion of multiple markers simultaneously and preserve spatial information, such as cyclic IHC/IF
33 and metal-based multiplex imaging technologies ^{5–8}. Imaging mass cytometry (IMC) ^{6,9}, one of
34 metal-based mIHC technologies, uses a high-resolution laser with a mass cytometer and enables
35 simultaneous measurement of up to 100 markers. Due to its high resolution and large number of
36 concurrent marker channels available, IMC has been proven to be highly effective in identifying
37 the complex cell phenotypes and interactions coupled with spatial locations, and has been utilized
38 in many biomedical and clinical studies on tumor or immune diseases ^{6,10–21}.

39 A number of methodological challenges must be overcome when applying IMC to clinical

40 applications in order to derive reliable cell quantification and phenotyping results from IMC. Im-
41 ages generated by a mass cytometry system are subject to noise and other acquisition artifacts
42 resulting from, e.g., sample protein degradation or signal spill-over between heavy metals ²². In-
43 strument performance can vary within a single sample, not to mention the technical variance among
44 different instruments. Besides, the antibody performance and antigen retrieval condition can dif-
45 fer between samples due to their storage time and environment, which result in protein variations
46 between and within samples. Therefore, specific data processing steps are needed to ensure mea-
47 surement of cellular markers with high resolution, quality, and reproducibility. Quality control and
48 data normalization have been incorporated into the standard operation procedures in the software
49 of the mass cytometers to convert raw signals to images ²³. Most IMC image quality control and
50 preprocessing steps are performed semi-automatically and tuned for individual datasets. Some
51 generic signal processing techniques have been applied to different datasets, including background
52 removal, removing hot pixels, and denoising by low pass filtering, etc ^{11,24}. Data normalization has
53 also been discussed to eliminate the variation between samples ^{25,26}. Despite the progress of IMC
54 data processing tools, in practice it is still possible to obtain IMC images with very poor signal-to-
55 noise ratios that exceed the processing capabilities of existing tools. In such cases, it remains as
56 an intricate issue to identify true positive cells from strong background noise, and harmonize their
57 protein expression levels across slices for downstream analysis .

58 In this paper, we present IMCell, a method of protein quantification for single cell from IMC
59 images. IMCell is able to reliably identify positive cells from highly noisy channels of an IMC
60 image, and perform expression quantification for those cells. To this end, IMCell uses monte carlo

61 method to create decoy cells randomly on the image, and computed the distribution of the protein
62 expression of the decoy cells to derive the background noise level of the image. The positive cells
63 are then identified with false discovery rate (FDR) control by comparing the protein expression
64 distribution of decoy cell with that of the segmented true cells. To reduce the effect of background
65 noise to to the quantification results, IMCell further performs noise reduction on the IMC image
66 the identified background noise level. Finally, the protein expression of the positive cells are
67 normalized to mitigate the variations of pixel values across different IMC images. Our evaluation
68 results show that IMCell can retain real signal with a user-defined confidence level and eliminate
69 sample variations, improves IMC images quality and benefits the downstream analysis.

70 **2 Results**

71 **IMCell identifies true positive cells from noise** IMCell identifies positive cells on each protein
72 channel based on FDR control with a distribution of permuted decoy cells. First, IMCell ran-
73 domly generates a large number of decoy cells on a potential noise region on each protein channel
74 (Methods, Figure 1). With the generated decoy cells, IMCell identifies positive cells by compar-
75 ing the distribution of image intensities (i.e., cell protein expressions at that pixel location) of all
76 segmented cells and decoy cells, from which the detection threshold can be set based on the target
77 FDR (Methods, Figure 1). Once the positive cells are identified on each protein channel, IMCell
78 further estimates the background noise level (Methods), which is then removed from the respective
79 IMC channel to generate a clean image for each channel.

80 We compared the performance for background noise removal of IMCell with two commonly
81 used methods, the percentile method and the median filter. The percentile method defines a lower
82 threshold T_l and an upper threshold T_h . It then removes outliers by setting pixel value to zero for
83 those lower than T_l , and setting pixel values to T_h for those higher than T_h . Here we used 1%
84 as the lower threshold T_l and 99% as the upper threshold T_h . Results show that the percentile
85 method removes outliers but cannot deal with noise of similar intensity values as the signal, such
86 as salt-and-pepper noise. On the other hand, the median filter removes salt-and-pepper noise but
87 does not work with other types of noise. By estimating the background noise level from decoy
88 cells randomly drawn from the noise areas of the image, IMCell successfully removed backgrond
89 noise and improved the signal-to-noise ratio, resulting in a cleaner image with true cells presented
90 (Figure 2a, 2b).

91 A clearer example is shown in Figure 2c, 2d). We compared the co-expression pattern of
92 CD45, CD3 and CD4 from different methods and observed that IMCell can retain true CD3 signal
93 since most CD4 T cells expressed CD3 and CD4. While median filter over-removed CD3 signal
94 and percentile method failed to remove noise in the CD3 channel.

95 **IMCell reduces variations in pixel intensity and cell protein expression across IMC images.**
96 Analysis of the raw images and segmented cells show that the range of pixel intensity and the level
97 of signal to noise ratio vary significantly among samples (Figure 3a). The difference is conspicuous
98 even after performing the variance stabilizing transform, e.g., the inverse sinh transform ²⁷, on the
99 IMC images to reduce the overall range of the pixel intensities (Figure 3b). The distribution plots

100 demonstrate that the variation across samples exists not only at pixel level but also at cell level,
101 if the cell protein expressions were calculated directly from the raw images. Large inter-sample
102 distribution variation could be misleading in downstream data analysis, as the cells may cluster by
103 samples but not by cell types. In IMCell, protein expression levels are normalized across the entire
104 dataset based on the identified positive cells (Methods). Figure 3c shows the variation of intensity
105 across three samples at both pixel and cell levels after intensity normalization by IMCell.

106 **IMCell enables clustering with biological significance** To investigate the effects of different
107 IMC image preprocessing methods on downstream analysis, we applied unsupervised clustering
108 on cells generated from raw IMC images. The clustering was performed using a same subset of
109 proteins as features. After clustering, the cell type of each cluster can be identified based on its
110 marker expression pattern compared to that of known immune and tumor cell types (Figure 4).
111 The cell types of the cell clusters obtained from raw IMC images or processed using the percentile
112 method can hardly be identified. As the heatmap shows, some clusters have more than one rela-
113 tively high cell-type-specific protein expressions (Figure 4a). For example, Cluster 1 from the raw
114 IMC images contains similar protein expression level for both lymphoid (e.g. CD4) and myeloid
115 cells (e.g. CD14, CD68), causing confusion in cell type identification. The percentile method also
116 shows a confusing heatmap where the cell types cannot be ascertained (Figure 4b). Alternatively,
117 by applying the median filter or IMCell on the raw images, the cell clustering results are more bio-
118 logically significant (Figure 4c, 4d). For the clustering results obtained from cells of IMC images
119 preprocessed by the median filter, we can annotate Cluster 12 as B cell, but still have difficulty to
120 determine other two clusters (Cluster 1 and 10) because they contain T cell markers (e.g. CD4,

121 CD8) and a certain amount of myeloid cell markers such as CD68 and CD14. On the other hand,
122 we will able to obtain highly specific cell clusters from clustering results obtained from cells quan-
123 tified with IMCell, e.g., CD4 T cell (Cluster 4), CD8 T cell (Cluster 1), B cell (Cluster 3) and
124 myeloid cell (Cluster 12, 13, 15). **Cluster 3: B cell? not that obvious??**

125 **3 Discussion**

126 In this work, we developed IMCell which enables efficient and more accurate cell quantification
127 from IMC images. Our work is based on the notion of statistical testing by contrasting the distri-
128 butions of both foreground cells (true cells identified by image segmentation software) and decoy
129 cells. As decoy cells are drawn from potential noise-only regions of IMC image with random
130 shapes and locations, it can be anticipated that its distributions will highly resemble those of nega-
131 tive cells (i.e., cells that don't express target proteins). Therefore, the positive cells can be reliably
132 identified with proper FDR control on the distributions of both cells. Note that the successful ap-
133 plication of IMCell depends on the availability of information on true cell segmentations. In this
134 work we used Dice-XMBD²⁸, a deep neural network based IMC cell segmentation tool that is able
135 to perform automatic cell segmentation from IMC images without manual annotation. However, it
136 is also possible to use other segmentation tools, e.g., Ilastik or CelProfiler, to perform such a task.

137 Normalization across different images is critical to align the protein expressions to the same
138 sea-level such that they can be compared in downstreaming data analysis. However, such normal-
139 ization can only be performed if the positive cells (i.e., cells expressing certain target proteins)

140 can be reliably identified. Otherwise, the normalization can falsely amplify negative cells located
141 at noise regions of the image, resulting in severe false positive issues that plague the downstream
142 biological analysis. For this reason, expression normalization is rarely performed in existing IMC
143 processing pipelines although significant inter-slice variations of marker protein expressions can
144 easily happen in IMC studies. In IMCell, by rigorous FDR control, expression normalization is
145 only performed on highly-confident positive cells, thus minimizing the risk of amplification of
146 false-positive cells. As validated by visual inspection and clustering analysis, cell quantification
147 by IMCell leads to much more consistent connections between cell phenotypes and marker protein
148 expressions.

149 IMCell is freely available as an open-source software at <https://github.com/xmuyulab/imcell>. We anticipate that IMCell could help to promote the better usage of IMC both in research
150 labs and in clinical settings.

152 4 Methods

153 **Patients and IMC data acquisition** Melanoma cancer formalin-fixed paraffin-embedded (FFPE)
154 tissues were stained with a customized panel (35 antibodies) to generate the IMC images used in
155 this study. We excluded images containing large areas with nonspecific background staining that
156 could be caused by nonspecific antibody binding²⁹ by manual inspection using the MCD viewer
157 (V1.0.560.6). The remaining 158 images were further analyzed in the following procedures.

¹⁵⁸ **Overview of the IMCell workflow.** IMCell consists of two main modules, i.e., denoising and nor-
¹⁵⁹ malization (Figure 5). Firstly, raw IMC images are preprocessed and segmented by any cell seg-
¹⁶⁰ mentation methods. Then we randomly generate a number of decoy cells on the high-confidence
¹⁶¹ noise region of each protein channel image. The protein expressions of the decoy cells are used
¹⁶² to estimate the background noise of the protein image. After that the protein expression distri-
¹⁶³ butions of all segmented cells and decoy cells are compared to identified positive cells with FDR
¹⁶⁴ control. Next in the normalization part, to fairly compare positive cells across images, we scale the
¹⁶⁵ mean expression of positive cells from each image to the same level. More details are described as
¹⁶⁶ following step by step.

¹⁶⁷ **Cell segmentation using Dice-XMBD** Single cells were identified by Dice-XMBD ²⁸ using a
¹⁶⁸ pretrained deep-learning based model, and referred to as segmented cells in this paper. Note that
¹⁶⁹ other cell segmentation methods can also be used in the IMCell pipeline, for example, interactive
¹⁷⁰ cell segmentation by using Ilastik and CellProfiler. For quality control, the segmented cells that
¹⁷¹ cover less than 5 pixels are discarded. The cell protein expressions are extracted as the mean of the
¹⁷² pixel values in each cell mask region.

¹⁷³ **Preprocessing and hot pixel removal** We first applied the hyperbolic inverse sine function (arc-
¹⁷⁴ sinh) on all the pixel values for each channel. The raw marker intensities output from cytometers
¹⁷⁵ tend to have strongly skewed distributions with varying ranges of expression. It is thus a com-
¹⁷⁶ mon practice to transform the raw marker intensities using arcsinh to make the distributions more
¹⁷⁷ symmetric and to map them to a comparable range of expression ^{27,30}.

178 Hot pixels were removed by filtering with a 5×5 pixel² window. If the center pixel of the
179 window was in the top 2% of all pixel values in the channel and was at least $4\times$ above the median
180 value of all pixels in the window, it will be identified as a hot pixel and its value will be replaced by
181 the median value in the window. This step reduces the scattered hot pixels' noise on quantification
182 of protein expression values for the cells.

183 **Generating decoy cells** We established the distribution of noise for each channel by generating
184 a large number (N) of decoy cells using monte carlo method. To this end, we first identified
185 regions on the image that potentially contain noise-only signals without real protein expression by
186 excluding pixels with values above $0.05 \times Q_{99}$, where Q_{99} is the 99th percentile (Q_{99}) of the pixel
187 values. After that, we set the value of remaining pixels to zero and smooth the noise regions by
188 applying a 5×5 median filter on the image.

189 We then fit each segmented cell as an ellipse. For each image, the mean and variance of the
190 major axis, the minor axis, and the orientation angle of all the segmented cells were calculated, and
191 these three parameters were fit using individual Gaussian models. Random parameters are drawn
192 from the distributions of the major axis, the minor axis, and the orientation, respectively, to form
193 an ellipse as a decoy cell. The decoy cell was randomly placed in the noise region of the channel
194 image, such that the center of the decoy cell was at least 5 pixels away from image boundaries.
195 The decoy cell should only lies in noise regions, i.e., all of its pixels lie in noise regions as in the
196 noise region mask. When the decoy cell lies on the border of the image, it must cover more than 5
197 pixels in the image, otherwise it will be discarded. We further filter out the decoy cell if the area

198 it covered exceeded the size range of all segmented cells. Then, the protein expression value for
199 each decoy cell was calculated as the mean of its pixel values in the preprocessed IMC image.

200 **Background noise removal** To eliminate the effect of different background noise profiles and
201 levels between different proteins in an IMC dataset, we removed background noise using the decoy
202 cells generated from the noise regions. For each channel image, the mean of protein expression
203 values of all generated decoy cells was calculated, and subtracted from each pixel value to remove
204 channel-specific background noise.

Positive cells identification by FDR control The segmented cells may include both positive cells and negative cells. We used a permutation test to compare the protein expression distributions between segmented cells and randomly drawn decoy cells from the noise regions, and use FDR control to identify positive cells. The FDR value can be adjusted to obtain positive cells with acceptable error-tolerant rate. The FDR of true cell identification is calculated by

$$FDR = \frac{FP}{FP + TP}, \quad (1)$$

205 where TP and FP refers to true positive and false positive, respectively. More specifically, TP refers
206 to the number of segmented cells with protein expression values larger than the threshold, while
207 FP refers to the number of decoy cells with protein expression values larger than the threshold.
208 The default value of FDR was set to 0.01, and the threshold for positive cell identification can be
209 then determined to satisfy the FDR level.

210 **Normalization of cell protein expressions** The data processing steps above are all performed on
211 individual channel images. As the antibody performance and the signal-to-noise ratio can differ
212 considerably between FFPE tissues due to variations in tissue processing, we further normalized
213 the cell protein expression values across different samples within one IMC dataset for each protein
214 separately. Denote the channel image of protein p for sample i as $I_i^{(p)}$ and the mean of the protein
215 expression values for all identified positive cells as $\mu_i^{(p)}$. Let $m^{(p)}$ denote the maximum protein
216 expression value among all identified positive cells for protein p in all samples. The cell protein
217 expression values for sample i were then scaled by factor $\frac{m^{(p)}}{\mu_i^{(p)}}$.

218 **Single cell clustering and phenotyping** High-dimensional single cell protein expression data
219 were clipped at the 99th percentile followed by min-max normalization. We selected 20 mark-
220 ers to perform cell clustering: CD45, CD3, CD4, CD8a, FoxP3, CD20, CD68, CD14, CD16,
221 CD11c, CD11b, IDO, Vimentin, α -SMA, E-cadherin, EpCAM, CA9, VEGF, PDGFRb and Col-
222 lagen. The clustering analysis consists of two consecutive steps: first a self-organizing map (50 x
223 50 nodes) implemented in FlowSOM (R package, v1.18.0) was used to generated several groups,
224 then a community detection algorithm using Phenograph (R package, v0.99.1) was used on the
225 mean expression values of each group from FlowSOM clustering. Cell phenotyping was deter-
226 mined by the mean of each cluster protein expression compared with the known cell types' protein
227 expression patterns.

228 **Conflict of Interest Statement**

229 RY and WY are shareholders of Aginome Scientific. The authors declare no other conflict of
230 interest.

231 **Author Contributions**

232 WY and RY discussed the ideas and supervised the study. NS implemented the denoising method
233 and XX conducted experiments in evaluation and biological analysis. All authors wrote and dis-
234 cussed on the manuscript.

235 **Data Availability**

236 **Code Availability**

237 **References**

- 238 1. Lähnemann, D. *et al.* Eleven grand challenges in single-cell data science. *Genome Biology*
240 **21**, 1–35 (2020).
- 241 2. Stubbington, M. J., Rozenblatt-Rosen, O., Regev, A. & Teichmann, S. A. Single-cell tran-
242 scriptomics to explore the immune system in health and disease. *Science* **358**, 58–63 (2017).
- 243 3. Potter, S. S. Single-cell rna sequencing for the study of development, physiology and disease.
244 *Nature Reviews Nephrology* **14**, 479–492 (2018).

- 245 4. Papalex, E. & Satija, R. Single-cell rna sequencing to explore immune cell heterogeneity.
- 246 *Nature Reviews Immunology* **18**, 35 (2018).
- 247 5. Tan, W. C. C. *et al.* Overview of multiplex immunohistochemistry/immunofluorescence tech-
- 248 niques in the era of cancer immunotherapy. *Cancer Communications* **40**, 135–153 (2020).
- 249 6. Giesen, C. *et al.* Highly multiplexed imaging of tumor tissues with subcellular resolution by
- 250 mass cytometry. *Nature Methods* **11**, 417–422 (2014).
- 251 7. Zrazhevskiy, P. & Gao, X. Quantum dot imaging platform for single-cell molecular profiling.
- 252 *Nature Communications* **4**, 1–12 (2013).
- 253 8. Angelo, M. *et al.* Multiplexed ion beam imaging of human breast tumors. *Nature Medicine*
- 254 **20**, 436–442 (2014).
- 255 9. Chang, Q. *et al.* Imaging mass cytometry. *Cytometry Part A* **91**, 160–169 (2017).
- 256 10. Damond, N. *et al.* A map of human type 1 diabetes progression by imaging mass cytometry.
- 257 *Cell Metabolism* **29**, 755–768 (2019).
- 258 11. Wang, Y. J. *et al.* Multiplexed in situ imaging mass cytometry analysis of the human endocrine
- 259 pancreas and immune system in type 1 diabetes. *Cell Metabolism* **29**, 769–783 (2019).
- 260 12. Ramaglia, V. *et al.* Multiplexed imaging of immune cells in staged multiple sclerosis lesions
- 261 by mass cytometry. *Elife* **8**, e48051 (2019).

- 262 13. Böttcher, C. *et al.* Single-cell mass cytometry reveals complex myeloid cell composition in
263 active lesions of progressive multiple sclerosis. *Acta Neuropathologica Communications* **8**,
264 1–18 (2020).
- 265 14. de Vries, N. L., Mahfouz, A., Koning, F. & de Miranda, N. F. Unraveling the complexity
266 of the cancer microenvironment with multidimensional genomic and cytometric technologies.
267 *Frontiers in Oncology* **10**, 1254 (2020).
- 268 15. Brähler, S. *et al.* Opposing roles of dendritic cell subsets in experimental gn. *Journal of the*
269 *American Society of Nephrology* **29**, 138–154 (2018).
- 270 16. Aoki, T. *et al.* Single-cell transcriptome analysis reveals disease-defining t-cell subsets in
271 the tumor microenvironment of classic hodgkin lymphoma. *Cancer Discovery* **10**, 406–421
272 (2020).
- 273 17. Jackson, H. W. *et al.* The single-cell pathology landscape of breast cancer. *Nature* **578**, 615–
274 620 (2020).
- 275 18. Ali, H. R. *et al.* Imaging mass cytometry and multiplatform genomics define the
276 phenogenomic landscape of breast cancer. *Nature Cancer* **1**, 163–175 (2020).
- 277 19. Dey, P. *et al.* Oncogenic kras-driven metabolic reprogramming in pancreatic cancer cells
278 utilizes cytokines from the tumor microenvironment. *Cancer Discovery* **10**, 608–625 (2020).
- 279 20. Zhang, Y., Gao, Y., Qiao, L., Wang, W. & Chen, D. Inflammatory response cells during acute
280 respiratory distress syndrome in patients with coronavirus disease 2019 (covid-19). *Annals of*
281 *Internal Medicine* (2020).

- 282 21. Schwabenland, M. *et al.* Deep spatial profiling of human covid-19 brains reveals neuroin-
- 283flammation with distinct microanatomical microglia-t-cell interactions. *Immunity* **54**, 1594–
- 284 1610.e11 (2021).
- 285 22. Chevrier, S. *et al.* Compensation of signal spillover in suspension and imaging mass cytometry.
- 286 *Cell Systems* **6**, 612–620 (2018).
- 287 23. Lee, B. H. & Rahman, A. H. Acquisition, processing, and quality control of mass cytometry
- 288 data. In *Mass Cytometry*, 13–31 (Springer, 2019).
- 289 24. Baharlou, H., Canete, N. P., Cunningham, A. L., Harman, A. N. & Patrick, E. Mass cytometry
- 290 imaging for the study of human diseases—applications and data analysis strategies. *Frontiers*
- 291 in *Immunology* **10**, 2657 (2019).
- 292 25. Ijsselsteijn, M. E., Somarakis, A., Lelieveldt, B. P., Hollt, T. & de Miranda, N. F. Semi-
- 293 automated background removal limits loss of data and normalises the images for downstream
- 294 analysis of imaging mass cytometry data. *bioRxiv* (2020).
- 295 26. Keren, L. *et al.* A structured tumor-immune microenvironment in triple negative breast cancer
- 296 revealed by multiplexed ion beam imaging. *Cell* **174**, 1373–1387 (2018).
- 297 27. Bendall, S. C. *et al.* Single-cell mass cytometry of differential immune and drug responses
- 298 across a human hematopoietic continuum. *Science* **332**, 687–696 (2011).
- 299 28. Xiao, X. *et al.* Dice-XMBD: Deep learning-based cell segmentation for imaging mass cytom-
- 300 etry. *Frontiers in Genetics* **12**, 1532 (2021).

- 301 29. Buchwalow, I., Samoilova, V., Boecker, W. & Tiemann, M. Non-specific binding of antibodies
 302 in immunohistochemistry: fallacies and facts. *Scientific Reports* **1**, 1–6 (2011).
- 303 30. Bruggner, R. V., Bodenmiller, B., Dill, D. L., Tibshirani, R. J. & Nolan, G. P. Automated
 304 identification of stratifying signatures in cellular subpopulations. *Proceedings of the National
 305 Academy of Sciences* **111**, E2770–E2777 (2014).

306 **Figure captions**

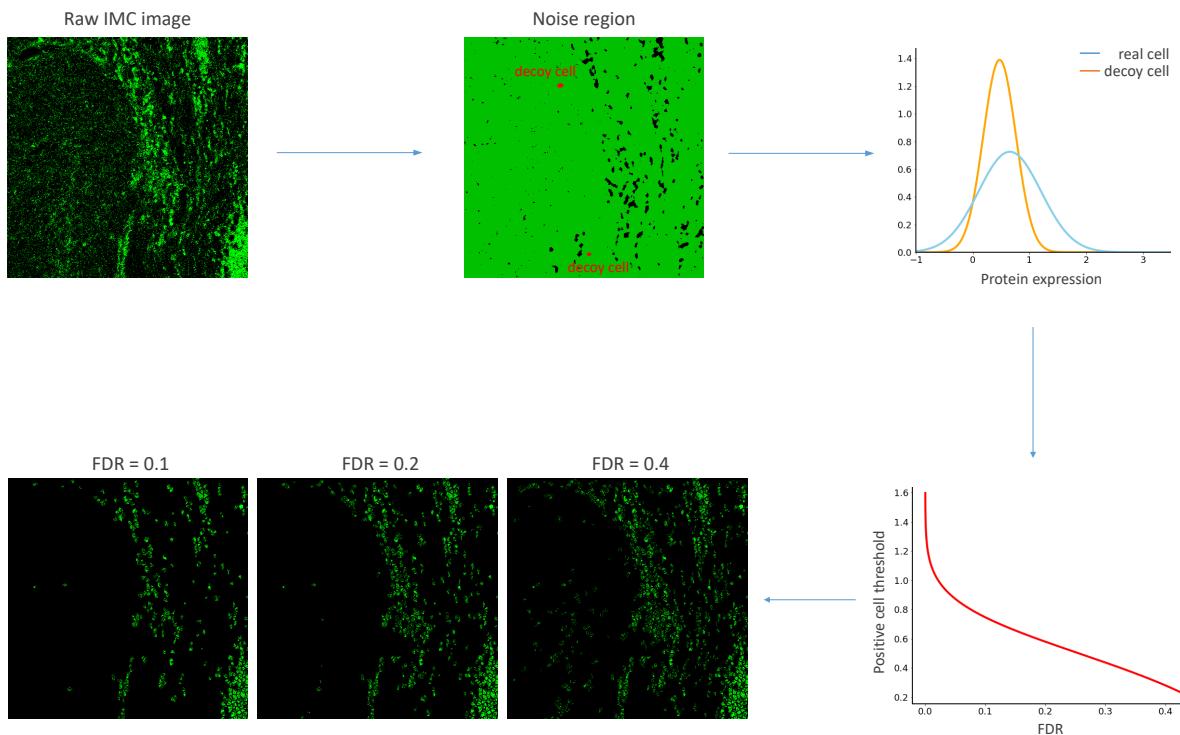


Figure 1: Positive cells identified by IMCell with different FDR control (sample: 76 ROI18, protein: CD74).

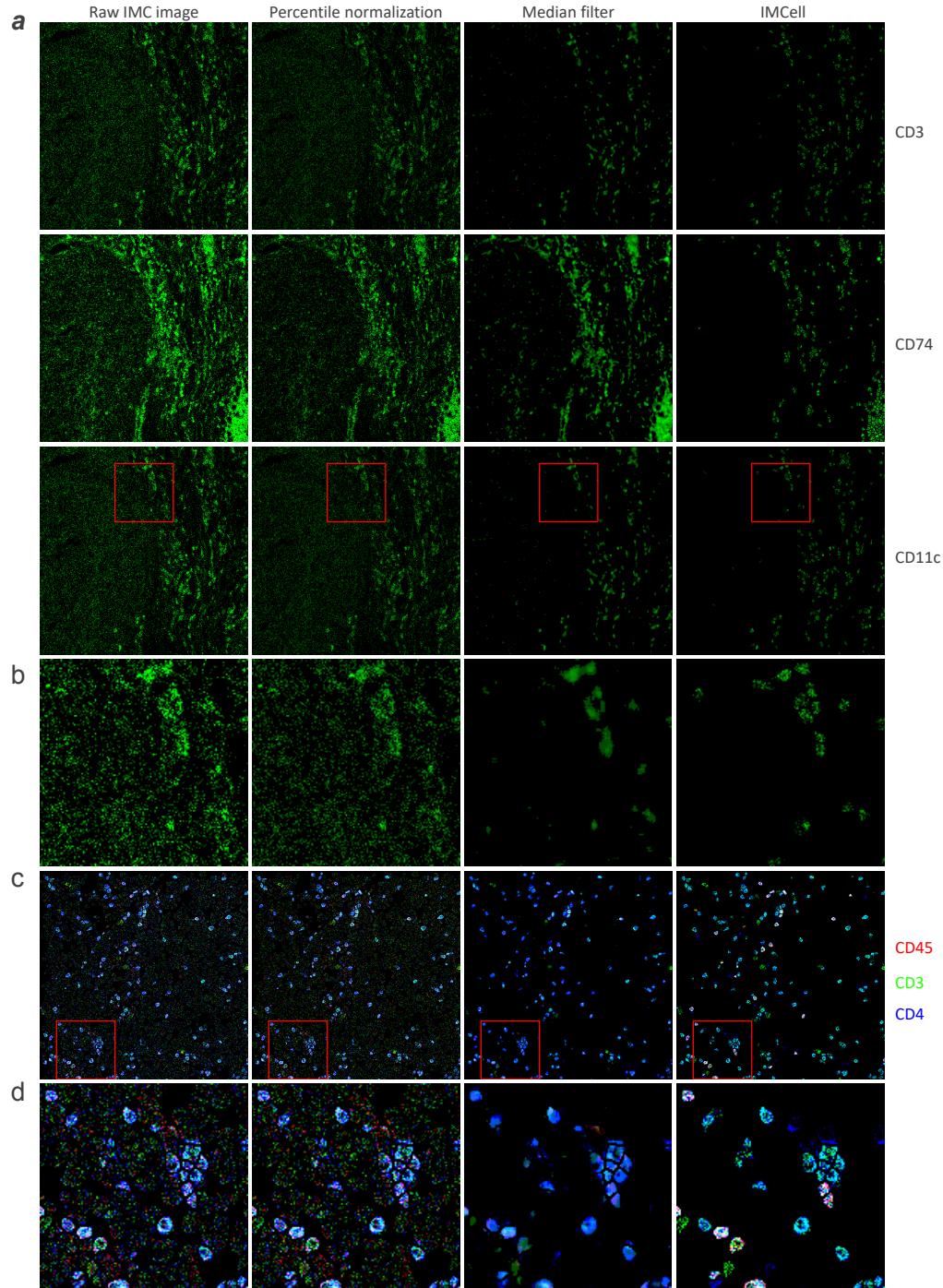


Figure 2: Performance evaluation of different methods. (a) Comparisons of cell identification results from raw IMC images, with 1st-99th percentile method to remove outliers, with median filter to remove salt-and-pepper noise, and with IMCell (sample: 76 ROI18). The red box marks the zoomed in areas on the below side (b) depicting the CD11c marker. Expression pattern of multi-markers (CD45, CD3, and CD4) in the whole images (c) and zoom-in areas (d).

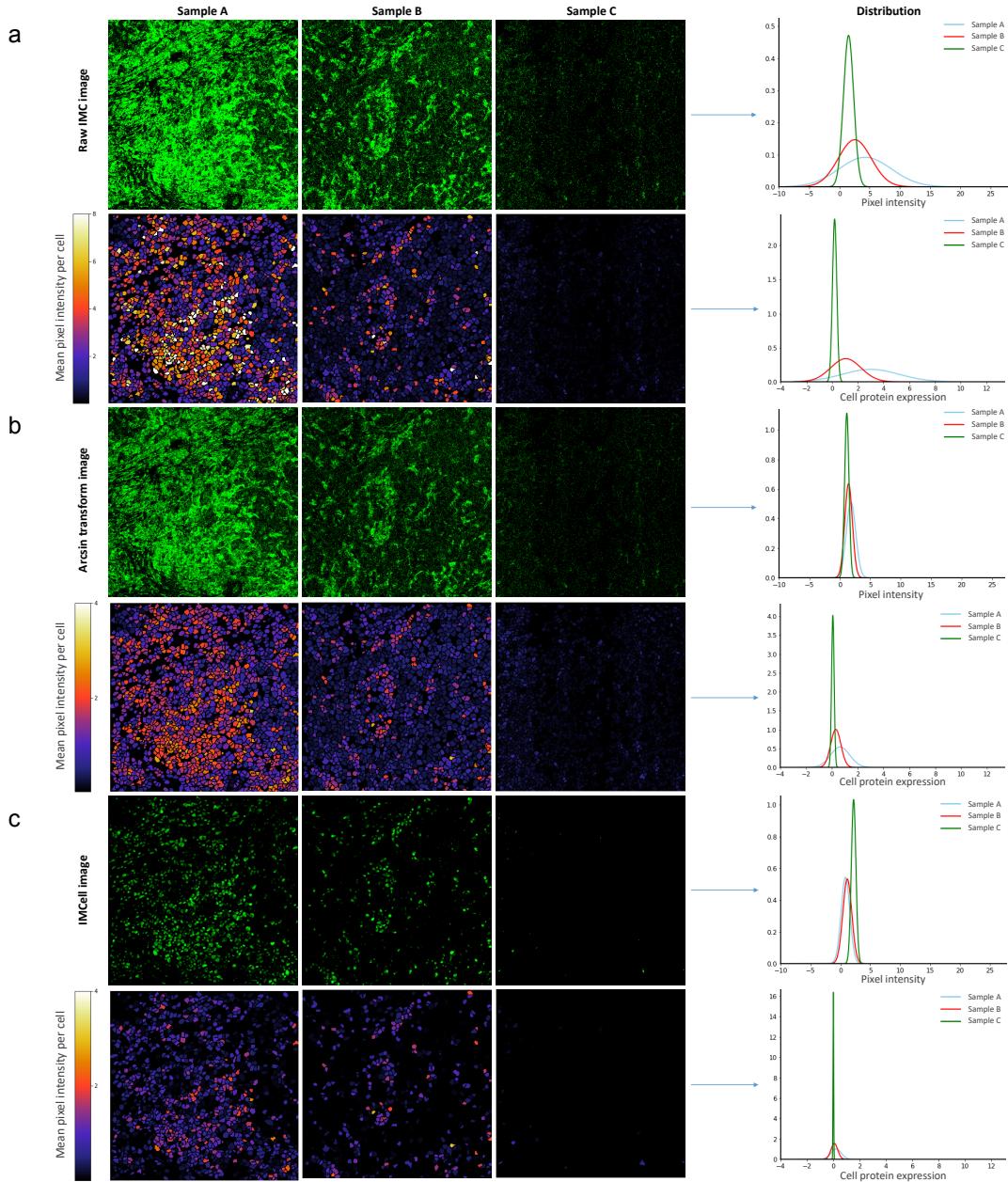


Figure 3: Variation of pixel intensity and cell protein expression across three samples (sample A: 65 ROI13, sample B: 65 ROI18, sample C: 33 ROI11). The left column shows (a) the pixel intensity (first row) and cell protein expression (second row) from the raw images, (b) the pixel intensity (first row) and cell protein expression (second row) from arcsinh-transformed images, and (c) the pixel intensity (first row) and cell protein expression (second row) from images processed by IMCell. The right column plots the distribution of the corresponding value (i.e., pixel intensity and cell protein expression).

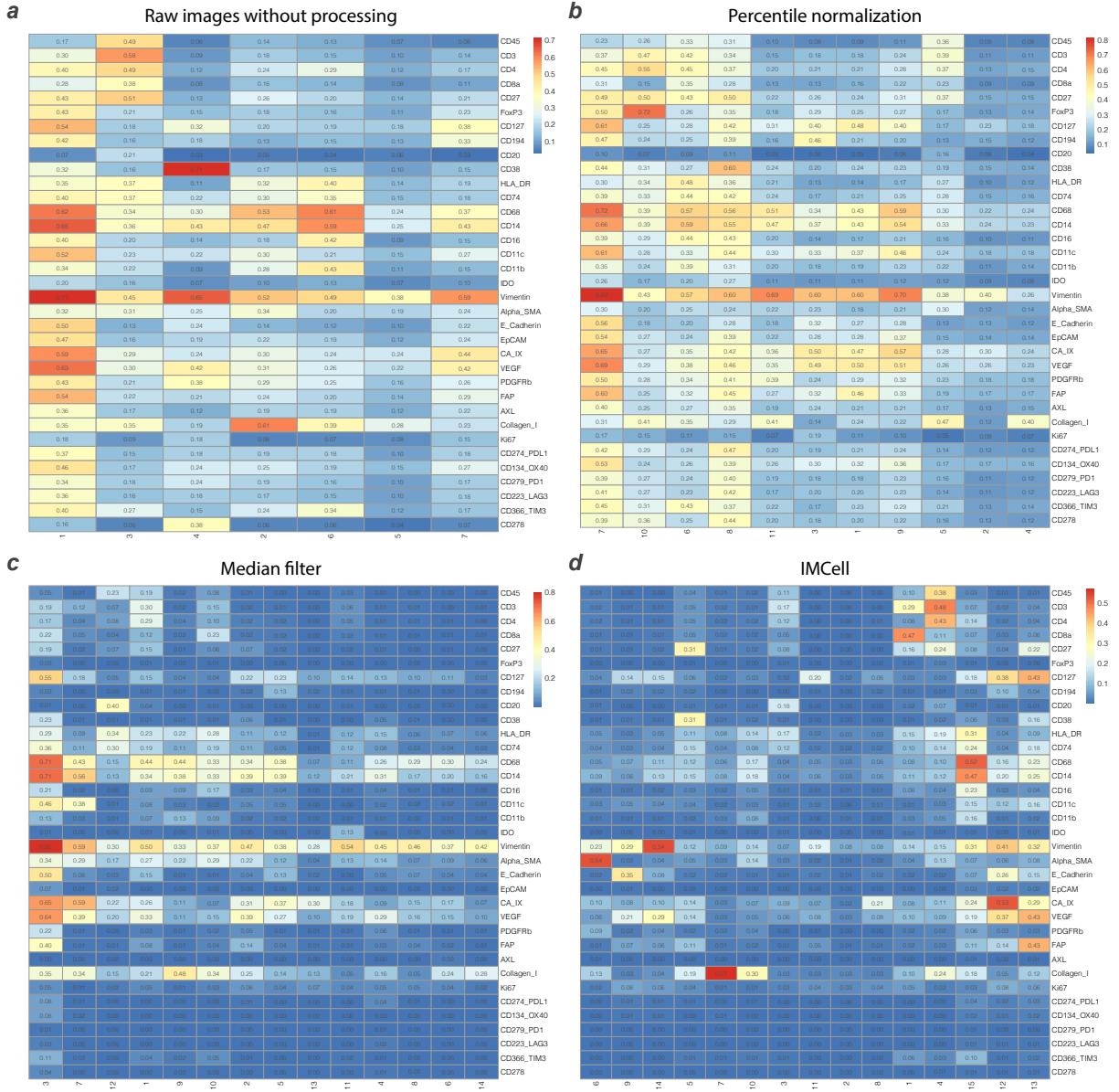
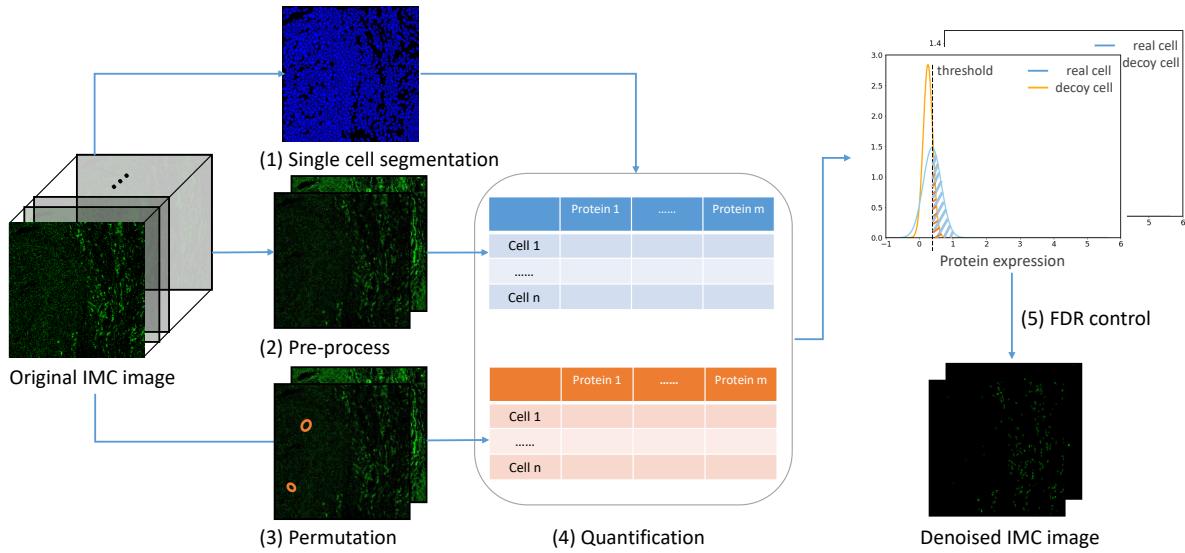


Figure 4: Clustering result from different methods. Heatmap showing mean value of normalized protein expression in each cluster. The high-dimensional single cell expression data were generated from (a) raw IMC images, (b) with 1st-99th percentile method to remove outliers, (c) with median filter to remove salt-and-pepper noise, and (d) with IMCell.

(a) Denoise



(b) Normalization

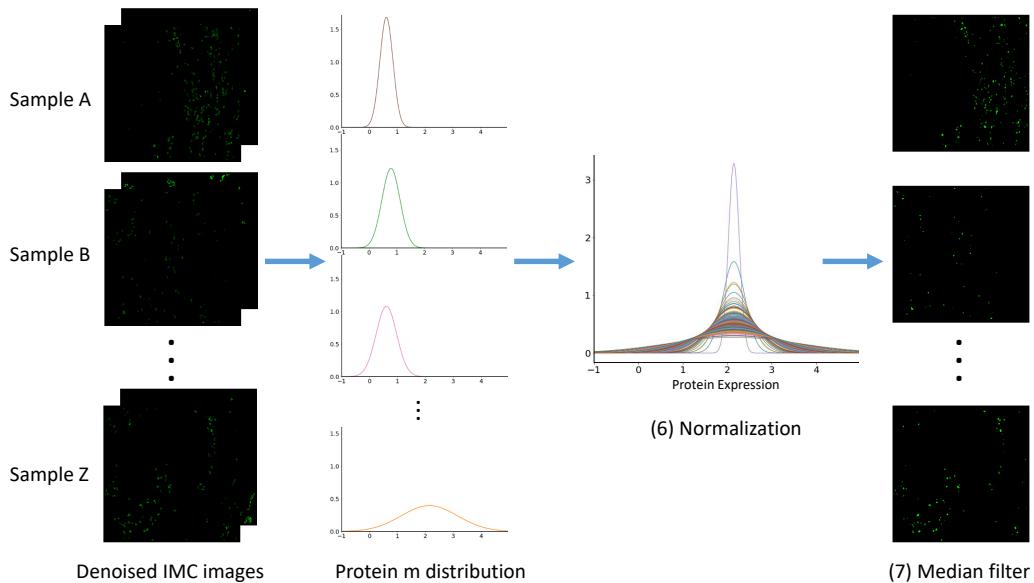


Figure 5: The workflow of IMCell consists of (a) denoising and (b) normalization. The workflow includes the following procedures, (1) single cell segmentation by Dice-XMBD, (2) image pre-processing and hot pixel removal, (3) random generation of decoy cells in high-confidence noise regions, (4) protein quantification for segmented cells and decoy cells, (5) identifying positive cells with FDR control, (6) normalization by scaling using the mean of protein expression of positive cells, and (7) apply the median filter on the denoised and normalized images.