

1   **STHD: probabilistic cell typing of single Spots in whole Transcriptome spatial  
2   data with High Definition**

3   Chuhanwen Sun<sup>1\*</sup>, Yi Zhang<sup>1-4\*#</sup>

4   1. Department of Neurosurgery, Duke University

5   2. Department of Biostatistics and Bioinformatics, Duke University

6   3. Brain Tumor Omics Program, The Preston Robert Tisch Brain Tumor Center, Duke University

7   4. Duke Cancer Institute

8   \*These authors contributed equally.

9   #Corresponding author: Yi Zhang, Ph.D. [yi.zhang@duke.edu](mailto:yi.zhang@duke.edu)

10   **Abstract**

11   Recent spatial transcriptomics (ST) technologies have enabled sub-single-cell resolution  
12   profiling of gene expression across the whole transcriptome. However, the transition to high-  
13   definition ST significantly increased sparsity and dimensionality, posing computational  
14   challenges in discerning cell identities, understanding neighborhood structure, and identifying  
15   differential expression - all are crucial steps to study normal and disease ST samples. Here we  
16   present STHD, a novel machine learning method for probabilistic cell typing of single spots in  
17   whole-transcriptome, high-resolution ST data. Unlike current binning-aggregation-deconvolution  
18   strategy, STHD directly models gene expression at single-spot level to infer cell type identities.  
19   It addresses sparsity by modeling count statistics, incorporating neighbor similarities, and  
20   leveraging reference single-cell RNA-seq data. We demonstrated that STHD accurately predicts  
21   cell type identities at single-spot level, which automatically achieved precise segmentation of  
22   global tissue architecture and local multicellular neighborhoods. The STHD labels facilitated  
23   various downstream analyses, including cell type-stratified bin aggregation, spatial  
24   compositional comparison, and cell type-specific differential expression analyses. These high-  
25   resolution labels further defined frontlines of inter-cell type interactions, revealing direct cell-cell  
26   communication activities at immune hubs of a colon cancer sample. Overall, computational

27 modeling of high-resolution spots with STHD uncovers precise spatial organization and deeper  
28 biological insights for disease mechanisms.

29

30 **Main**

31 Spatial transcriptomics (ST) technologies have enabled gene expression profiling in the original  
32 spatial context of tissues. It provides a systematic approach to profile cell type distribution,  
33 regional variations, and cell-cell communications in normal and disease samples. Current ST  
34 technologies are based on next-generation sequencing or in situ hybridization imaging, in short,  
35 seq-ST and image-ST<sup>1</sup>. Often, seq-ST approaches can cover whole-transcriptome with  
36 compromised resolution or depth (e.g. 10X Visium, Slide-seq<sup>2</sup>, Stereo-seq<sup>3</sup>), whereas image-  
37 based approaches can provide sub-cellular resolution for genes in a designed panel (e.g. 10X  
38 Xenium, STARmap<sup>4</sup>, MERFISH<sup>5</sup>).

39

40 Recent ST technologies have rapidly advanced with image-ST increasing gene coverage and  
41 seq-ST enhancing resolution to subcellular scale. For example, VisiumHD from 10X Genomics  
42 recently become available, capable of whole transcriptome profiling in spots of size 2x2um,  
43 which is sub-cellular for many cell types and tissue contexts. However, analyses of this data  
44 type present significant computational challenges. The read coverage in each HD spot is often  
45 highly sparse and usually requires binning of adjacent spots at a chosen size (e.g. aggregating  
46 4x4 spots into 8x8um bins, single-cell size). The binning process then computationally merges  
47 spots for sufficient read depth in downstream analytical tasks<sup>6</sup>. However, binning compromises  
48 resolution and mixes cell types, which further requires cell type deconvolution as in low-  
49 resolution seq-ST (using e.g. RCTD<sup>7</sup>, CARD<sup>8</sup>, Cell2Location<sup>9</sup>). Additionally, the vast number of  
50 spots or bins in high-resolution ST data poses challenges in computational efficiency.

51

52 We here develop STHD (probabilistic cell typing of single Spots in whole Transcriptome spatial  
53 data with High Definition), a machine learning method for cell typing of individual high-resolution  
54 spots in whole-transcriptome spatial data. The STHD model leverages cell type-specific gene  
55 expression from reference single-cell RNA-seq data, constructs a statistical model on spot gene  
56 counts, and employs regularization from neighbor similarity. The STHD model outputs cell type  
57 labels and probabilities for each single spot, which addresses the compromised resolution  
58 resulted from the current binning approach. Our results demonstrate that STHD labels can  
59 further stratify cell types in bin aggregation, which facilitates cell type-specific analyses like  
60 compositional comparison and differential expression analyses. Additionally, STHD identifies  
61 frontlines of inter-cell type interactions to enable direct cell-cell communication analysis.

62

63 An overview of STHD model is illustrated in **Fig.1a**. STHD infers latent cell type identities of  
64 each spot using a loss function that simultaneously optimizes two components: likelihood of  
65 spot gene counts and similarity with neighbor spots (**Fig.1a, Methods**). The first part optimizes  
66 log likelihood of spot gene counts following Poisson distribution, where parameters are based  
67 on normalized gene expression derived from a reference scRNA-seq dataset with cell type  
68 annotation from a matched disease or system. The second part denoises sparsity using cross  
69 entropy loss based on neighborhood similarity in cell type probabilities. The contribution from  
70 neighbors is controlled by tuning the neighbor parameter  $\beta$ . For each high-resolution spot,  
71 STHD outputs cell type probabilities and labels based on Maximum a Posterior (MAP). The spot  
72 cell types enable multiple downstream analytical tasks, including cell type-stratified bin  
73 aggregation, cell type-specific compositional and differential gene expression analyses, and  
74 cell-cell interaction analyses. STHD implements fast optimization enabled by efficient gradient  
75 descent<sup>10,11</sup>. STHD also includes local low-count detection to filter cell and tissue gaps  
76 (**Methods, Supplementary Fig.1**). To be compatible with the vast number of spots in one  
77 sample, we parallelized patch-level inference followed by whole-sample integration

78 (**Supplementary Fig.2**). Finally, to enable close examination of cellular neighborhoods, we  
79 developed STHDviewer for interactive and scalable visualization of spots in the entire sample  
80 (**Fig.2**).  
81

82 Different from current standards that start from binning, we propose that analyses of whole-  
83 transcriptome, high-resolution ST data can begin with modeling the subcellular HD spots as  
84 STHD. This strategy offers advantages compared to alternative approaches. First, direct  
85 unsupervised clustering of spot-level gene expression results in false grouping of spots due to  
86 extreme sparsity (**Fig.1b**, left). Second, the current bin-aggregation approach can increase  
87 depth, but it mixes cell types and compromise resolution (**Fig.1b**, middle). Lastly, segmenting  
88 cell areas based on the hematoxylin and eosin (H&E) histological images may provide precise  
89 cell boundaries, but whole-cell segmentation has remained a challenging task<sup>12</sup> (**Fig.1b**, right).  
90

91 We demonstrated STHD on the recent available VisiumHD sample from 10X Genomics, where  
92 a 6.5x6.5mm of human colorectal tumor tissue was profiled by 8,726,600 HD spots covering  
93 18,085 genes. We constructed the normalized gene expression reference from an independent  
94 human colon cancer atlas study that profiled 43,113 genes from 370,115 single cells annotated  
95 into 98 fine-grained cell types from GSE17834<sup>13</sup> (**Supplementary Fig.3**). The reference gene  
96 profile construction step of STHD selected 4618 cell type-informative genes (**Methods**,  
97 **Supplementary Fig.3**). **Fig.1c** illustrated STHD label for a patch in the colorectal cancer  
98 sample containing 22,336 spots in a 300x300um area or 1100x1100 pixels in the full-resolution  
99 H&E image. First, we observed that STHD-predicted cell type labels collectively captured the  
100 compartmentalization of colon crypt-like structure with tumor-like epithelial cell type signatures  
101 (**Fig.1c**). In detail, stem-like proliferating cells were positioned inside, circled by the abundant  
102 enterocyte cells with scattered tumor stem and transit amplifying (TA)-like cells<sup>13,14</sup>. The  
103 intercrypt regions were composed of abundant endothelial and fibroblast cells, with diverse

104 immune cell populations such as plasma cells, monocytes, macrophages, dendritic cells, and T  
105 cells (**Fig.1c**). Interestingly, the cell type probabilities displayed high confidence for abundant  
106 epithelial spots as well as for rarer cells like macrophages and plasma IgG B cells (**Fig.1d**,  
107 **Extended Fig.1**). By default, spots with maximum posterior probability below 0.8 were assigned  
108 as ambiguous, which mainly located at boundaries of different cell types and tissue structures,  
109 which could overlap cell mixtures or cell gaps (**Fig.1c, d**). In comparison, firstly, clustering of  
110 spot-level gene expression vaguely reflected colon crypt structure with mingled cell identities  
111 (**Fig.1e, Supplementary Fig.4a**). Secondly, clustering of gene expression of bins aggregated  
112 from 4x4 spots (8x8um bins) within patch resulted in less noise in crypt center and intercrypt  
113 regions (**Fig.1f left, Supplementary Fig.4b**). Similar bin clustering for the entire sample could  
114 clearly group the stemness center, enterocytes, and intercrypt regions, but the boundaries are  
115 assigned as one separated cluster that mixed intercrypt immune cells (**Fig.1f right**). Further  
116 deconvolution on the aggregated bins using RCTD separated epithelial structure and immune  
117 cells but lost spot-level resolution (**Fig.1g**). Finally, the H&E image contains densely organized  
118 cells where sensitive whole-cell segmentation has been challenging (**Fig.1h**). In contrast,  
119 STHD-guided binning aggregated spots of the same cell type and separated cell types from the  
120 same bin areas; thus, tumor epithelial cells and intercrypt endothelial, fibroblasts and immune  
121 cells were clearly separated on the STHD-guided binning spatial maps (**Fig.1i, Supplementary**  
122 **Fig.4c, d**).

123

124 Since high-resolution ST data lacks ground truth of spot labels, we evaluated performance from  
125 aspects of histopathology, comparison to other methods, simulation, and marker gene  
126 expression. First, the spatial distribution of STHD cell type labels aligns with cell organization  
127 morphology based on H&E staining (**Fig.1h**). Second, we compared to cell type deconvolution  
128 on 8x8um bins by RCTD doublet mode or full mode using the same reference human colon  
129 cancer scRNA-seq reference (**Methods**). RCTD also resolved cell type proportions that match

130 the colon crypt organization and the distribution of STHD labels for enterocytes and Stem/TA-  
131 like proliferating cells (**Fig.1g**, **Extended Fig.2a,b**). At the raw 2x2um spot level, we next tested  
132 computational methods that were able to label cell types by reference-based deconvolution like  
133 RCTD<sup>7</sup> and CARD<sup>8</sup>, Cell2Location<sup>9</sup>, Celloscope<sup>15</sup>, Stereoscope<sup>16</sup>, or by neighborhood-  
134 augmented clustering like Banksy<sup>17</sup>. We noted that the common gene count filters in these tools  
135 need to be removed to avoid filtering of most spots; for example, 20436 of the 22336 spots are  
136 removed by RCTD default UMI cutoff 100 (**Extended Fig.2c**). At single-spot level, RCTD,  
137 CARD, Cell2Location, and Stereoscope can vaguely reveal the crypt structure while with highly  
138 mixed cell types from different lineages (**Extended Fig.3a-e**). The unsupervised Banksy  
139 efficiently separated epithelial crypt and intercrypt endothelial structures, but mixed the rarer  
140 immune cells inside the endothelial cluster and called the boundary as a separate cluster  
141 (**Extended Fig.3f**). Moreover, STHD has highest computational efficiency among all methods  
142 that enables whole-sample inference (**Extended Fig.4**).  
143

144 To quantitatively assess performance of STHD, we simulated a comprehensive *in silico* whole-  
145 transcriptome, high-resolution ST data (**Methods**, **Fig.1j**, **Supplementary Fig.5**). As shown in  
146 **Fig.1j**, we simulated 500 cells with heterogeneous cell sizes and cell types on 22,500 spots,  
147 each with a true cell type label and gene counts following expression profile from the same  
148 human colon cancer reference. STHD achieved average area under the curve (AUC) 0.9977  
149 with overall accuracy 93.97% for this multi-class classification task with 98 categories (**Fig.1j**,  
150 **Supplementary Fig.5**). Next, we systematically investigated cell type-specific expression of  
151 marker genes based on STHD prediction across the entire colon cancer sample. Despite spot  
152 sparsity (**Supplementary Fig.4e**, **Supplementary Fig.3d**), STHD labels displayed similar cell  
153 type-specific expression patterns of marker genes as in scRNA-seq reference data in normal  
154 epithelial cells (**Fig.1k**, left two panels, **Supplementary Fig.7a**) and all other cell lineages:  
155 tumor epithelial cells (**Supplementary Fig.7b**), myeloid cells (**Supplementary Fig.8a**), B and

156 plasma cells (**Supplementary Fig.8b**), CD8T, CD4T, and other T cells (**Supplementary Fig.9**,  
157 **Supplementary Fig.10a**), fibroblasts (**Supplementary Fig.10b**), and endothelial cells and  
158 pericytes (**Supplementary Fig.11**). STHD-guided binning effectively increased marker gene  
159 coverage and meanwhile retained cell-type specificity, shown in two example bin sizes: 8x8um  
160 bins (**Fig. 1k**, third panel, **Supplementary Figs.7-11**) and 16x16um bins (**Fig.1k**, fourth panel,  
161 **Supplementary Figs.7-11**).

162  
163 STHD enabled scalable inference across the entire ST sample through parallelized inference of  
164 spatial patches followed by whole-sample integration (**Supplementary Fig.2**). To facilitate close  
165 examination of spot identities and scalable investigation of cellular neighborhoods, we  
166 developed STHDviewer, an interactive and scalable visualization toolkit capable of rendering  
167 millions of spots in a webpage (**Fig.2a**, **Extended Fig.5**). As demonstrated in **Fig.2a**, STHD  
168 prediction automatically achieved precise segmentation of global tissue architecture and local  
169 multicellular neighborhoods. Specifically, fibroblasts separated regions of normal-like area with  
170 colon gland morphology (named Epi-Normal-like, **Extended Fig.5**, **Extended Fig.6**) and tumor  
171 epithelial areas at the margins (Epi-Tumor-left, Epi-Tumor-top, Epi-Tumor-right) compared to  
172 tumor in deposit (Epi-Tumor-inside)<sup>18,19</sup> (**Fig.2b**). Across tumor regions of interest (ROI), STHD  
173 labels delineated compositional heterogeneity of cell lineages and types (**Fig.2c**). Notably,  
174 spatial heterogeneity of immune cell populations was also observed; for instance, the Epi-  
175 Tumor-inside region contained macrophages<sup>20,21</sup> but few T cells compared to tumor regions at  
176 the margin, indicating distinct “hot” verses “cold” T cell infiltration within the same sample  
177 (**Fig.2c**, top, **Extended Fig.7**).

178  
179 Besides cell type compositional comparisons, STHD facilitates genome-wide and cell type-  
180 specific differential expression (DE) analyses across spatial regions. Comparing different tumor  
181 ROI, epithelial cells in Epi-Tumor-inside expressed highest stemness gene LGR5<sup>22</sup>, whereas

182 epithelial cells in Epi-Tumor-top highly express CCL20<sup>23</sup> and enrich pathways like TNF-alpha  
183 signaling (**Fig.2e, Supplementary Fig.12, Extended Fig.8**). Since the Epi-Tumor-inside region  
184 enrich macrophages, we compared differentially expressed genes specific to macrophage spots  
185 and observed that macrophage infiltrating inside tumor expressed higher extracellular matrix  
186 gene SPP1 and lower metalloelastase MMP12, cytokine CXCL8, and antigen presentation gene  
187 B2M, indicating that macrophage states segregated spatially in the same tumor<sup>21,24,25</sup> (**Fig.2f,**  
188 **Supplementary Fig.13**). The cell type-specific DE analyses can also be performed on the level  
189 of STHD-guided bins with higher depth, which resulted in a higher number of significant DE  
190 genes with well correlated effects (**Fig.2f, Extended Fig.8, Supplementary Fig.13**).

191  
192 We next demonstrated STHD in analyzing the less abundant immune cells after using  
193 STHDviewer to examine immune cell hubs at different locations in the colon cancer sample. An  
194 advantage of STHD is that rarer cell types and fine-grained cellular neighborhoods are  
195 maintained on top of the denoising effect for those abundant cell types forming continuous  
196 tissue architecture; an example of immune-endothelial cell interactions was demonstrated in  
197 **Extended Fig.9. Fig.2g** highlighted two T cell-rich areas marked in **Fig.2b**. The Immune-Trich-  
198 left region contained CD4T and CD8T cells with abundant macrophages and C1Q-high dendritic  
199 cells (DC)<sup>26</sup>, while the Immune-Trich-right area contained CD4T, CD8T cells, mature DCs  
200 enriched in immunoregulatory molecules (mregDC)<sup>27</sup>, and Plasma IgG B cells (**Fig.2h**). In  
201 neighborhood enrichment analyses, each area displayed unique co-occurrence patterns among  
202 immune cell types; for example, T cells in Immune-Trich-left area frequently interacted with  
203 macrophage and DCs, while T cells in Immune-Trich-right area co-occurred with DCs (**Fig.2g,i**).  
204 Utilizing cell type-specific DE analysis, we observed that T cells in the Immune-Trich-left region  
205 expresses higher cytotoxicity genes like GZMB comparing to T cells in the Immune-Trich-right  
206 region, which can be associated with the immunosuppressive function of mregDCs in tumor<sup>28,29</sup>  
207 (**Extended Fig.10**).

208

209 Cell-cell communications between cell types in disease microenvironment could drive disease  
210 progression like tumorigenesis<sup>30,31</sup>. Molecular interactions between cells types in single-cell  
211 transcriptomics analyses are traditionally inferred by evaluating averaged expression of  
212 communicating gene pairs, which lacks support on direct and physical interactions<sup>32</sup>. With the  
213 sub-cellular resolution spots labeled, we observed that STHD could pinpoint the precise  
214 frontlines of inter-cell type interaction. As shown in **Fig.2j**, we highlighted the frontlines of T cell-  
215 macrophage interaction and T cell-DC cell interaction in each of the T cell-rich regions. We  
216 hypothesized that cell-cell interactions near-frontline could be stronger than that in far-frontline  
217 spots. For a given cell type pair, we identified spots at boundary of two cell types based on  
218 posterior probabilities and considered STHD-labeled spots within 4 spots (comparable to single  
219 cell size) in distance as near-frontline spots (**Methods**). We calculated ligand-receptor  
220 interaction scores (LRIS)<sup>21,22</sup> between the near-frontline spots and compared with background  
221 interaction scores from far-frontline spots, conditioning on the same cell type pair. As shown in  
222 **Fig.2k**, near-frontline spots (grouped as 3,4) harbored higher LRIS activities compared to the  
223 far-frontline spots (grouped as 1,2) in both cell type pairs of T-Macrophage (**Fig.2j,k** top) and T-  
224 DC interactions (**Fig.2j,k** bottom). Specifically, the T-Macrophage communications in Immune-  
225 Trich-left were largely driven by CD74-IFNG interaction, while the T-DC communications in  
226 Immune-Trich-right were driven by CCL17/CCL19-CCR7, which indicated the immunoregulatory  
227 effect of mature DCs on T cells<sup>33</sup>.

228

229 In summary, we present STHD, a method for cell typing of individual spots for genome-wide  
230 spatial transcriptomics with high resolution. STHD addresses the computational challenge  
231 caused by data sparsity and proposes a fundamental strategy to enhance resolution compared  
232 to current standards. We demonstrated the accuracy, precision, and efficiency of STHD and its  
233 usability in multiple downstream analyses for global disease architecture, local cellular

234 neighborhoods, cell type-specific differential expression, cell type-augmented binning, and cell-  
235 cell interaction studies. STHD and STHDviewer provided efficient modeling, analyses, and  
236 examination for the entire sample. STHD allows incorporation of external tissue gap masks,  
237 better reference cell type signatures, and additional cell mixture deconvolution of ambiguous  
238 spots. Overall, we propose that modeling of individual high-resolution spots by borrowing  
239 neighborhood information can computationally eliminate the dilemma in resolution-throughput  
240 tradeoff for emerging ST technologies, which enables genome-wide and high-resolution  
241 discoveries in any tissues and diseases of intricate cellular heterogeneity.

242

### 243 **Data Availability**

244 The reference scRNA-seq atlas data and cell type annotation from Human Colon Cancer Atlas  
245 was downloaded from Pelka et. al through GSE178341:  
246 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE178341>). The reference data is also  
247 available through <https://portals.broadinstitute.org/crc-immune-hubs/> and Broad Institute Single  
248 Cell Portal (Human Colon Cancer Atlas c295) at  
249 [https://singlecell.broadinstitute.org/single\\_cell/study/SCP1162/human-colon-cancer-atlas-c295](https://singlecell.broadinstitute.org/single_cell/study/SCP1162/human-colon-cancer-atlas-c295).  
250 We downloaded the publicly available VisiumHD data demonstrated at  
251 <https://www.10xgenomics.com/support/software/loupe-browser/latest/tutorials/assay-analysis/lb-hd-spatial-gene-expression> from <https://www.10xgenomics.com/datasets/visium-hd-cytassist-gene-expression-libraries-of-human-crc>, named Visium HD Spatial Gene Expression Library,  
253 Human Colorectal Cancer (FFPE), by Space Ranger 3.0.0, 10x Genomics (2024, Mar 25).

255

### 256 **Code Availability**

257 The STHD software was implemented in Python and available on GitHub at  
258 <https://github.com/yi-zhang/STHD>. An interactive STHDviewer of the high-resolution cell type

259 result for the human colon cancer sample is available at [https://yi-zhang-compbio-lab.github.io/STHDviewer\\_colon\\_cancer\\_hd/STHDviewer\\_crchd.html](https://yi-zhang-compbio-lab.github.io/STHDviewer_colon_cancer_hd/STHDviewer_crchd.html).

261

## 262 **Acknowledgements**

263 C.S. and Y.Z. was supported by start-up funds to Y.Z. from Brain Tumor Omics Program,  
264 Preston Robert Tisch Brain Tumor Center, Department of Neurosurgery and Department of  
265 Biostatistics and Bioinformatics at Duke University School of Medicine. We acknowledge 10X  
266 Genomics to share VisiumHD datasets publicly available. We acknowledge authors of the  
267 Human Colon Cancer Atlas to make scRNAseq datasets publicly available.

268

## 269 **Author Information**

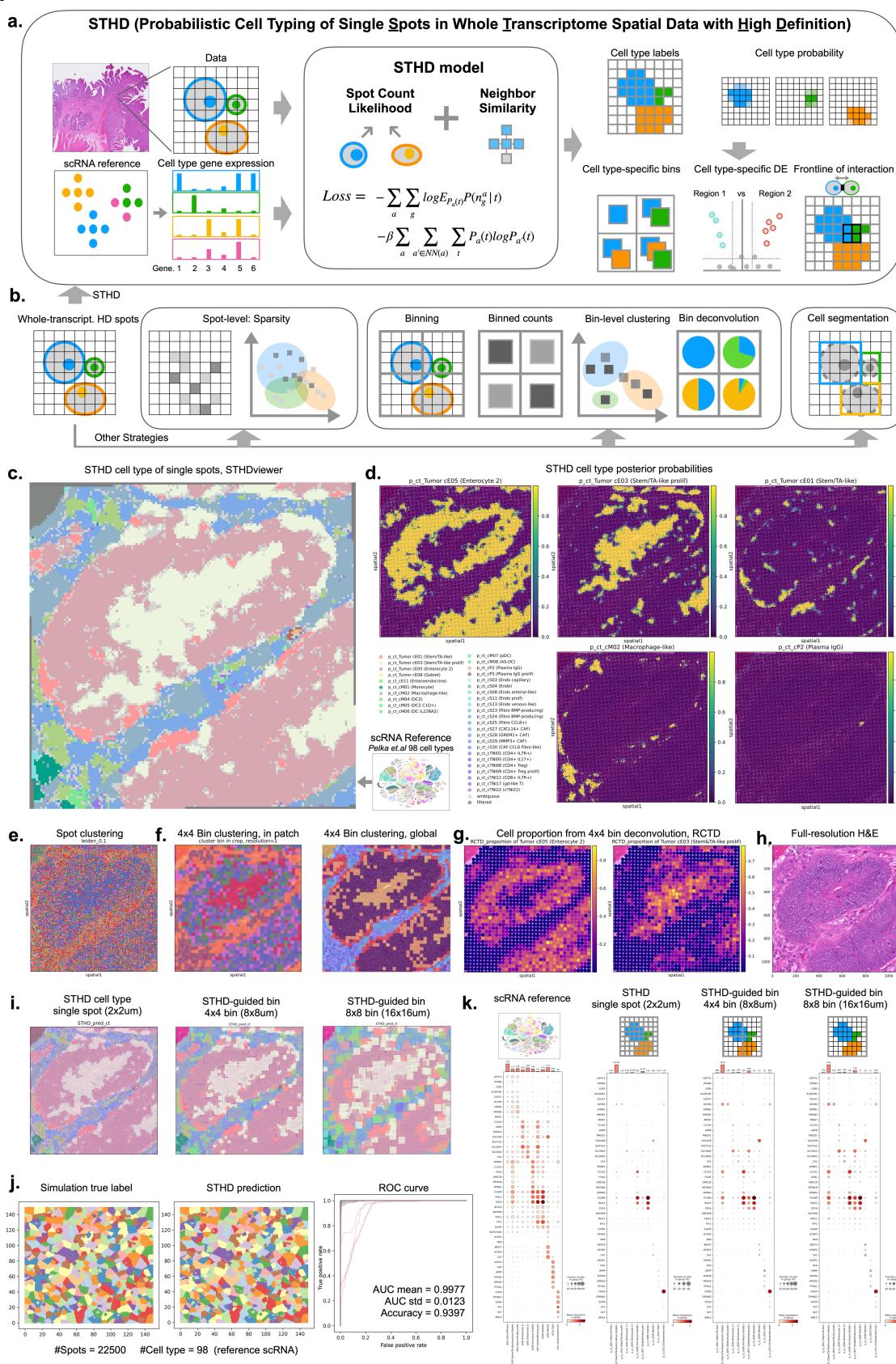
270 This study was conceived and led by Y.Z. Y.Z. designed the model and algorithm, implemented  
271 the STHD software, led the data analyses, and wrote the manuscript. C.S. analyzed the data  
272 and helped with manuscript writing.

273

## 274 **Figures**

275

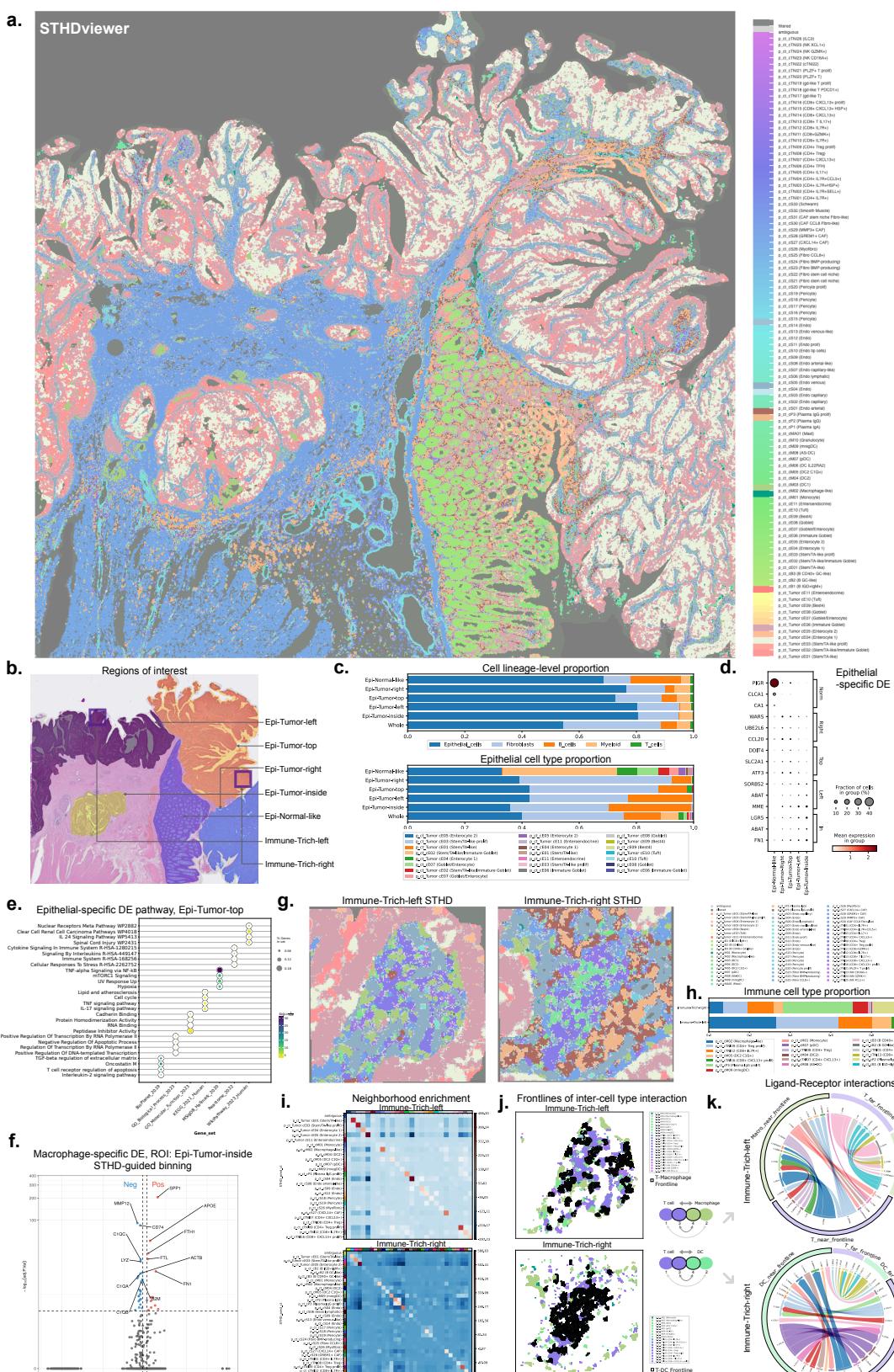
276 Fig.1



278 **Fig.1. Summary of STHD model. a.** STHD, (Probabilistic Cell Typing of Single Spots in Whole  
279 Transcriptome Spatial Data with High Definition) takes two inputs: whole-transcriptome spatial  
280 data of high resolution, and gene expression profiles by cell type from reference single-cell  
281 dataset. The STHD model is a machine learning method that infers latent cell type identity of  
282 individual spots by simultaneously modeling gene counts and neighbor similarity. STHD outputs  
283 cell type identities for each single spot in formats of cell type labels and posterior probabilities.  
284 The spot-level cell identities can then be used for multiple downstream analyses, including  
285 STHD-guided binning with cell type stratification, spatial comparison of cell type-specific  
286 compositional and differential expression analyses, and cell-cell interaction frontline analyses. **b.**  
287 Alternative approaches to analyze the HD datatype. Left, direct unsupervised clustering of spot-  
288 level gene expression is affected by sparsity. Middle, current binning-deconvolution approach  
289 increases coverage but compromise resolution. Right, ideal cell-level aggregation based on  
290 segmentation. **c.** Example of STHD results of a patch from the human colon cancer VisiumHD  
291 sample, revealing tumor-like epithelial cells and intercryptic immune population, visualized using  
292 the interactive STHDviewer. **d.** Cell type posterior probabilities at spot-level for the same patch,  
293 visualized using Squidpy. Top row, three epithelial cell types; Bottom row, two immune cell  
294 types. **e.** A spatial plot of unsupervised leiden clustering on spot gene expression. **f.** Spatial  
295 plots of unsupervised leiden clustering of gene expression aggregated by bins in size of 4x4  
296 spots. Left, clustering bins within patch, right, global cluster group of bins. **g.** Cell type proportion  
297 decomposed by RCTD on bins of size of 4x4 spots, demonstrating two tumor epithelial classes,  
298 enterocytes, and stem-like cells. **h.** The H&E histopathology image in full resolution of the same  
299 crop. **i.** The STHD predicted labels at spot-level and STHD-guided binning in bins of 4x4 spots  
300 and 8x8 spots, visualized using Squidpy. **j.** Simulated high-resolution spatial data. Left, ground  
301 truth cell type labels; middle, STHD predicted cell type labels; right, receiver operating  
302 characteristic curve (ROC) curve for all cell types. **k.** Expression dot plots for marker genes for  
303 normal epithelial cells. Left to right: normal epithelial cell types in reference human colon cancer

304 sample, normal epithelial cell types based on STHD predicted spots, normal epithelial cell types  
305 from STHD-guided bins of size 4x4 spots, normal epithelial cell types in STHD-guided bins of  
306 size 8x8 spots. AUC, area under the curve.  
307

308 Fig.2.



310 **Fig.2. STHD prediction facilitates segmentation of global structures and local**  
311 **neighborhoods.** **a.** Interactive spatial visualization by STHDviewer of predicted cell type labels  
312 for spots in the entire human colon cancer VisiumHD sample. **b.** Annotation of regions of  
313 interest, separating epithelial regions and two T cell-infiltrated regions. **c.** Cell type proportion  
314 bar plots at cell lineage and cell type level for epithelial regions and whole map. Upper:  
315 Proportions of major cell lineages including epithelial cells, fibroblasts, B cells, myeloid cells,  
316 and T cells. Bottom: Proportions of finer cell types of epithelial lineage. **d.** Spot-level gene  
317 expression for epithelial-specific differential genes across regions. Top three genes by log fold  
318 change were shown. **e.** Pathway enrichment of epithelial cell-specific genes differentially  
319 expressed in the Epi-Tumor-top region. **f.** Macrophage-specific differential genes for the Epi-  
320 Tumor-inside region, where differential analyses were performed at bin level of size 4x4 spots  
321 with cell type stratification from STHD. **g.** STHD cell type labels for two T cell-rich regions  
322 marked in panel b. **h.** Proportions of immune cell types in the two immune regions. **i.** Cell  
323 neighborhood interaction patterns of the two T cell-rich regions showing neighborhood  
324 enrichment z-scores where spots within 6 spots are considered neighbors. **j.** T cell and myeloid  
325 cells in the two immune-rich regions, with frontlines of inter-cell type interactions. Top, T cell-  
326 macrophage frontline in Immune-Trich-left region; Bottom, T cell-dendritic cell frontline in  
327 Immune-Trich-right region. Considering a pair of cell types, spots within 4 spots (8um) from the  
328 inter-cell type interaction location are highlighted as near-frontline spots, with the rest as far-  
329 frontline spots. **k.** The ligand-receptor communications at near-frontline spots compared to far-  
330 frontline communication scores of the same cell type pair as background, where the ligand-  
331 receptor interaction scores were calculated by squidpy.ligrec. Top 10 ligand-receptor pairs  
332 different between near-frontline and far-frontline communications were visualized where width  
333 represents average expression level of ligand receptor pairs. Spots are classified into four  
334 groups 1,2,3,4 based on their cell type and near- or far-frontline categories.  
335

336 **Methods**

337 **STHD algorithm**

338 The STHD model is a machine learning method that combine modeling of per-spot gene counts  
339 based on latent cell type identities with regularizing similarity among neighboring spots in terms  
340 of cell type probability distribution. The overall loss function  $\mathcal{L}$  composes of negative log  
341 likelihood  $-LL$  that models gene count of all spots and the cross entropy loss  $CE$  measuring  
342 similarity of neighboring spots. The amount of contribution from neighbors is controlled by the  
343 parameter  $\beta$ . The cell type identities at every single spot are thus estimated by Expectation  
344 Maximization.

345 
$$\mathcal{L} = -LL + \beta CE$$

346 For the negative log likelihood part, we model counts of each gene at each spot to follow  
347 Poisson distribution. Let the normalized average gene expression level for gene  $g$  within the  
348 latent cell type  $t$  is  $\lambda_g^t$ , read depth of spot  $a$  is  $d_a$ , and spot  $a$ 's probability of belonging to cell  
349 type  $t$  is  $P_a(t)$ . Thus, the count of gene  $g$  at spot  $a$ :  $n_g^a \sim Poission(d_a \lambda_g^t)$ .

350 
$$LL = \log \prod_a \prod_g P(n_g^a)$$
  
351 
$$= \sum_a \sum_g \log P(n_g^a)$$
  
352 
$$= \sum_a \sum_g \log E_{P_a(t)} P(n_g^a | t)$$
  
353 
$$\geq \sum_a \sum_g E_{P_a(t)} \log P(n_g^a | t)$$
  
354 
$$= \sum_a \sum_t P_a(t) \sum_g (n_g^a \log(d_a \lambda_g^t) - d_a \lambda_g^t) + const$$

355 Considering high sparsity of single spots, the second part of the loss function is related to cross  
356 entropy loss between spot  $a$  and nearest neighboring spots  $a' \in NN(a)$  in terms of distribution  
357 of cell type probability.

358 
$$CE = - \sum_a \sum_{a' \in NN(a)} \sum_t (P_a(t) \log P_{a'}(t))$$

359 Note that one constraint is  $\sum_t P_a(t) = 1$ . In calculation, we use the following for inherent

360 incorporation of constraint:  $P_a(t) := \frac{e^{w_a^t}}{\sum_{s \in T} e^{w_s^t}}$ , where  $T$  is the full set of cell types. We thus

361 optimize  $P_a(t)$  or  $w_a^t$  to minimize the total loss:

362 
$$Loss = - \sum_a \sum_t P_a(t) \sum_g (n_g^a \log(d_a \lambda_g^t) - d_a \lambda_g^t) - \beta \sum_a \sum_{a' \in NN(a)} \sum_t (P_a(t) \log P_{a'}(t))$$

363

364 **Construct normalized gene expression profile from reference single-cell atlas dataset**

365 The average normalized gene expression  $\lambda_g^t$  can be estimated from a pre-defined single-cell

366 RNA-seq reference dataset. A reference single-cell dataset shall match the tissue or disease

367 origin of the high resolution, whole-transcriptome spatial transcriptomics data. It shall also

368 contain cell type annotation that is usually based on multi-sample harmonization and

369 unsupervised clustering of single cells into various cell type or groups. The granularity of cell

370 type annotation directly influence the number of predicted cell type and downstream analyses.

371 For each cell type, average normalized gene expression profiles are estimated with  $\lambda_g^t = \frac{n_g^t}{\sum_g n_g^t}$

372 so that  $\sum_g \lambda_g^t = 1$ . Specifically for human colon cancer, we downloaded the raw counts file of the

373 single-cell RNA-seq data from the human colon cancer atlas from Pelka et.al<sup>13</sup> through

374 GSE178341. The data contains 370115 tumor and adjacent normal single cells from 62 primary

375 treatment-naïve colorectal cancer patients, with 43113 genes profiled. Quality control of single

376 cells was performed by filtering cells with minimum 100 counts and 100 genes, filtering genes

377 with minimum 3 cells, and removing cells with mitochondrial read ratio larger than 27.7%, which

378 is mean plus one standard deviation of the mitochondrial read ratio across all cells. After the

379 above quality control, 303358 cells and 31,873 genes remain. The cell annotations were

380 obtained from the study that provided 98 fine-grained cell types and 7 cell lineages. We used

381 the 98 high-granularity cell type annotation throughout this study. To select genes informative in  
382 distinguishing the 98 cell types, we followed the calculation of cell type expression profile as in  
383 RCTD<sup>7</sup>. First, we calculated the estimated cell type expression profiles  $\lambda_g^t$  for all genes and cell  
384 types. We then selected genes with average expression profile above 0.000125, which is  
385 0.0625 counts per 500, and at least 0.5 log-fold-change compared to the average expression  
386 across all cell types, using log2 scale. Moreover, we removed the mitochondrial genes and  
387 ribosome genes in the selection. The selection resulted in 4618 genes that are cell type-  
388 informative, for each we calculated  $\lambda_g^t$  in all cell types. Note that STHD uses the 4618 genes to  
389 infer cell type identity of spots, and later can consider whole-transcriptome for downstream  
390 analyses.

391

## 392 **STHD hyperparameter tuning**

393 To select hyperparameters that balance between statistical likelihood and neighborhood  
394 similarity, we comprehensively tested hyperparameter tuning and provided an automatic  
395 stopping criterion and optional tuning tutorial. The main parameters in STHD includes the  
396 neighborhood parameter  $\beta$ , which controls the contribution of cross entropy loss among local  
397 neighbors to the total loss, the optimization steps n\_iteration, and the learning rate that can be  
398 fixed in practice. We provided a tutorial on parameter tuning based on performance in random  
399 sampled patches; for the human colon cancer sample, the optimized parameters are 23 iteration  
400 steps and  $\beta=0.1$ . In detail, for the colon cancer sample, we randomly selected six different  
401 patches in size of 1100x1100 pixels (in full resolution, around 22500 spots). For the optimization  
402 of each patch, we tested a range of  $\beta$  from 0, 0.03, 0.1, 0.3, 1, and 3, each outputting optimal  
403 iteration steps based on the early stopping criterion, where the drop of loss in the past 10 steps  
404 falls within 1%. Since the regularization term from neighborhood contribution  $\beta$  will lower log  
405 likelihood, as  $\beta$  increases, we select the optimal parameter pairs where drop in cross entropy

406 slows down to around 10% while a high log likelihood is maintained. The optimal iteration steps  
407 when fixing  $\beta$  of 0, 0.03, 0.1, 0.3, 1, 3, is each 24, 23, 22, 23, 23, 24. The optimal iteration steps  
408 is robustly 23.17 on average. Therefore, we chose  $\beta=0.1$  and  $n\_iteration=23$  as optimal  
409 parameters for the whole sample modeling that balance the count-based likelihood and the  
410 contribution from neighbors.

411

#### 412 **STHD low-count region detection and filter**

413 Throughout the spatial transcriptomics data, regions of low gene counts exist, usually arising  
414 from tissue boundaries, tissue gaps, and regions low in DNA/RNA such as dead cells, necrosis,  
415 vessel enriching red blood cells, and collagen rich area. We included an automatic local low-  
416 count region detection and filtering functionality in STHD. By default, STHD computed neighbor  
417 graphs within 2 rings of spots (13 spots for the grid system), aggregated local connected spots,  
418 and detected low-count local regions not exceeding 50 counts, where the model will skip  
419 modeling due to low confidence in cell type inference. Other customized tissue gap and  
420 boundary masks can also be provided to further filter spots.

421

#### 422 **Assessing model accuracy with simulated spatial transcriptomics data**

423 To assess the accuracy of STHD model prediction, we simulated the high-dimensional, high-  
424 resolution spatial transcriptomics data *in silico*. We simulated a region with 22500 (150x150)  
425 grid spots and position 500 simulated cells overlapping the region with cell sizes spanning 1-20  
426 spots. Each cell is randomly assigned a cell type from the 98 cell types annotated in colon  
427 cancer atlas, and spots overlapping each whole cell are assigned the corresponding cell type.  
428 To mimic the real-world scenario in gene expression profiles and spot sparsity, we adapted the  
429 same normalized gene expression profile from the same colon cancer atlas data with 98 cell  
430 types. The depth of each spot was simulated following the distribution of exponential decay with  
431 mean read depth 47 matching the colon cancer VisiumHD sample. After simulating counts of

432 each gene at each spot in the region, we applied STHD with default optimal parameters ( $\beta=0.1$ )  
433 and also STHD no-neighbor model ( $\beta=0$ ) to the simulated spatial data, and measured prediction  
434 accuracy, standard error, and area under the receiver operating characteristic curve (ROC), for  
435 this task of multi-class classification of 98 categories at each spot.

436

#### 437 **Comparison with other cell typing methods**

438 To benchmark the STHD results at  $2 \times 2 \mu\text{m}$  resolution, we compared STHD to six methods that  
439 can output cell typing, including the deconvolution-based cell typing tools RCTD<sup>7</sup>, CARD<sup>8</sup>,  
440 Cell2Location<sup>9</sup>, Celloscope<sup>15</sup>, Stereoscope<sup>16</sup>, and the unsupervised method Banksy<sup>17</sup>. Because  
441 of computational efficiency and data size limit of the existing methods, we used the tumor region  
442 as in Fig.1 to perform benchmarking, which contains 22,336 cells and 18,085 genes. For methods  
443 requiring the reference scRNA-seq data as input, we used the same human colorectal cancer  
444 scRNA-seq datasets (GSE178341) with 98 cell types. Since many methods in comparison  
445 requires a maximum number of cells from reference scRNaseq dataset – for example, CARD  
446 requires below 50000 cells - we down-sampled the reference data to include maximum 600 cells  
447 for each cell type, resulting in a total of 49,351 cells and 41,195 genes. For running RCTD on  
448 single spots, we used RCTD doublet mode, testing both the default minimum UMI 100 and also  
449 by removing the minimum UMI requirement due to loss of most raw spots. For testing CARD, we  
450 followed guidelines from <https://github.com/YMa-lab/CARD>. Similarly, to avoid most spots being  
451 filtered, the minimum spot count and minimum gene numbers were set to 0. Cell2location and  
452 Stereoscope were ran with default parameters according to the tutorials  
453 [https://cell2location.readthedocs.io/en/latest/notebooks/cell2location\\_tutorial.html](https://cell2location.readthedocs.io/en/latest/notebooks/cell2location_tutorial.html), and  
454 <https://github.com/Almaan/Stereoscope>. In Cell2location, by default, the reference data was  
455 filtered using the following parameters: cell\_count\_cutoff = 5, cell\_percentage\_cutoff = 0.03, and  
456 nonz\_mean\_cutoff = 1.12. The reference single-cell regression model was trained with default  
457 parameter max\_epochs = 250, lr = 0.002. The spatial cell2location model was obtained with

458 default parameters max\_epoch = 30,000. In Stereoscope, by default, the reference data was  
459 filtered using a minimum count threshold of min\_counts = 10, and all mitochondrial genes were  
460 removed. The reference single-cell model was trained with parameters max\_epochs = 100; the  
461 spatial stereoscope model was obtained using parameters max\_epochs = 10,000. Finally, for  
462 Celloscope, a binary matrix with prior knowledge about marker genes is required among input,  
463 for which we selected top 200 genes in normalized expression from our informative gene matrix  
464 for each cell type. Another required input is estimated number of cells in each spot, for which we  
465 set to 1 due to sub-cellular-resolution nature of the data. Celloscope was then tested following  
466 guidelines at <https://github.com/szczurek-lab/Celloscope>. Finally, Banksy was tested with default  
467 parameters following the tutorial at  
468 [https://github.com/prabhakarlab/Banksy\\_py/blob/main/slideseqv2\\_analysis.ipynb](https://github.com/prabhakarlab/Banksy_py/blob/main/slideseqv2_analysis.ipynb). For all toolkits,  
469 the weights or cell type abundance for each spot were extracted and normalized as cell  
470 proportions for subsequent visualizations.  
471

#### 472 **Comparison with deconvoluted cell proportions in aggregated bins**

473 Following current standards of binning-deconvolution, we binned the spatial data and aggregated  
474 the bins containing 4x4 spots. We applied RCTD deconvolution of doublet mode on the same  
475 spatial patch as in Fig.1, following the default filtering and parameters of RCTD as on spacexr  
476 (<https://github.com/dmcable/spacexr>), and using the same subsampled colon cancer reference  
477 scRNA-seq data as in comparison with other methods. To compare cell proportions across bins,  
478 we counted STHD cell type labels and obtained cell proportions for each 4x4 bin after excluding  
479 filtered and ambiguous spots. We also computed cosine similarity between STHD's spot count-  
480 based proportions and RCTD's deconvoluted proportions and visualize the proportions and  
481 cosine similarity on a spatial plot.  
482

#### 483 **Timing computational efficiency**

484 To assess computational efficiencies, different cell typing tools are applied onto the tumor  
485 epithelial patch using a server with 8 CPU cores and maximum 128GB memory, running time  
486 recorded. For tools supporting GPU acceleration (cell2location and stereoscope) the time cost  
487 was evaluated on a single GPU, and the spatial training time was recorded without including  
488 scRNA training time. For computational efficacy, we also estimated the time cost for the entire  
489 sample by linearly scaling running time based on number of spots in patch (22,336 spots) and in  
490 the entire sample (8,726,224 spots).

491

#### 492 **STHD for full-size ST sample**

493 The STHD software contains several components to flexibly split the whole sample into patches  
494 for inferences and parallelization. The input/output functionalities (STHDio) expanded from  
495 Squidpy Visium interface to enable cropping of full-resolution image. Second, for the full-size  
496 sample, STHD will patchify the full sample into multiple regions with overlaps, followed by low-  
497 count filtering, per-patch inference, and result aggregation where probabilities of overlapping  
498 spots are averaged. The STHD cell type label results are determined by Maximum A Posterior  
499 (MAP) with a preset cutoff (default 0.8). An optional fast implementation for STHD-guided  
500 binning with a predefined bin size is provided, which outputs per-gene count for each bin  
501 stratified by STHD cell type labels. Finally, the predicted cell labels are exported to STHDviewer  
502 for a rasterized and interactive scatter plot implemented using Bokeh<sup>34</sup>. In STHDviewer, the spot  
503 grid of array\_row and array\_column were used for x-axis and y-axis for visualization efficiency.

504

#### 505 **Cell type-specific differential gene expression analysis across spatial regions of interest**

506 To study regional differences within the cancer sample, we first focused on regions enriching  
507 epithelial cell types and manually separated four different regions and extracted spot barcodes.  
508 Based on STHD predicted labels, we performed differential expression analyses specific for

509 epithelial cell types. Gene counts were normalized towards the spot count depth (target total  
510 counts 10000) and log-transformed followed by Wilcoxon test in Scanpy. Genes passing  
511 adjusted p-value cutoff  $10^{-6}$  and absolute log2 fold change 0.5 were selected as significant, and  
512 top genes passing log2 fold change 0.5 are exported to pathway enrichment analyses by  
513 GSEAp<sup>35</sup> using EnrichR and a few pathway databases, including GO Molecular Function 2023,  
514 GO Biological Process 2023, MSigDB Hallmark 2020, Reactome 2022, WikiPathway 2023  
515 Human, KEGG 2021 Human, and BioPlanet 2019. The FDR cutoff of 0.1 was used to filter the  
516 pathways and top four pathways were visualized for tumor epithelial pathway comparison. We  
517 also tested cell type-specific differential expression using both STHD-labeled spot and STHD-  
518 guided bin aggregates of size 4x4 spots. At STHD-guided bin level, we followed the same step  
519 of normalization and performed differential expression with Wilcoxon with the same cutoff for  
520 epithelial and macrophage comparison. For T cell-specific differential expression analyses  
521 between the two T cell rich regions, the cutoff was set log2 fold change 0.5 and adjusted p-  
522 value cutoff 0.05 due to small number of spots from patches. For differential gene dotplots,  
523 significant genes are sorted by log fold change and the top three genes for the tumor  
524 comparison and top ten genes for the immune comparison were plotted.

525

## 526 **Neighborhood interaction enrichment analysis for immune-rich regions**

527 To investigate the neighborhood interaction patterns across various cell types in the two immune-  
528 infiltrated regions, local neighborhoods were calculated using the 4 nearest cells in the grid and  
529 considering 6 rings of adjacent spots to obtain a local neighborhood structure using  
530 squidpy.gr.spatial\_neighbors in Squidpy. We used STHD predicted cell type labels and excluded  
531 spots in the filtered and ambiguous categories. The neighborhood enrichment Z-scores were  
532 calculated using squidpy.gr.nhood\_enrichment, which follows CellChat<sup>32</sup> with ligand-receptor

533 databases and permutation of 1000 times. Heatmaps of neighborhood interaction patterns were  
534 plotted based on the co-occurrence score of pairwise cell types in each region of interest.

535

536 **Inter-cell type frontline and cell-cell communication analysis** **frontline computation**

537 For a given pair of cell types, we first defined their precise frontlines to be spots where those two  
538 cell types are adjacent to each other. We further restrict frontline spots to be those whose neighbor  
539 only contains that pair of cell types based on maximum posterior to eliminate potential noises  
540 from other cell types. Considering a single cell is comparable to the size of 4 spots, we assumed  
541 spots within 4 spots from the frontlines represented one cell with direct interaction with another  
542 cell type. Thus, spots were categorized based on their cell type and proximity to frontline. We  
543 examined interactions between T cells and Dendritic cells, and T cells and Macrophages  
544 separately for the two immune regions of interest, Immune-Trich-left and Immune-Trich-right. In  
545 the chord plots, For a pair of cell types A and B, the labels were assigned as follows: '1' for far-  
546 frontline cell type A, '2' for far-frontline cell type B, '3' for near-frontline cell type A, and '4' for near-  
547 frontline cell type B. We then hypothesized that cell-cell communications like ligand-receptor  
548 expression levels were higher between group 3 and 4 compared to group 1 and 2. To perform the  
549 ligand-receptor interaction analysis, normalized data was applied the squidpy.gr.ligrec function  
550 which permuted 1000 times to detect significant ligand-receptor pairs. Top 10 ligand-receptor  
551 pairs different between near-frontline and far-frontline communications were visualized in the  
552 chord diagram showing the mean expression of each ligand-receptor pair and sender-receiver  
553 directions.

554

555 **References**

- 556 1. Tian, L., Chen, F. & Macosko, E. Z. The expanding vistas of spatial transcriptomics. *Nat*  
557 *Biotechnol* **41**, 773–782 (2023).

- 558 2. Stickels, R. R. *et al.* Highly sensitive spatial transcriptomics at near-cellular resolution with  
559 Slide-seqV2. *Nat Biotechnol* **39**, 313–319 (2021).
- 560 3. Chen, A. *et al.* Spatiotemporal transcriptomic atlas of mouse organogenesis using DNA  
561 nanoball-patterned arrays. *Cell* **185**, 1777-1792.e21 (2022).
- 562 4. Wang, X. *et al.* Three-dimensional intact-tissue sequencing of single-cell transcriptional  
563 states. *Science* **361**, eaat5691 (2018).
- 564 5. Moffitt, J. R. *et al.* Molecular, spatial, and functional single-cell profiling of the hypothalamic  
565 preoptic region. *Science* **362**, eaau5324 (2018).
- 566 6. Hao, Y. *et al.* Dictionary learning for integrative, multimodal and scalable single-cell analysis.  
567 *Nat Biotechnol* **42**, 293–304 (2024).
- 568 7. Cable, D. M. *et al.* Robust decomposition of cell type mixtures in spatial transcriptomics. *Nat*  
569 *Biotechnol* **40**, 517–526 (2022).
- 570 8. Ma, Y. & Zhou, X. Spatially informed cell-type deconvolution for spatial transcriptomics. *Nat*  
571 *Biotechnol* **40**, 1349–1359 (2022).
- 572 9. Kleshchevnikov, V. *et al.* Cell2location maps fine-grained cell types in spatial transcriptomics.  
573 *Nat Biotechnol* **40**, 661–671 (2022).
- 574 10. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. Preprint at  
575 <https://doi.org/10.48550/arXiv.1412.6980> (2017).
- 576 11. Lam, S. K., Pitrou, A. & Seibert, S. Numba: a LLVM-based Python JIT compiler. in  
577 *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC* 1–6  
578 (Association for Computing Machinery, New York, NY, USA, 2015).  
579 doi:10.1145/2833157.2833162.
- 580 12. Hartman, A. & Satija, R. Comparative analysis of multiplexed *in situ* gene expression  
581 profiling technologies. 2024.01.11.575135 Preprint at  
582 <https://doi.org/10.1101/2024.01.11.575135> (2024).

- 583 13. Pelka, K. *et al.* Spatially organized multicellular immune hubs in human colorectal  
584 cancer. *Cell* **184**, 4734-4752.e20 (2021).
- 585 14. Merlos-Suárez, A. *et al.* The intestinal stem cell signature identifies colorectal cancer  
586 stem cells and predicts disease relapse. *Cell Stem Cell* **8**, 511–524 (2011).
- 587 15. Geras, A. *et al.* Celloscope: a probabilistic model for marker-gene-driven cell type  
588 deconvolution in spatial transcriptomics data. *Genome Biology* **24**, 120 (2023).
- 589 16. Andersson, A. *et al.* Single-cell and spatial transcriptomics enables probabilistic  
590 inference of cell type topography. *Commun Biol* **3**, 1–8 (2020).
- 591 17. Singhal, V. *et al.* BANKSY unifies cell typing and tissue domain segmentation for  
592 scalable spatial omics data analysis. *Nat Genet* 1–11 (2024) doi:10.1038/s41588-024-01664-  
593 3.
- 594 18. Höppener, D. J. *et al.* The relationship between primary colorectal cancer histology and  
595 the histopathological growth patterns of corresponding liver metastases. *BMC Cancer* **22**,  
596 911 (2022).
- 597 19. Sirinukunwattana, K. *et al.* Gland segmentation in colon histology images: The glas  
598 challenge contest. *Med Image Anal* **35**, 489–502 (2017).
- 599 20. Ozato, Y. *et al.* Spatial and single-cell transcriptomics decipher the cellular environment  
600 containing HLA-G+ cancer cells and SPP1+ macrophages in colorectal cancer. *Cell Rep* **42**,  
601 111929 (2023).
- 602 21. Zhang, Y. *et al.* MetaTiME integrates single-cell gene expression to characterize the  
603 meta-components of the tumor immune microenvironment. *Nature Communications* **14**, 2634  
604 (2023).
- 605 22. Xu, L., Lin, W., Wen, L. & Li, G. Lgr5 in cancer biology: functional identification of Lgr5 in  
606 cancer progression and potential opportunities for novel therapy. *Stem Cell Res Ther* **10**, 219  
607 (2019).

- 608 23. Kapur, N. *et al.* CCR6 expression in colon cancer is associated with advanced disease  
609 and supports epithelial-to-mesenchymal transition. *Br J Cancer* **114**, 1343–1351 (2016).
- 610 24. Cheng, S. *et al.* A pan-cancer single-cell transcriptional atlas of tumor infiltrating myeloid  
611 cells. *Cell* **184**, 792-809.e23 (2021).
- 612 25. Matusiak, M. *et al.* Spatially Segregated Macrophage Populations Predict Distinct  
613 Outcomes In Colon Cancer. *Cancer Discovery* (2024) doi:10.1158/2159-8290.CD-23-1300.
- 614 26. Baruah, P. *et al.* C1q enhances IFN-gamma production by antigen-specific T cells via  
615 the CD40 costimulatory pathway on dendritic cells. *Blood* **113**, 3485–3493 (2009).
- 616 27. Maier, B. *et al.* A conserved dendritic-cell regulatory program limits antitumour immunity.  
617 *Nature* **580**, 257–262 (2020).
- 618 28. Kvedaraite, E. & Ginhoux, F. Human dendritic cells in cancer. *Science Immunology* **7**,  
619 eabm9409 (2022).
- 620 29. Puleston, D. J. *et al.* Polyamine metabolism is a central determinant of helper T cell  
621 lineage fidelity. *Cell* **184**, 4186-4202.e20 (2021).
- 622 30. Wang, X. *et al.* In vivo CRISPR screens identify the E3 ligase Cop1 as a modulator of  
623 macrophage infiltration and cancer immunotherapy target. *Cell* **184**, 5357-5374.e22 (2021).
- 624 31. Burdziak, C. *et al.* Epigenetic plasticity cooperates with cell-cell interactions to direct  
625 pancreatic tumorigenesis. *Science* **380**, eadd5327 (2023).
- 626 32. Jin, S. *et al.* Inference and analysis of cell-cell communication using CellChat. *Nat  
627 Commun* **12**, 1088 (2021).
- 628 33. Li, J. *et al.* Mature dendritic cells enriched in immunoregulatory molecules (mregDCs): A  
629 novel population in the tumour microenvironment and immunotherapy target. *Clin Transl Med*  
630 **13**, e1199 (2023).
- 631 34. Bokeh Development Team. *Bokeh: Python Library for Interactive Visualization*. (2018).
- 632 35. Fang, Z., Liu, X. & Peltz, G. GSEAp: a comprehensive package for performing gene set  
633 enrichment analysis in Python. *Bioinformatics* **39**, btac757 (2023).