

# SEAM: a Spatial single nucleAr metAboloMics method for dissecting tissue microenvironment

Michael Zhang (✉ [michael.zhang@utdallas.edu](mailto:michael.zhang@utdallas.edu))

University of Texas at Dallas

Zhiyuan Yuan

Tsinghua University

Qiming Zhou

Tsinghua University

Lesi Cai

Tsinghua University

Lin Pan

China-Japan Friendship Hospital

Weiliang Sun

China-Japan Friendship Hospital

Shiwei Qumu

China-Japan Friendship Hospital

Si Yu

Peking Union Medical College Hospital

Yongchang Zheng

Peking Union Medical College Hospital

Shao Li

Tsinghua University

Yang Chen

Tsinghua University

Xinrong Zhang

Tsinghua University

---

## Article

**Keywords:** SEAM, tissue organization, spatial metabolome

**Posted Date:** August 28th, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-63938/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Nature Methods on October 4th, 2021. See the published version at <https://doi.org/10.1038/s41592-021-01276-3>.

1 **SEAM: a Spatial single nucleAr metAboloMics method for dis-**  
2 **secting tissue microenvironment**

3 Zhiyuan Yuan<sup>1,9</sup>, Qiming Zhou<sup>2,9</sup>, Lesi Cai<sup>3,9</sup>, Lin Pan<sup>6</sup>, Weiliang Sun<sup>6</sup>, Shiwei Qumu<sup>7</sup>, Si  
4 Yu<sup>8</sup>, Yongchang Zheng<sup>8</sup>, Shao Li<sup>1</sup>, Yang Chen<sup>1\*</sup>, Xinrong Zhang<sup>3\*</sup>, Michael Q. Zhang<sup>1,4,5\*</sup>

5 1, Ministry of Education Key Laboratory of Bioinformatics, Center for Synthetic and Sys-  
6 tems Biology, Department of Automation, BNRist, Tsinghua University, Beijing 100084,  
7 China

8 2, School of Life Sciences, Tsinghua University, Beijing 100084, China

9 3, Department of Chemistry, Tsinghua University, Beijing, 100084, China

10 4, Department of Biological Sciences, Center for Systems Biology, The University of Texas,  
11 Richardson, TX 75080-3021, USA

12 5, Department of Basic Medical Sciences, School of Medicine, Tsinghua University, Beijing  
13 100084, China

14 6, Institute of Clinical Medicine, China-Japan Friendship Hospital; National Clinical Re-  
15 search Center for Respiratory Diseases; Institute of Respiratory Medicine, Chinese Acad-  
16 emy of Medical Science, Beijing 100029, China

17 7, Department of Pulmonary and Critical Care Medicine, China-Japan Friend Hospital; Na-  
18 tional Clinical Research Center for Respiratory Diseases, Beijing 100029, China

19 8, Department of Liver Surgery, Peking Union Medical College Hospital, Chinese Academy  
20 of Medical Sciences and Peking Union Medical College, Beijing 100730, China

21 9, These authors contributed equally: Zhiyuan Yuan, Qiming Zhou and Lesi Cai

22

23 \*email:

24 [yc@mail.tsinghua.edu.cn](mailto:yc@mail.tsinghua.edu.cn)

25 [xrzhang@mail.tsinghua.edu.cn](mailto:xrzhang@mail.tsinghua.edu.cn)

26 [michael.zhang@utdallas.edu](mailto:michael.zhang@utdallas.edu)

27

28

## 29 **Abstract**

30 Spatial metabolomics can reveal intercellular heterogeneity and tissue organization. To achieve highest  
31 spatial resolution, we reported a novel Spatial single nucleAr metAboloMics (*SEAM*) method, a scalable  
32 platform combining high resolution imaging mass spectrometry (IMS) and a series of computational algo-  
33 rithms, that can display multiscale/multicolor tissue tomography together with identification and clustering  
34 of single nuclei by their *in situ* metabolic fingerprints. We firstly applied *SEAM* to a range of wild type  
35 mouse tissues, then delineate a consistent pattern of metabolic zonation in mouse liver. We further studied  
36 spatial metabolome in human fibrotic liver. Intriguingly, we discovered novel subpopulations of hepato-  
37 cytes with special metabolic features associated with their proximity to fibrotic niche, which was further  
38 validated by spatial transcriptomics with Geo-seq. These demonstrations highlight how *SEAM* may be  
39 used to explore the spatial metabolome and tissue anatomy at single cell level, hence leading to a deeper  
40 understanding of the tissue metabolic organization.

## 41 **Introduction**

42 The hierarchical organization of multicellular organisms is stably maintained by homeostasis at different  
43 levels. At the tissue level, such homeostasis is often further modulated by the combination of intracellular  
44 gene expression network and extracellular (microenvironmental) signals<sup>1-4</sup>. Cell and its extracellular en-  
45 vironment interact dynamically through various signaling mediators, including metabolites, secretome,  
46 and ligand-receptor interactions. Metabolites from extracellular environment can significantly influence  
47 cell behavior or even transform its identity. For instance, extensive alcohol intake not only activates the  
48 detoxification activity of hepatocytes but also alters the epigenetic landscape of hepatocytes<sup>5</sup>. Conversely,  
49 cell releasing metabolites can also have impact on its microenvironment. One classic example is baso-  
50 phils and mast cells releasing histamine to increase the permeability of the capillaries when encountering  
51 infection<sup>6</sup>. To facilitate a deeper and more systematic understanding of the multi-scale nature of biological  
52 processes (e.g. organ development or tumor microenvironment), various single cell omics-technologies  
53 have been rapidly developed and utilized<sup>7</sup>. Currently, advanced imaging mass spectrometry (IMS) based  
54 techniques are also being made possible to profile a large number of metabolites spatially and/or tempo-  
55 rally, providing new dimensional insights to those hierarchical processes<sup>8,9</sup>.

56 In spatially resolved metabolomics studies, different techniques have been developed including ma-  
57 trix-assisted laser desorption/ionization (MALDI-MS)<sup>10</sup>, desorption electrospray ionization (DESI-MS)<sup>11</sup>,  
58 laser ablation inductively coupled plasma (LA-ICP-MS)<sup>12</sup>, and secondary ion mass spectrometry  
59 (SIMS)<sup>13</sup>. MALDI-MS utilized t-MALDI ion source for imaging of phospholipids and a few other biomole-  
60 cule classes in thin, matrix-coated tissue sections and cell cultures at a pixel size of about 1–2  $\mu\text{m}$ <sup>14</sup>.  
61 With further improvement, MALDI-2 was introduced by adapting a t-MALDI-2 ion source to an Orbitrap  
62 mass analyzer and a pixel size of 600 nm was achieved on brain tissue<sup>15</sup>. DESI-MS has been utilized to  
63 visualize tissue level metabolomic alterations in 256 esophageal cancer patients<sup>11</sup>. Recently, based on  
64 SIMS, 3D OrbiSIMS, a label-free IMS with subcellular lateral resolution, and high mass-resolving power,  
65 has been developed<sup>16</sup>. These techniques will increasingly be used in future spatial metabolomics appli-  
66 cations.

67 Although the above techniques achieved unprecedented subcellular resolution, several analytical  
68 complications still exist, e.g. single cell segmentation and cell fingerprint extraction. Previous studies  
69 typically segmented cells using hematoxylin-eosin (H&E) staining, which suffered from either inaccurate  
70 segmentation due to imperfect registration of adjacent slides, or labeling on the same slides, which  
71 might bring exogenous substances leading to sabotaging sample integrity<sup>17</sup>. Another cell segmentation  
72 strategy exploited convolution neural network (CNN) trained on pixel-wise annotated cells, demanding  
73 for huge human expert labour<sup>18</sup>. As for cell fingerprint extraction, the common practice that took the av-  
74 erage of pixel profiles within each cell caused the impairment of distributive information<sup>19,20</sup>. These defi-  
75 ciencies hinder the efforts for the quantification of single cell metabolome while preserving spatial infor-  
76 mation. Consequently, although there have been instrumental-wise improvements for IMS, the down-  
77 stream analytical methods still require further development for users to fully exploit spatial metabolomic  
78 features.

79 To overcome those deficiencies, we proposed Spatial single nucleAr metAboloMics (*SEAM*), a novel  
80 platform leveraging the spatial metabolome provided by SIMS and a comprehensive series of computa-  
81 tional algorithms for delineating *in situ* single cell level metabolome and tissue microenvironment. To our  
82 knowledge, this is the first study capable of segmenting and analyzing single nuclear metabolic profiles  
83 directly on tissue sections. Importantly, *SEAM* is label-free and only requires minimal experimental prep-  
84 aration, which avoids the introduction of exogenous substance and preserves samples' native state. As  
85 a proof of principle, we comprehensively calibrated *SEAM* using popular cell cultures, and then system-  
86 atically scaled up to various mouse tissues, including wild type mouse lung, kidney, small intestine, and  
87 liver. Finally, we discovered different hepatocyte metabolic subpopulations and their spatial network or-  
88 ganization within the tissue microenvironment in human fibrotic liver.

89

## 90 **Results**

### 91 **Overview**

92 *SEAM* is an integrated platform for qualitative and quantitative analysis of tissue metabolic cell typing and  
93 *in situ* microenvironment. The whole pipeline is composed of two main parts: IMS assay and computa-  
94 tional analysis suite (Fig. 1a).

95 As an IMS technology, time-of-flight secondary ion mass spectrometry (TOF-SIMS) provides both mass  
96 spectra (chemical information) and ion images (spatial information), of biomolecules on tissue sections  
97 (Fig 1a, top left). Typically, hundreds of peaks in a mass spectrum could be extracted from a 400 × 400  
98 μm<sup>2</sup> scan area on a tissue section. Every experiment outputs multiplex SIMS data with 256×256 pixels in  
99 spatial resolution, and each pixel is associated with a vector of over 200 selected *m/z* peaks (Fig. 1a and  
100 see Methods). With the reference of H&E staining, to facilitate users with quickly viewing of the metabolic  
101 spatial pattern across the full spectrum, rather than manually reading hundreds of *m/z* images one by one,  
102 *SEAM* provides *SIMS-View* to compress the multiplex SIMS images from hundreds of channels into three,  
103 while preserving local and global structures in the feature space (Fig. 1a, bottom left and middle). Then

104 the three-channel images are mapped to CIELAB color spaces<sup>21</sup> and can be rapidly surveyed by human  
105 vision.

106 To compensate for the potential information loss of dimensionality reduction by taking the advantage of  
107 compositional characteristics and spatial continuity, *SEAM* can further build a spatial single nucleus map  
108 and delineate the organization of metabolically distinct *in situ* cell subpopulations (Fig. 1a bottom right).  
109 More specifically, *SEAM* provides three additional data analysis modules (see Methods): single nucleus  
110 segmentation (*SIMS-Cut*, Fig. 2a), single nucleus representation (*SIMS-ID*, Fig. 2b) and differential me-  
111 tabolite analysis (*SIMS-Diff*).

112

### 113 ***SEAM* can resolve metabolomic profiles at single cell resolution on various tissues with different** 114 **cell densities**

115 To demonstrate the universality and as a sanity check, we tested *SEAM* using mouse liver (Fig. 1a bottom  
116 row), lung, kidney, and small intestine samples (Fig. 1b). Qualitative visualization of *SIMS-View* may illus-  
117 trate the corresponding tissue structures: e.g. in the liver, the metabolites show gradual changes spread-  
118 ing out from the central vein (CV)<sup>22</sup>; in the lung and kidney, the specific structure of the local metabolic  
119 niches, such as bronchioles and glomerulus<sup>23</sup>; and in the small intestine, the characteristic anatomic pat-  
120 tern along the intestinal villus axis<sup>24</sup> (Supplementary Figs. 1-3).

121 In addition to the spectral projection by UMAP in *SIMS-View*, one can selectively add more histological  
122 or functional information back by using those different *SEAM* modules through quantitatively characteriz-  
123 ing the spatial and compositional information within the single nuclear metabolome. Compared with the  
124 *SIMS-View*, clustering results using the single nuclear representation module *SIMS-ID* can mark strong  
125 correspondence to the well-established cell types, for example, hepatocytes and endothelial cells in the  
126 liver, Clara cells in the lung, as well as enterocytes and lamina propria in the small intestine (Supplemen-  
127 tary Figs. 1-3).

### 128 **Algorithms design and modular data analysis for *SEAM***

129 *SIMS-View* is a fast visualization tool designed for SIMS data, which takes advantage of the efficiency as  
130 well as the local and global structure preservation of UMAP<sup>25</sup>. It takes multiplex SIMS data as input and  
131 outputs a single human-readable image using three steps. First, SIMS data is regarded as 256×256 in-  
132 dependent pixels, each represented by a fixed-length vector, and each pixel is feature-wise normalized to  
133 avoid feature bias. Next, the 65536 pixels are fed into UMAP to reduce the dimensionality to 3. Finally,  
134 each of the three resulting dimensions is scaled and color-coded by CIELAB color space, and all pixels  
135 are mapped back to their original positions. *SIMS-View* provides a global view of all the ion distribution  
136 features in one single image at the pixel level.

137 To solve the cell segmentation problem, various *in situ* works used different methods. Some used  
138 matched H&E stain<sup>17</sup>, others took one simple measurement as input<sup>18,26</sup>. And most of them used super-

139 vised segmentation, via either pixel-wise classification or modeling the whole image using CNN. Interest-  
140 ingly, based on the visualization of *SIMS-View* results on different samples, the nuclei of cells showed  
141 similar color for most cells yet different from other non-nuclear areas (Fig. 1a, b). Therefore, we decided  
142 to isolate the nucleus to demarcating every single cell. To avoid extra staining and heavy annotation labor  
143 which would sabotage the original metabolic state of samples, we developed *SIMS-Cut*, an unsupervised  
144 label-free algorithm, to segment regions of interest (ROIs) using corresponding metabolic markers, for  
145 example, adenine ( $m/z$  134) as the nuclear marker<sup>16</sup>. The input data format is multiplex by selecting those  
146 ion species highly co-localized with nuclei, which is highly consistent across different samples (Supple-  
147 mentary Fig. 4a). And the core of *SIMS-Cut* is an expectation-maximization (EM) algorithm, aiming to  
148 solve an optimization problem of a probabilistic graphical model (PGM)<sup>27</sup> which combines a restricted  
149 Boltzmann machine (RBM)<sup>28-31</sup>, and a Potts model<sup>32,33</sup> (Fig. 2a, Supplementary Fig. 4d). The RBM (Sup-  
150 plementary Fig. 4c) is suitable for modeling the appearance of a multi-image pixel given its label (fore-  
151 ground/background), and the Potts model (Supplementary Fig. 4b) encourages the resulting segmenta-  
152 tion masks to be smooth.

153 To demonstrate the superior performance of *SIMS-Cut*, we compared with several popular unsuper-  
154 vised segmentation algorithms (Supplementary Fig. 5c), using different cell cultures with adenine ( $m/z$   
155 134) as the ground truth (Supplementary Fig. 5a, and see Methods). The results showed that *SIMS-Cut*  
156 could consistently outperform contestants in all cases visually and quantitatively (Supplementary Figs.  
157 5a-c, 6). To test the suitability on tissue samples, we also applied *SIMS-Cut* on various wild type mouse  
158 tissues, ranging from lung, kidney, small intestine, and liver (Supplementary Figs. 7, 8). For the more  
159 challenging case, where cells might display distinct sizes and densities, *SIMS-Cut* was finally applied on  
160 human liver fibrosis tissues from multiple patients, and all resulted in consistent and satisfactory perfor-  
161 mance (Supplementary Figs. 9, 10).

162 After segmentation, the metabolic fingerprint of each segmented nucleus needs to be extracted and  
163 represented. Given the fact that SIMS captures the cumulative intensities along the z-axis for each pixel,  
164 extracting the metabolic fingerprint of each cell (both nucleus and cytoplasm) can be done by combining  
165 its segmentation mask and corresponding SIMS data. Existing works often represented cells by compu-  
166 ting the average of all the pixels containing within each cell<sup>19,20</sup>, which required strong assumptions like  
167 Gaussian or unimodal, and suffered from loss of pixel variation (Supplementary Fig. 11c). To obtain better  
168 results, *SIMS-ID* represents cells using the bipartite graph of pixels and cells constructed by a self-super-  
169 vised learning algorithm<sup>34-36</sup>, which can soften the hard labeling produced by *SIMS-Cut* (Fig. 2b and Sup-  
170 plementary Fig. 11a, b). The resulting representation showed superior discriminative power, noise robust-  
171 ness, and pixel distribution preservation.

172 To test the distinguishing features mentioned above, we constructed 11 datasets (See Methods) con-  
173 taining both mixed cell populations simulated based on single cell line cultures (Supplementary Figs. 5a,  
174 12), and mixed-cultured cells (Supplementary fig. 13). To compare the discriminative power between  
175 *SIMS-ID* and the conventional mean representation, we tested supervised classification using KNN  
176 equipped with cross-validation and unsupervised clustering using several standard algorithms (K-Means<sup>37</sup>,  
177 SC3<sup>38</sup>, SIMLR<sup>39</sup>, T-SNE<sup>40</sup> followed by K-Means, and UMAP followed by HDBSCAN<sup>41</sup>), then applied them

178 on both representation methods to compare on datasets 4,5,6,7, each containing 4 cell clusters (Supple-  
179 mentary Fig. 12a), whose ground truth is naturally derived *in silico*; and on datasets 10,11, two mixed-  
180 cultured datasets, whose ground truth is provided by BrdU/IdU labeling<sup>42</sup> (Supplementary Fig. 13, and  
181 see Methods) without affecting on cell metabolic fingerprint (Supplementary Fig. 21). The results showed  
182 superior performance of *SIMS-ID* in both supervised (Supplementary Fig. 17) and unsupervised (Supple-  
183 mentary Fig. 18) cases, even in cases of minor fold changes on two feature dimensions (Supplementary  
184 Figs. 12a, 17). To evaluate the sensitivity of capturing pixel distribution of cells, we first tested *SIMS-ID*  
185 with dataset 3, where it could identify the change of the pixel distribution from the original data to Gaussian  
186 (Supplementary Fig. 16), then on dataset 8 and 9 (Supplementary Fig. 12b,c), where *SIMS-ID* could dis-  
187 tinguish cell types with unimodal and multimodal distributions (Supplementary Fig. 19), or different joint  
188 distributions even on two feature dimensions (Supplementary Fig. 20). To test the robustness to inaccur-  
189 ate segmentation and pixel-wise multiplicative noise, *SIMS-ID* was applied on dataset 1 and 2, and  
190 showed consistently better performance than the mean representation (Supplementary Figs. 14, 15). The  
191 *SEAM* analyses of datasets 10,11 are shown in Supplementary Figs. 22, 23.

192 The resulting representation of *SIMS-ID* lies in high dimensional feature space. *SIMLR*<sup>39</sup> is a popular  
193 single cell clustering algorithm, which automatically learns cell to cell affinity with multiple kernel ensemble  
194 learning, and shows satisfactory performances when combined with *SIMS-ID* (Supplementary Fig. 18).  
195 We simply adopted *SIMLR* as our clustering method.

196 To characterize the key metabolites differentiating clusters, and account for the variation of pixels within  
197 cells, we developed *SIMS-Diff* as our differential analysis algorithm. *SIMS-Diff* regards cells as distribu-  
198 tions of pixels and uses earth mover's distance (EMD, see Methods)<sup>43</sup> as the dissimilarities among cells.  
199 Using this, the discriminative power of one feature with respect to a given cluster partition can be meas-  
200 ured as the ratio of between cluster variation (BCV) and within cluster variation (WCV).

201

## 202 ***SEAM* reveals cell spatial metabolic states in wild type mouse liver.**

203 Liver is an important metabolic organ consisting of repeating hexagonal-shaped units called lobules<sup>44</sup>.  
204 Spatial heterogeneity of metabolic mechanism has been thoroughly investigated using immunohistochem-  
205 istry (IHC) analyses<sup>45</sup>, transcriptome<sup>22</sup>, and epigenome<sup>46</sup>, but, to our knowledge, single cell level of direct  
206 spatial metabolome has not been reported. This allows us to fill up the gap by a proof-of-concept demon-  
207 stration of *SEAM*.

208 To this end, wild type mice were used to obtain sequential liver sections, and CV centered regions were  
209 selected for *SEAM* analysis. The SIMS data consists of approximately 200~300 ion species after spectral  
210 peak selection and filtering (See Methods), and *SIMS-Cut* detected 724 nuclei in the square. To extract  
211 metabolic cell fingerprint, we used *SIMS-ID* to represent each cell with a fixed-length vector, which was  
212 fed into *SIMLR* to obtain metabolically distinct cell subpopulations. *SIMLR* reached an optimal  $k = 8$ , and the  
213 resulting 8 metabolically distinct subpopulations correspond to major liver cell types, including Kupffer  
214 cells, 2 subpopulations of endothelial cells, and 4 subpopulations of hepatocytes (Fig. 2c).



215 The identified subpopulations showed specific spatial patterns consistent with the known liver organi-  
216 zation (Fig. 2c). Kupffer cells are specialized macrophages in the liver, which typically line on the walls of  
217 the sinusoids. Endothelial cells correspond to vascular endothelial cells and liver sinusoidal endothelial  
218 cells, typically lying between the crevices of hepatocytes and receiving blood from both the hepatic artery  
219 and the portal veins into the hepatic parenchyma<sup>47</sup>. Hepatocytes (the parenchymal cells) constitute 80%  
220 of the mass and 60% of cell composition in a healthy mammalian liver, performing various metabolic  
221 functions strongly associated with their positions<sup>44</sup>. *SIMS-Diff* identified differential ion species among the  
222 subpopulations (Fig. 3a, b). We found *m/z* 60, 76, and 77 as metabolic markers of endothelial cells, while  
223 *m/z* 134, 181, and 91 enriched in Kupffer cells (*m/z* 134 is reported to be adenine, reflecting the higher  
224 nucleus-to-cytoplasm ratio). Hepatocytes, which differ from liver non-parenchymal cells, were character-  
225 ized by *m/z* 255, 279, and 281, corresponding to the fatty acid metabolism of parenchymal tissue. Inter-  
226 estingly, hepatocyte may be sub-classified by C1, C2, C3, and C4 each showing different metabolic fin-  
227 gerprints (Fig. 3a, b).

228

#### 229 **Hepatocyte metabolic clusters show a consistent but complementary spatial pattern with liver** 230 **zonation**

231 Having identified the metabolic heterogeneity among hepatocytes in wild type mouse liver lobule, we  
232 searched for differential gene expression corroboration in the literature. Hepatocyte C1 was visually lo-  
233 calized around CV, and quantitative analysis revealed that the cells in Hepatocyte C1 showed significantly  
234 smaller distances from CV compared with the other hepatocytes ( $P < 10^{-9}$ , one-side Wilcoxon rank sum  
235 test) (Fig. 3d). We also found 6 ion species markers and observed the gradual changes along the liver  
236 lobule (Fig. 3c), as well as the zonation pattern of each representative metabolite in single cell level,  
237 showing consistent pattern with reported spatial transcriptome<sup>22</sup> (Fig. 3e). Additionally, replicate experi-  
238 ments on different CV regions also showed consistent metabolic patterns and cluster-specific metabolites,  
239 indicating the robustness and effectiveness of our method (Supplementary Figs. 24a, 25a-f). We reported  
240 *SEAM* results of the liver portal node (PN) as our negative control (Supplementary Fig. 25g, h). Consistent  
241 with the spatial expression of *GLUL*<sup>22</sup>, the spatial pattern of *m/z* 58, 59, 69, 71, 87, and 101 showed higher  
242 expression in the nearest 1~2 layers of hepatocytes from CV (Fig. 3c, e). We further conducted the IHC  
243 of two liver zonation markers, Glutamine synthetase (GS), the protein encoded by *GLUL*, and Cytochrome  
244 P450 2E1 (Cyp2e1), at the adjacent slides and confirmed liver zonation pattern (Supplementary Fig. 24b-  
245 d). This example provided *SEAM* with a positive control that it can accurately and comprehensively char-  
246 acterize the spatial heterogeneity within a well-studied tissue microenvironment.

247

#### 248 **SEAM identified metabolically different hepatocyte subpopulations associated with the fibrotic** 249 **niche.**

250 Liver cirrhosis has been a major killer, and progressive liver fibrosis often results in liver cirrhosis<sup>48</sup>. Having

251 been proven effective in the case of wild type mouse liver, *SEAM* was applied to human liver fibrosis to  
252 characterize the metabolic microenvironment around a fibrotic niche. We hypothesized that there should  
253 be metabolic alterations of hepatocytes around the fibrotic niche, and such alteration might be associated  
254 with the distance between hepatocytes and fibrotic boundary (FBD) at a local scale.

255 To test this hypothesis, we collected 10 non-tumor tissue regions from 3 liver cancer patients (Supple-  
256 mentary Table. 2) and made a sequential 10 $\mu$ m slides for SIMS and other assays. We selected 4 regions  
257 from one sample, each containing a fibrotic niche, and conducted SIMS experiments (Fig. 4a, b). The  
258 resulting data consists of approximately 200~300 ion species after spectral peak selection and filtering  
259 (See Methods). The color-coded pixel visualizations produced by *SIMS-View* depicted a qualitative spatial  
260 pattern within each region (Fig. 4c left column). To quantitatively characterize the cell composition and  
261 spatial organization, *SIMS-Cut* detected 902, 716, 546, and 682 nuclei in four square regions respectively.  
262 *SIMS-ID* and *SIMLR* were subsequently performed to get metabolically distinct cell subpopulations. The  
263 consistent manifolds and clusters shown by UMAP (Fig. 4c middle column) and the spatial single nucleus  
264 map (Fig. 4c right column) confirmed the reliability and robustness of *SEAM*. The identified subpopulations,  
265 corresponding to Kupffer cells, immune cells, fibroblasts, endothelial cells, and 3 subpopulations of  
266 hepatocytes, exhibited the specific spatial distributions (Fig. 4d) and the matching metabolic fingerprints  
267 (Fig. 4e). The correspondence and incongruity between cell subpopulations of human and mouse liver  
268 samples were also analyzed (See Methods and Supplementary Fig. 26).

269 Intriguingly, we observed that Hepatocyte C1 was visually localized near the FBD, and its associated  
270 metabolic markers, e.g. *m/z* 69, 55, and 57, showed the consistent spatial pattern across 10 regions (Fig.  
271 4f, g, and Supplementary Figs. 27, 28). To quantify the association between the hepatocyte metabolic  
272 alteration and the distance to the FBD, we separately conducted two statistical analyses on 10 regions of  
273 3 patients (Supplementary Table. 2) given defined FBD (see Methods and Supplementary Fig. 30 second  
274 column): the distance from FBD to hepatocyte C1/C2 (distance-based analysis), and the normalized count  
275 ratio between hepatocyte C1 and C2 (count-based analysis). Using R1 as a demonstration, we first de-  
276 fined 5 zones (zone 0~4) with increasing areas (Fig. 4h left), each representing an accumulative territory  
277 between the FBD and the corresponding parallel strip (parallel strips are indicated by gray solid lines, and  
278 the accumulative territories of zones are indicated by gray dotted brackets), then the distances from FBD  
279 to Hepatocyte C1/C2 within the 5 zones were subsequently summarized by a series of paired boxplots  
280 (Fig. 4h right, n=10). Meanwhile, we calculated the normalized count ratio between Hepatocyte C1 and  
281 C2 within an area as a function of the distance from the outer edge (indicated by the gray solid line in Fig.  
282 4i left) to the FBD (Fig. 4i right, n=10). The result of the distance-based analysis showed that Hepatocyte  
283 C1 was significantly closer to FBD than C2 to FBD within the 5 zones (one-side Wilcoxon rank sum test,  
284 Fig. 4h right, n=10), and the relative proximity exhibited high similarity across 10 regions (Supplementary  
285 Fig. 30 third column). Complementarily, the count-based analysis showed that the normalized count of C1  
286 is consistently higher than C2, specifically, C1 was about ~30-50% denser than C2 within 100 $\mu$ m (a typical  
287 hepatocyte size is ~25 $\mu$ m) to the FBD and reduced quickly to about the same level as C2 after ~350 $\mu$ m  
288 (Fig. 4i right, n=10), and this trend was highly similar across 10 regions (Supplementary Fig. 30 fourth  
289 column). Detail of FBD determination, zone partition, distance, and normalized count ratio calculation, as  
290 well as other necessary terms definition, is exactly described in Methods. The above statistical analyses  
291 verified our hypothesis that the metabolic alteration of the hepatocyte subpopulations might be associated  
292 with the spatial proximity to the fibrotic niche. To verify the variation of microenvironment was not only  
293 reflected at the metabolic level, we subsequently performed Geo-seq, a spatial transcriptome assay at  
294 the same ROIs of different hepatocyte subpopulations.

295

296 **Spatial transcriptome validated metabolism associated gene expression alteration in heterogene-**  
297 **ous hepatocytes identified by SEAM**

298 To get a deeper understanding of *SEAM* results, we performed Geo-seq with a modified protocol (See  
299 Methods) of the transcribed RNA samples isolated from the tissues of the corresponding ROIs from the  
300 adjacent slides (Fig. 5a, b, and Supplementary Fig. 31). To increase reproducibility, multiple adjacent  
301 slides were used (Supplementary Fig.32-35). The Geo-seq slides showed high continuity with the corre-  
302 sponding SIMS slides in terms of spatial histology (Fig. 5b). Hepatocyte C1 from *SEAM*'s result, which  
303 was proximal to fibrotic niche and enriched with ions species *m/z* 69 series were defined as Hepa<sup>69-high</sup>,  
304 whereas Hepatocyte C2, which were distal and not enriched with ions species *m/z* 69 series were defined  
305 as Hepa<sup>69-low</sup>. We also collected the fibrotic regions as the FB samples. In total, 15 cDNA libraries were  
306 constructed successfully (Hepa<sup>69-high</sup> n=6, Hepa<sup>69-low</sup> n=5, and FB n=4). Principle component analysis  
307 (PCA) plot indicated that two different groups (Hepa<sup>69-high</sup> -proximal and Hepa<sup>69-low</sup> -distal) of hepatocytes  
308 shared higher similarity relative to FB samples (Fig. 5c). More importantly, Hepa<sup>69-high</sup> samples were con-  
309 sistent closer to FB samples than Hepa<sup>69-low</sup> to FB samples in PCA space (Fig. 5c and Supplementary  
310 Fig. 36). To validate the expression pattern of each group, we first compared gene expression profiles  
311 between hepatocytes (i.e. Hepa<sup>69-high</sup>/ Hepa<sup>69-low</sup>) and FB, then performed gene ontology (GO) enrichment  
312 for both up-regulated and down-regulated differentially expressed genes (DEGs) (See Methods and Sup-  
313 plementary Fig. 37, 38). Up-regulated DEGs were mainly involved in liver biosynthesis pathways for both  
314 Hepa<sup>69-high</sup> and Hepa<sup>69-low</sup> groups and down-regulated DEGs were highly enriched in lymphocyte activa-  
315 tion and humoral immune response pathways. We further looked at the well-known marker genes specific  
316 for hepatocytes (*ASL*, *HP* & *SAA1*), fibrosis (*TGFB1*, *PDGFB* & *COL4A1*), and immune response (*IGHM*,  
317 *IGHG3* & *IGHV4-59*). Both hepatocytes groups showed high levels of hepatocyte marker genes. Whereas  
318 genes typically activated in fibrotic regions for fibrosis and immune response were highly expressed in FB  
319 samples (Supplementary Fig. 39). There were 718 differentially expressed genes (DEGs) fitting into the  
320 criteria of adjust P-value < 0.05 and log fold change (LFC) standard error < 3. The expression heatmap  
321 indicated that these genes had different expression patterns between the proximal hepatocytes (Hepa<sup>69-</sup>  
322 high) and the distal (Hepa<sup>69-low</sup>) (Fig. 5d). We inputted the DEGs for GO enrichment analysis (Fig. 5e).  
323 There were 17 genes enriched in the first GO entry, 16 of them were consistently higher in Hepa<sup>69-high</sup>  
324 than Hepa<sup>69-low</sup> (Fig. 5f). Genes of solute carrier transporters families with different functions were enriched  
325 in the fibrosis proximal (Hepa<sup>69-high</sup>) group, indicating the corresponding metabolite transmembrane ex-  
326 change activities were elevated.

327

328 **Discussion**

329 In this study, we have developed *SEAM*, a platform combining experiments, and computational algorithms  
330 to quantitatively characterize the metabolic intra- or inter-cellular features with multiscale spatial resolution.  
331 Unlike other IMS instruments such as DESI (40–60µm)<sup>11</sup>, SIMS can provide a high spatial resolution

332 allowing one to visualize detailed metabolic structures in tissue histology. With fast and minimal sample  
333 processing, SIMS maximumly preserves the native state of samples. Given the nature of SIMS, although  
334 it breaks most of the molecules into fragments, making it more difficult to annotate (a common challenging  
335 issue for MS studies), it produces high multiplexity of metabolic features with the potential of characteriz-  
336 ing cell and fine tissue microenvironment. Benefiting from both high spatial resolution and high multiplexity  
337 of SIMS, the algorithms of *SEAM* start solely from the features generated by SIMS and run a pipeline  
338 enabling metabolic analysis from pixels to single nuclei, then to the selected metabolic molecules with  
339 spatial information annotated. Previously, there have been reports on spatial metabolic features at tissue  
340 level or *in vitro* single cell level<sup>16</sup>. But, to our knowledge, this is the first study capable of segmenting and  
341 analyzing single nuclear metabolic profiles directly on tissue sections. In addition, this algorithmic pipeline  
342 is principally scalable to other spatial omics studies based on other IMS platforms, transcriptomics, and  
343 proteomics with minimum adjustments, and it's also easy to work together with bioinformatics tools such  
344 as CIPHER to predict and prioritize disease-related metabolic molecules<sup>49</sup>.

345         Apart from the scalability of *SEAM*'s algorithms, we have demonstrated that the range of *SEAM*  
346 applications could cover from *in vitro* cell culture assays to various tissue samples. Firstly, in the mixed  
347 cell-cultured assay, *SEAM* could easily deconvolute the different cell lines co-cultured together. Addition-  
348 ally, in different wild-type murine tissue samples, *SEAM* successfully segmented single nuclei without  
349 extra labeling required. The single nuclear metabolic profile analysis was also consistent with conventional  
350 tissue histological characterization (Supplementary Figs. 1-3). Specifically, in the liver, a spatially well-  
351 orchestrated but complex organ, the CV-PN axis zonation has been well-established at single cell tran-  
352 scriptome level in wild type mouse<sup>22</sup>. We observed consistent zonation patterns at single cell level in CV  
353 centered region with the gradational decrease of certain characteristic metabolites. Lastly, we found that  
354 hepatocyte subpopulations (among which, to our knowledge, the novel C1 has never been reported before)  
355 differentiated by different metabolic features were also transcriptionally distinct shown by Geo-seq (Fig.  
356 5c-f). The elevated expression level of solute carrier genes can potentially explain the enrichment of a list  
357 of metabolite species found by *SEAM* (Fig. 4). These genes are involved amino acid transport  
358 (*SLC36A4*, *SLC3A2* & *SLC38A9*)<sup>50-52</sup>, phosphate transport (*SLC17A2* & *SLC17A4*)<sup>53</sup> and Gamma-Ami-  
359 nobutyric Acid (GABA) transport (*SLC6A12*)<sup>54</sup>. *SLC3A2* has already been reported to play a central role  
360 in fibronectin matrix assembly, which also concurs with our result as the proximal samples were more  
361 close to the fibrotic region<sup>51</sup>. It indicates that spatial microenvironment differences could influence cellular  
362 metabolic homeostasis, which may in turn further alter the gene regulation and downstream response due  
363 to cell adaptation and genetic/epigenetic feedback.

364         In summary, *SEAM* provides a high spatial resolution single nuclear metabolome profiling pipeline  
365 requiring minimal sample preparation and labeling. It is automatically scalable to different biological sam-  
366 ples ranging from cell culture assays to complex tissue samples. It can have a great impact on differenti-  
367 ating subtle tissue metabolic changes undetectable for or complementary to other conventional assays.  
368 With future improvement of IMS resolution and molecule annotation capability, *SEAM* would be able to  
369 provide more detailed spatial metabolome profiles with higher resolution and broader functionality.

370 **ONLINE METHODS**

371 **IMS experiments**

372 TOF-SIMS 5 (ION-TOF GmbH, Münster, Germany) equipped with a Bi liquid metal ion gun (LMIG)  
373 is used in this study, collected TOF-SIMS spectra and images of tissue samples using a 30 keV  
374 Bi<sub>3</sub><sup>+</sup> LMIG with a high spatial resolution (HSR) mode. The Bi<sub>3</sub><sup>+</sup> current in the HSR mode was 0.1  
375 pA (100 ns pulse width, unbunched beam). The total Bi<sub>3</sub><sup>+</sup> accumulated ion dose was about 2.0 ×  
376 10<sup>10</sup> ions/cm<sup>2</sup>, the typical probe sizes of the Bi<sub>3</sub><sup>+</sup> LMIG was ~200 nm in HSR mode. The secondary  
377 ion images were acquired using Bi<sub>3</sub><sup>+</sup> LMIG rastering over a 400 × 400 μm<sup>2</sup> area with 256 × 256  
378 pixels. The Bi<sub>3</sub><sup>+</sup> LMIG was operated at a cycle time of 150 μs (mass range: 0 ~2000 u). Negative  
379 spectra were mass-calibrated using CH<sub>2</sub><sup>-</sup>, O<sup>-</sup>, OH<sup>-</sup>, PO<sub>2</sub><sup>-</sup>. A flood gun with low energy electrons was  
380 used to compensate for charge buildup on sample surface. A 10-keV Ar<sub>2500</sub><sup>+</sup> commercial gas cluster  
381 ion gun (GCIB) was used as a sputter gun (rastering over a 550 × 550 μm<sup>2</sup> area, incident angle  
382 45°) to carry out the depth profiling. A final 2D image was an overlay of 80~120 layers of depth  
383 profiling scan images.

384 In initial cell analysis, a high mass resolution (HMR) mode was used with 0.8 pA (<1 ns pulse width,  
385 bunched beam) Bi<sub>3</sub><sup>+</sup> current, the mass resolutions (measured at C<sub>2</sub>H<sup>+</sup>) were typically >6000. The  
386 total Bi<sub>3</sub><sup>+</sup> LMIG accumulated ion dose was between 10<sup>11</sup> and 10<sup>12</sup> ions/cm<sup>2</sup>, rastering over a 300  
387 × 300 μm<sup>2</sup> area with 256 × 256 pixels. The Bi<sub>3</sub><sup>+</sup> LMIG was operated at a cycle time 150 μs (mass  
388 range: 0 ~2000 u). Negative spectra were mass-calibrated using CH<sub>2</sub><sup>-</sup>, O<sup>-</sup>, OH<sup>-</sup>, PO<sub>2</sub><sup>-</sup>. A flood gun  
389 with low energy electrons was used to compensate for charge buildup on sample surface. A 10-  
390 keV Ar<sub>2500</sub><sup>+</sup> commercial gas cluster ion gun (GCIB) was used as a sputter gun (rastering over a  
391 450 × 450 μm<sup>2</sup> area, incident angle 45°) to carry out the depth profiling. A final 2D image was an  
392 overlay of 50-80 layers of depth profiling scan images.

393 **Peak selection.** To avoid noise interference and improve follow-up analysis efficiency and accu-  
394 racy, picking out peaks from a full spectrum was necessary. A Peak Search process in SurfaceLab  
395 was carried out with the parameters as bellow: mass range 50-500; minimum counts 10000; min-  
396 imum signal/noise ratio 1000. Typically, 200-500 peaks were picked out from a full spectrum.

397 **SIMS data preprocessing.** Each peak corresponds to a highly spatially resolved and spectrally  
398 filtered ion image: the former originated from a specific one or a class of chemical substances in  
399 the tissue sample while the latter shows its characteristic spatial distribution features in this tissue  
400 square (Fig. 1a, top right). For further data analysis, each ion image can be exported as an Amer-  
401 ican Standard Code for Information Interchange (ASCII) mode data file by the SIMS built-in data  
402 processing software SurfaceLab, which contains three columns corresponding to the X-axis, Y-  
403 axis coordinates and signal intensity values.

404

405 **Biological experiments**

406 **Cell culture.** Human non-small cell lung cancer cell line A549, human cervix carcinoma cell line  
407 HeLa, murine hepatoma cell line Hepa 1-6 and murine liver epithelial cell line NCTC 1469 cell lines  
408 were grown on microscope cover glass (CITOGLAS, China) with Dulbecco's Modified Eagle Me-  
409 dium (DMEM) (Gibco, USA) containing high glucose, L-glutamine, sodium pyruvate and 10% dia-  
410 lyzed, heat-inactivated FBS (Gibco, USA). Human mammary gland cell line MCF 10A was grown  
411 on microscope cover glass (CITOGLAS, China) with DMEM/F12 (1:1) (Gibco, USA) containing  
412 insulin 10ug/ml, EGF 20ng/ml, cholera toxin 100ng/ml, hydrocortisone 0.5mg/ml and 5% equine  
413 serum. Human breast adenocarcinoma cell line MDA-MB-468 cell line was grown on microscope  
414 cover glass (CITOGLAS, China) with L-15 medium containing 10% FBS (Gibco, USA) and free air  
415 exchange.

416 **BrdU cell mix-culture experiment.** Following protocol from the previous study, A549 and HeLa  
417 cell lines were both cultured with and without 20 $\mu$ M BrdU (Sigma, USA) for 48 hours before seed-  
418 ing. A549 with BrdU were then replated with non-BrdU HeLa at the same density on microscope  
419 cover glass (CITOGLAS, China) for 20 hours and vice versa for non-BrdU A549 and HeLa with  
420 BrdU. The same mix-culture procedure for IdU (Sigma, USA) was applied at Hepa 1-6 and NCTC  
421 1469 cell lines.

422 **Mice.** C57BL/6N mice were purchased from Charles River. All mice were housed in isolated ven-  
423 tilated cages (maxima six mice per cage) barrier facility at Tsinghua University. The mice were  
424 maintained on a 12/12-hour light/dark cycle, 22-26°C with sterile pellet food and water ad libitum.

425 The laboratory animal facility has been accredited by AAALAC (Association for Assessment and  
426 Accreditation of Laboratory Animal Care International) and the IACUC (Institutional Animal Care  
427 and Use Committee) of Tsinghua University approved all animal protocols used in this study (Ani-  
428 mal Welfare Assurance Number F16-00228 (A5061-01)).

429 **Intrahepatic cholangiocarcinoma (ICC) patient non-tumor liver tissues.** The ICC non-tumor  
430 liver tissues were obtained from leftover pieces from surgery. The protocol of this study was com-  
431 pliant with the principles of the Declaration of Helsinki and was also approved by the Institutional  
432 Review Board (IRB) and Ethics Committee (EC) of Peking Union Medical College Hospital  
433 (PUMCH) (JS-2492).

434 **Tissue section preparation.** Mouse and human tissues were isolated individually and embedded  
435 in Optimum Cutting Temperature (O.C.T) compound (SAKURA, USA), then snap-frozen in liquid  
436 nitrogen. Cryo-section were performed using CM1900 Cryostat (Leica, Germany) to obtain 3 $\mu$ m ~  
437 10 $\mu$ m continuously adjacent sections.

438 **Histology staining.** Tissue cryo-sections were thawed at room temperature for 5 min then washed  
439 in PBS twice, 5min each time. Slides were fixed in 4% paraformaldehyde (PFA) for 20 min at room  
440 temperature then washed in PBS once. H&E stainings were then performed using the H&E staining  
441 kit (Leagene, China). Images were obtained from Axio Scan. Z1 (ZEISS, Germany) or Cytation5

442 (Biotek, USA).

443 **Immunohistochemistry.** Tissue cryo-sections were thawed at room temperature for 5 min then  
444 washed in PBS twice, 5min each time. Samples were permeabilized and blocked in 5% BSA solu-  
445 tion (Sigma, USA) with 0.4% Triton-X100 (AMRESCO, USA) for 2h at room temperature. Dilute  
446 and apply primary antibody in PBS with 0.1% Triton-X100 with suited concentration according to  
447 each antibody and incubate in a humid dark chamber at 4°C overnight. Wash three times in PBS  
448 with 0.1% Triton-X100, 10min each. Dilute and apply secondary antibody in PBS with 0.1% Triton-  
449 X100 and incubate in a humid dark chamber at room temperature for 2h. Wash three times in PBS  
450 with 0.1% Triton-X100, 10min each. Slides were mounted using ProLong™ Gold Antifade  
451 Mountant (ThermoFisher, USA). Images were captured either by LSM780 confocal microscope  
452 (ZEISS, Germany) or Cytation5 (Biotek, USA).

453 **Modified Geo-seq.** A spatial transcriptome analysis method, Geo-seq, previously described by  
454 Chen, Jun, et al<sup>55</sup>. A modified version was adopted. Tissue cryosections were mounted on the PEN  
455 membrane slide and stored at -80 degree freezer for short term storage. Slides were stained in  
456 0.5% cresyl violet and dehydrated in serial ethanol. Tissue blocks were obtained in a 0.2 ml PCR  
457 tube by LMD7000 (Leica, Germany). Buffer RLT (Qiagen, Germany) with DTT (Sigma, USA) were  
458 added and shaken vigorously for tissue lysis and RNA release. RNA Clean beads (Vazyme, China)  
459 1.8x were added to isolate total RNA. Prepare annealing procedure in the same tube with 3ul H<sub>2</sub>O,  
460 1ul dNTP, 1ul Oligo(dT), and 0.5ul RNase Inhibitor (RI) (Life Technologies, USA). Incubated at 72  
461 degrees 3min and immediately transfer in ice for 2min. Prepare reverse transcription reaction in  
462 the same tube with 2ul 5x RT buffer, 0.5ul DTT, 0.5ul RI, 0.5ul Template Switch Oligo (TSO, Sangon  
463 Biotech, China), 1ul Maxima reverse transcriptase (Life Technologies, USA). Incubate with 50 de-  
464 grees with 1 hour and deactivate reverse transcriptase with 85 degrees for 5 min. Amplified the  
465 first strand product with 12.5ul 2x KAPA HIFI HotStart ReadyMix (Sigma, USA), 0.5ul TSO-PCR  
466 primer (Sangon Biotech, China) and 2ul H<sub>2</sub>O. The reaction condition was 95 degrees 3min, 98  
467 degree 20s, 67 degree 15s, 72 degrees 6min for 21 cycles, and 72 degrees for 5min. PCR product  
468 was purified with 0.8x DNA Clean beads (Vazyme, China). The next generation sequencing (NGS)  
469 library was then constructed by TruePrep DNA Library Prep Kit V2 for Illumina (Vazyme, China).  
470 Libraries were sequenced by Illumina Xten Pair-end 150bp by Annoroad.

471

## 472 **RNA-seq data processing and analysis**

473 RNA-seq data were firstly performed with adaptor removal and quality filtering by Trim Galore<sup>56</sup>.  
474 The qualified reads were then mapped to the human gencode reference genome using STAR and  
475 generated BAM files<sup>57,58</sup>. Duplication was removed by PICARD ([http://broadinstitute.github.io/pi-  
476 card/](http://broadinstitute.github.io/picard/)) for all the BAM files. Read count for each gene was performed by HTSeq-count with refer-  
477 ence to gencode human gene annotation, release 32 (GRCh38.p13)<sup>57,59</sup>. Different gene expres-  
478 sion analyses were analyzed using DESeq2 in R<sup>60</sup>.

479

480 **SIMS-Cut framework**

481 Given an  $M \times N \times N$  SIMS data, with  $M$  filtered metabolic peaks and  $N \times N$  image as input, *SIMS-*  
482 *Cut* first select  $m$  metabolites co-localizing with nucleus (Supplementary Fig. 4a), and then itera-  
483 tively solves a maximum a posteriori (MAP) problem (Supplementary Fig. 4d) to get an  $N \times N$  bi-  
484 nary matrix  $Y$  that indicates a nucleus.

$$Y_{ij} = \begin{cases} 1 & \text{nuclei region} \\ 0 & \text{otherwise} \end{cases} \quad i, j \in [1, N] \quad (1)$$

485 Since the SIMS data is superimposed of a certain thickness of biological slice in its nature, we  
486 regard the segmented nuclei region as a cell containing molecular fragments in both cytoplasm  
487 and nucleus. The main part of *SIMS-Cut* can be formulated as finding an optimal  $Y^*$ :

$$Y^* = \operatorname{argmax}_Y p(Y|X) \quad (2)$$

488 where

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)} \propto P(X|Y)P(Y) \quad (3)$$

489  $X = [x_{ij}]$ ,  $i, j \in [1, N]$ , and  $x_{ij} \in \mathbb{R}^m$ , which is the  $m$  dimensional metabolic density at the coordinate  
490 of  $(i, j)$ . This Bayesian formulation aims to find the optimal label assignment  $Y^*$  that produces the  
491 maximum posterior probability given  $X$ .

492 As with traditional hidden Markov random field (HMRF) based image segmentation<sup>61,62</sup>, *SIMS-Cut*  
493 uses a similar graphical model, consisting of  $P(Y)$ , the smoothing model for unknown label field  $Y$   
494 before guarantee spatial homogeneity, and  $P(X|Y)$ , the data model for the conditional distribution  
495 of pixel metabolic profiles  $X$  given corresponding pixel label.

496 **Smoothing model.** The label prior,  $P(Y)$  is modeled as a special Markov random field (MRF),  
497 called Potts model<sup>32</sup>. According to the Hammersley-Clifford theory<sup>63,64</sup>,  $P(X)$  follows a Gibbs distri-  
498 bution<sup>65</sup>:

$$P(Y) = \frac{1}{Z} \exp(-U(Y)) \quad (4)$$

499 Where  $U$  is called energy function, which is calculated by summing over the potential of all second-  
500 order cliques  $V$ , each clique corresponds to a pair of neighboring pixels(e.g. the 4-neighborhood  
501 system).  $Z$  is called a partition function, making  $P(Y)$  a valid probability density function (pdf).

502

$$U(Y) = \sum_{(i_1, j_1), (i_2, j_2) \in \text{doubletons}} V(y_{i_1, j_1}, y_{i_2, j_2}) \quad (5)$$



503  $V$  is defined on doubleton, penalizing the heterogeneity of labels.

$$V(y_{i_1,j_1}, y_{i_2,j_2}) = \begin{cases} -1, & \text{if } y_{i_1,j_1} = y_{i_2,j_2} \\ +1, & \text{if } y_{i_1,j_1} \neq y_{i_2,j_2} \end{cases} \quad (6)$$

504 **Data model.** According to the graphical model (Supplementary Fig. 4b), and d-separate<sup>27</sup>,

$$P(X|Y) = \prod_{i,j \in [1,N]} P(x_{ij}|y_{ij}) \quad (7)$$

505 While the multivariate Gaussian distribution is typically suited for the data model of color image  
506 segmentation<sup>66,67</sup>, its model capacity is limited and its assumptions are too strong for SIMS data.  
507 Instead, we use Restricted Boltzmann Machines (RBM)<sup>28-31</sup> to model the conditional distribution of  
508 data intensities given label assignment.

509 RBM as a generative model is typically a two-layer bipartite undirected graph. It's composed of a  
510 visible layer which is  $m$  dimensional metabolic profile in our case and a hidden layer which is a  
511 kind of  $d$  dimensional memory providing model capacity. In theory, RBM is a Universal approxima-  
512 tion for any pdf with a large enough number of hidden layers<sup>30</sup>. Here we use two separate RBMs  
513 to model  $P(x_{ij}|y_{ij} = 0)$  and  $P(x_{ij}|y_{ij} = 1)$  respectively, and we describe one RBM in the following.  
514 For the sake of notation simplicity, in the following, we use  $V = [v_p], p \in [1, m]$  to denote  $x_{ij}$  (the  
515 subscript is removable thanks to the conditional independence given by (7)).

516 The graphical model of RBM is shown in Supplementary Fig. 4c.  $H = [h_q], q \in [1, d]$  is the hidden  
517 layer variable, and  $V$  is the visible layer variable.  $C = [c_q], q \in [1, d]$ ,  $B = [b_p], p \in [1, m]$ , and  $W =$   
518  $[w_{pq}], p \in [1, m], q \in [1, d]$  are parameters. The joint probability density function is:

$$P(V, H) = \frac{1}{Z} e^{-E(V, H)} \quad (8)$$

519 where  $E$  is the energy function:

$$E(V, H) = - \sum_{p=1}^m \sum_{q=1}^d w_{pq} h_q v_p - \sum_{p=1}^m b_p v_p - \sum_{q=1}^d c_q h_q \quad (9)$$

520 and  $Z$  is the partition function:

$$Z = \sum_{V, H} e^{-E(V, H)} \quad (10)$$

521 The probability that an RBM model assigns a vector  $V$ , e.g.  $x_{ij}$  is given by (8).

$$\begin{aligned} p(x_{ij}|y_{ij} = a) &= \text{RBM}(V; W^a, C^a, B^a) = \frac{1}{Z^a} \sum_H e^{-E(V, H)} \\ &= \frac{1}{Z^a} \prod_{p=1}^m e^{b_p^a v_p} \prod_{q=1}^d (1 + e^{c_q^a + \sum_{p=1}^m w_{pq}^a v_p}) \end{aligned} \quad (11)$$

522 Note that the superscripts indicates the parameters of specific RBM.

523 **Partition function of RBMs Estimation.** For a specific pixel given its segmentation label  $a$ , the  
 524 log probability that RBM assigns metabolic profiling  $x_{ij}$  is computed as:

$$\log P(x_{ij} | y_{ij} = a) = -F^a(x_{ij}) - \log Z^a \quad (12)$$

525 Here  $F^a(x_{ij})$  is the free energy of RBM corresponding to class  $a$ , which can be rapidly calculated.  
 526 To estimate the partition function  $Z$ , we build a softmax model to classify  $x_{ij}$  at every pixel to its  
 527 label  $y_{ij}$ :

$$\log P(y_{ij} = a | x_{ij}) = \frac{e^{-F^a(x_{ij}) - \log Z^a}}{\sum_{y_{ij}} e^{(-F^{y_{ij}}(x_{ij}) - \log Z^{y_{ij}})}} \quad (13)$$

528 **MAP.** Our objective can be an expression as:

$$\begin{aligned} \operatorname{argmax}_Y \log P(X|Y) + \log P(Y) &= \operatorname{argmax}_Y \sum_{i,j} \log P(x_{ij} | y_{ij}) + \log P(Y) \\ &= \operatorname{argmax}_Y \sum_{i,j} \log \operatorname{RBM}(x_{ij}; W^{y_{ij}}, C^{y_{ij}}, B^{y_{ij}}) + \log P(Y) \end{aligned} \quad (14)$$

529 It's a nonconvex problem, we develop an EM-style algorithm to alternating between two steps to  
 530 reach a locally optimal point iteratively.

531 Each iteration of *SIMS-Cut* consists of three sub-problems, each of which can be solved efficiently.  
 532 The input of each iteration is the segmentation mask output by the previous iteration, and the first  
 533 level's input is simply k-means clustering of an input image. The segmentation mask will converge  
 534 in no more than 20 levels according to our experiments.

535 In the first sub-problem, the parameters of two RBMs are estimated given the label of each pixel  
 536 input from the previous level. Estimated as the parameters of RBMs, solving the partition function  
 537 is time-consuming, thus the second sub-problem bypasses the obstacle and at the same time  
 538 controls the bias of each iteration with the help of a simple binary classification task<sup>68</sup>. And the third  
 539 and last sub-problem uses the well-known graph-cut algorithm<sup>32,33,69,70</sup> to obtain the pixel labels,  
 540 i.e. the segmentation mask for the current iteration. As the process of iterations, the intermediate  
 541 segmentation masks gradually shrinkage, while local homogeneity and nucleus centralization are  
 542 simultaneously kept. Finally, the reaping algorithm is used to salvage as many isolated nuclei as  
 543 possible during the shrinkage process. More details about solving these sub-problems are as fol-  
 544 lows:

545 **Initialization of Y, C, B, W.** The parameters of RBM, e.g.  $C, B, W$  is randomly initialized using a  
 546 Gaussian with zero mean and unit variance. The label assignment  $Y$  is initialized using k-means.

547 **Sub-problem 1: Fix Y to update C,B,W.** This step is Equivalent to learn two independent RBMs.  
 548 Since  $Y$  is given, the training data for the two RBMs can be extracted from  $X$ . An efficient learning  
 549 algorithm, persistent contrastive divergence (PCD)<sup>71,72</sup> can be applied. Also, PCD algorithm is  
 550 based on Maximum likelihood estimation, leading to an increase of objective.

551 **Sub-problem 2: Fix Y,C,B,W to update two partition functions.** Partition function estimation of

552 RBM is time consuming even if all its parameters are known. Based on the efficient way to deal  
553 with the unknown partition functions<sup>29,68</sup>, we build an auxiliary binary classification task and treat  
554 the two partition functions as parameters to estimate. Furthermore, a hyper-parameter beta can  
555 be tuned to control the process of iteration (see Details and online code).

556 **Sub-problem 3: Fix C, B, W, and two partition functions to update Y.** This step is equivalent to  
557 an energy minimization problem, and global optimized Y can be efficiently using a graph cut algo-  
558 rithm.

559 **Reaping.** Using our parameter setting, the above algorithm converges to all-zeros Y within 20  
560 iterations. Because of the spatially different contrast of SIMS image, some nucleus may be lost  
561 during the iteration. We develop an enhancement algorithm to maintain the intermediate identified  
562 nucleus.

563 Due to the bias, as the levels grow high, the region of within-nucleus gets smaller. But the MRF-  
564 based segmentation makes the intermediate segmentation mask of each level homogenous and  
565 evident. To get the final non-connected nucleus mask, a reaping algorithm is proposed in Algorithm  
566 1.

#### 567 **Algorithm 1**

568 *Input:*  $M_k$ : segmentation masks for each level;  $A_u$ : upper bound of nuclei area;  $A_l$ : lower bound  
569 of nuclei area;

570 *Step 1:* Create a queue Q to maintain isolated segments. Create an all-zeros mask  $M_{rst}$

571 *Step 2:* Initialize Q by putting all isolated segments of level 2 to the head of Q; initialize  $M_{rst}$  using  
572  $M_2$ .

573 *Step 3:* pop a segment q from head of Q, set the segment region of  $M_{rst}$  to zeros.

574 *Step 4:* for l from k+1 to K, where q belongs to  $M_k$

575     *if l reaches K*

576         *then set the q region of  $M_{rst}$  to ones*

577         *if two or more segments in  $M_l$  belongs to q*

578             *then push these new segments to tail of Q;*

579             *set these segment region of  $M_{rst}$  to ones;*

580             *break*

581 *Step 5:* return to step 3, until Q is empty

582 *Step 6:* return  $M_{rst}$

583 **Implementation details and parameters setting.** We use correlation distance to select top 20  
584 co-localized ions with Adenine ( $m/z$  134), whose conditional probabilities given labels are modeled

585 by two label-specific RBMs. K-means on a 134 intensity map is used to initialize the segmentation  
586 label, we set  $k=4$  and set clusters with the lowest center as background, other 3 clusters as fore-  
587 ground. For the smoothing model, we use the 4-neighborhood system. For the data model, we use  
588 two Generative RBMs, each with 20 visible nodes and 50 hidden nodes. For RBM training, persis-  
589 tent contrast divergence (PCD) is used for 10 epochs each level.

590 For convenience, we use a Matlab toolbox for RBM modeling and training<sup>73</sup>. When optimizing the  
591 energy minimization problem, we use the Matlab version of the Boykov-Kolmogorov algorithm<sup>69</sup>  
592 provided by <https://vision.cs.uwaterloo.ca/code/>. The original algorithm takes the smoothing model  
593 as a neighbor weights matrix, whose format is described in the code comment, but we modified  
594 the matrix by average filtering with a window size of 21 to provide more smooth quality (optional).  
595 To weight between the data model and the smoothing model, we divide the weights matrix by a  
596 constant (typically 5~10, we use 5 for best practice).

597 To bypasses the time-consuming partition function estimating problem of the two RBMs, a simple  
598 classifier is performed during each iteration. Note that the exact value of the two partition functions  
599 needn't be known<sup>68,74</sup>, the difference matters instead. We first calculate the free energy of all  $N \times N$   
600 pixels separately using the parameters of the two RBMs and sort the difference. Then sort the  
601 difference and take every  $N-1$  interval as classification cutoff. At the same time, one confusion  
602 matrix for each cut off is maintained, so  $N-1$  F measures controlled by beta corresponding to every  
603 interval can be calculated. Finally, the partition function difference with the best F measure is se-  
604 lected. The beta parameter (typically 0.5~1) is tuned to control the convergence process.

605 During the *SIMS-Cut* procedure provided in the methods section, due to the beta parameter, as  
606 the levels grow high, the region of within-nucleus gets smaller. But the MRF-based segmentation  
607 makes the intermediate segmentation mask of each level homogenous and evident. To get the  
608 final non-connected nucleus mask, a reaping algorithm is proposed. The detail is as follows: Sup-  
609 pose after L level's segmentation, *SIMS-Cut* converges to an all-background segmentation mask.  
610 Since each level is an intermediate segmentation mask given beta and upper level's estimated  
611 parameters. The hierarchical structure can be modeled as a tree, whose nodes are nucleus of all  
612 levels, root is a dummy node, the second highest level is the nucleus of first segmentation. Node i  
613 is the child of node j if i belong to the next level of j, and the segmentation region of i is a subset of  
614 segmentation region of j. The leaf nodes are nucleus in the lowest level, the last level of *SIMS-Cut*  
615 procedure. From top to bottom, nodes are split alongside the tree structure, and the reaping algo-  
616 rithm can capture nodes that are optimally split (i.e. according to  $m/z$  134 intensity).

617

## 618 ***SIMS-ID* Framework**

619 After *SIMS-Cut*, hundreds of separated nuclei has been detected from an  $N \times N$  image, each pixel  
620 containing M dimensional metabolic profiles. Thus, each nucleus contains a diverse number of

621 connecting pixels, represented by fixed dimensional vectors. *SIMS-ID* conducted an auxiliary clas-  
622 sification task to assign a single fixed dimensional vector to each nucleus, which is robust to  
623 over/under segmentation in *SIMS-Cut*. The representation learned by *SIMS-ID* compresses all the  
624 pixel metabolic information using a distilled softmax space<sup>75</sup>, regarding a nucleus as a whole while  
625 including distribution information of pixels. A fixed dimensional representation of the nucleus helps  
626 further analysis of single nuclei data analysis, like clustering, visualization, and so on.

627 **Data preprocessing.** Due to the variability of tissue thickness, and variation in ionization and de-  
628 tector efficiency, SIMS data need to be preprocessed. We use Variance-stabilizing normalization<sup>76</sup>,  
629 specifically, the median spectrum is used to estimate the normalization factor, and logarithm was  
630 used as variance-stabilizing transformation.

631 **Motivation.** *SIMS-ID* is based on the observation that the outputs of a trained neural network  
632 contain much richer information than just a one-hot classifier. Hinton, G. et al observe that mutual  
633 similarity between classes can be distilled from a trained softmax based neural network classifier,  
634 e.g. an image of a BMW, may only have a very small chance of being mistaken for a garbage truck,  
635 but that mistake is still many times more probable than mistaking it for a carrot<sup>75</sup>. Lu, Y. applies  
636 factor analysis to reveal the visual similarity of image classes<sup>77</sup>. Wu, Z. utilizes a similar concept  
637 to train an instance-level classifier as an auxiliary task for unsupervised representation learning<sup>78</sup>.

638 **Auxiliary classifier construction.** *SIMS-ID* first constructed a multiple-layer dense neural net-  
639 work armed with a softmax activation at the last layer for classification, then preprocessed pixel  
640 data are input to classify each pixel to the right nuclei, after training, the temperature of softmax  
641 output is raised to a user-set value to soften the probabilistic distribution, and finally the distilled  
642 softmax output of each input pixel can be considered a similarity between the nuclei to which that  
643 pixel belongs and other nuclei, from that pixel's point of view. Further experiments showed that the  
644 overfitting of the auxiliary classifier doesn't hurt the performance of afterward analysis.

645 **Interpretation.** The auxiliary classifier can naturally capture apparent similarity among classes, i.e.  
646 nucleus without being directed to do so. The distilled information, i.e. the high-temperature softmax  
647 output of each pixel can be expressed as a  $P \times C$  matrix PCM, where  $P$  is the number pixels within  
648 all nucleus, and  $C$  is the number of the identified nuclei. The matrix can be interpreted using three  
649 distinct ways.

650 **Nucleus Similarity measure from each pixel's view.** Each row of PCM can be considered as a  
651 similarity measure between the corresponding nuclei and other nuclei. If the  $i$ -th pixel belongs to  
652 the  $j$ -th nuclei, for the  $i$ -th row of PCM, after dividing each element by the  $j$ -th element of the row,  
653 we can get a normalized similarity vector, whose  $j$ -th element is 1. Moreover, in the auxiliary clas-  
654 sification phase, the more easily confused with the correct class, i.e. nuclei, the higher the corre-  
655 sponding element of normalized PCM is.

656 **Nucleus representation of multiple instance learning.** In the multiple instance learning (MIL)  
657 literature<sup>79,80</sup>, a bag of instances can typically be represented by similarities between this bag and  
658 all instances. A column of normalized PCM can be considered as the probability of each pixel

659 belonging to that nucleus.

660 **The adjacency matrix of nucleus-pixel bipartite graph.** The original one-hot pixel-nucleus rela-  
661 tionship doesn't provide any information between nuclei. After knowledge distillation, the one-hot  
662 relationship is shattered to a more smooth knowledge, from which nucleus relationship can be  
663 discovered. The normalized PCM can be interpreted as an adjacency matrix identifying to the bi-  
664 partite graph, and the (i,j)-th entry of PCM is the weight between the i-th pixel and j-th nucleus.

665 **Parameters setting and network structure.** The pixel classification network structure is shown  
666 in Supplementary Fig. 11b. We use multiple layer perceptrons except for the last layer, ReLU<sup>81</sup>  
667 activation function for each layer, softmax as probability output, and Adam<sup>82</sup> as an optimizer. The  
668 number of neurons of the first layer is M, the number of observed metabolites and the number of  
669 neurons of the last layer is the same as the number of the nucleus. Since overfitting doesn't hurt  
670 the representation performance according to our experiment, we set all the pixels as training data,  
671 and the number of training epochs is set to 100~300.

672

## 673 **Clustering**

674 Represented by fixed-length vectors, the nuclei can be straightforwardly clustered and visualized  
675 in low dimensional space. The number of cells that one SIMS experiment captures typically ranges  
676 from 400~1000, and the length of the representation vector for each cell is equal to the number of  
677 pixels within segmented cells, typically ranging from 5000~15000. With the consideration of both  
678 data characteristics and experimental performance (Supplementary Fig. 18), we apply SIMLR<sup>39</sup>, a  
679 single cell clustering algorithm, which automatically learns the low-rank similarity matrix by means  
680 of multiple kernel ensemble. Besides, SIMLR also provides means of estimating the number of  
681 clusters, which we can take as a guideline to explore populations of metabolic cell states in different  
682 scales.

683

## 684 **SIMS-Diff framework**

685 The goal of this algorithm is quantification the feature's discriminative power to tell clusters apart.  
686 Due to the nature of our data, the traditional two-sample test can't be directly applied. We assume  
687 that discriminative features can produce a similarity matrix with a block diagonal structure. There-  
688 fore, we use the ratio between BCV and WCV to evaluate the compactness of the similarity matrix,  
689 where BCV is between cluster variation, and WCV is within cluster variation. For each feature, we  
690 use EMD (earth mover's distance)<sup>43</sup> as a metric for two nuclei represented by histograms, and the  
691 variation can be simply evaluated by summing all pairwise distances.

692 **Earth mover's distance as a valid metric for histograms.** EMD originally arose in the field of

693 optimal transporting problems, recent studies show that it can be fruitfully applied to compare his-  
694 tograms. Thus, if one thinks of a histogram as a pile of dirt, then the EMD between two histograms  
695 is the minimum cost required to move the dirt in one pile to the other. Here, the cost is defined as  
696 the amount of dirt moved multiplied by the distance it is moved. Univariate EMD has several nice  
697 properties: (1) it's a true distance; (2) it doesn't need to assume the distribution form of histograms;  
698 (3) it's computationally efficient.

699 **Discriminative feature identification using EMD.** For each feature, a  $C \times C$  EMD matrix can be  
700 calculated, whose  $(i,j)$ -th entry is the distance between  $i$ -th nuclei histogram and  $j$ -th nuclei histo-  
701 gram. Then we use the given clustering result to sort the rows and columns, and discriminative  
702 features may pose a block diagonal EMD matrix. The ratio between BCV and WCV can be used  
703 to evaluate the feature's discriminative power between two clusters. BCV can be simply calculated  
704 by summing over all pairwise distance between the two clusters, and similarly, WCV can be simply  
705 calculated by summing over all pairwise distance within two clusters independently.

706

#### 707 **Multimodal intersection analysis between mouse and human liver samples.**

708 To access the correspondence between clusters identified in mouse and human samples, we  
709 adopted modified multimodal intersection analysis (MIA)<sup>83</sup>. Specifically, we ranked metabolites by  
710 the score computed using SCANPY<sup>84</sup>, which is z-score underlying the computation of a p-value  
711 (Student's t-test) for each gene for each cluster. Next gene sets of each cluster were defined as  
712 genes with the top 20 associated scores. And the significance of the intersection of gene sets  
713 between any pair of clusters was inferred using the hypergeometric distribution. The MIA map was  
714 finally displayed as a heatmap, with each element defined as the negative logarithm P-value (hy-  
715 pergeometric test) of the corresponding cluster pair.

#### 716 **Statistical analysis of human samples**

717 To exactly describe the statistical analysis in Fig. 4, we defined following terms:  $FBD_{R_i}$  is the fi-  
718 brotic boundary of region  $R_i$ ;  $PSP(j, FBD_{R_i})$  is a parallel strip whose distance to  $FBD_{R_i}$  is equal to  
719  $j \mu\text{m}$ ;  $AREA(j, i)$  is the territory between  $FBD_{R_i}$  and  $PSP(j, FBD_{R_i})$ ;  $Zone(j, i)$  is short for  
720  $AREA((j + 1) \times 100, i)$ ;  $CFBD(\text{cell}_i, Zone(j, k))$  is the distance ( $\mu\text{m}$ ) between  $\text{cell}_i$  and  $FBD_{R_k}$   
721 within  $Zone(j, k)$ ;  $NCC(\text{population}_i, \text{area}_1, \text{area}_2)$  is the ratio between the number of cells in  
722  $\text{population}_i$  within  $\text{area}_1$  and the number of cells in  $\text{population}_i$  within  $\text{area}_2$ .

723 The FBD is approximated according to *SIMS-View* and spatial single nucleus map (Supplementary  
724 Fig. 30). Coming to cases where FBD couldn't be well fitted by a single line segment, polylines are  
725 used, and the distance to FBD is simply adjusted to be the smallest among distances to all line  
726 segments.

727 The statistical analysis of Fig. 4h is conducted as following: In zone  $j$ ,  $j \in \{0,1,2,3,4\}$ , the red box-  
728 plot is the summarization of  $\{CFBD(\text{cell}_i, \text{Zone}(j, k)) \mid k \in \{1,2,3,4,5,6,7,8,9,10\}, \text{cell}_i \in$   
729  $\text{Hepatocyte C1}\}$ , and the green boxplot is the summarization of  $\{CFBD(\text{cell}_i, \text{Zone}(j, k)) \mid k \in$   
730  $\{1,2,3,4,5,6,7,8,9,10\}, \text{cell}_i \in \text{Hepatocyte C2}\}$ . The P-value is based on Wilcoxon rank sum test.

731 The statistical analysis of Fig. 4i is conducted as following: the x-axis is the distance between  
732  $\text{PSP}(j, \text{FBD}_{\text{Ri}})$  and corresponding fibrotic boundary ( $\text{FBD}_{\text{Ri}}$ ),  $i \in \{1,2,3,4,5,6,7,8,9,10\}, j \in [0,450]$ ;  
733 the y-axis is the normalized count ratio between C1 and C2, which is  $\frac{\text{NCC}(\text{C1}, \text{AREA}(j, i), \text{AREA}(j_{\max}, i))}{\text{NCC}(\text{C2}, \text{AREA}(j, i), \text{AREA}(j_{\max}, i))}$ ,  $i \in$   
734  $\{1,2,3,4,5,6,7,8,9,10\}, j \in [0,450]$ .

735 All parameters of boxplots are set as default using Seaborn (<https://seaborn.pydata.org>), a Python  
736 statistical data visualization toolbox.

737

738

## 739 Datasets

740 **Simulated datasets:** Four different human cell lines are cultured as a source of simulation (Sup-  
741 plementary Fig. 5a), and all the following datasets are manual alteration and a combination of the  
742 four cell lines.

743 Dataset 1: Use 4 cell lines as 4 clusters, for each cell, randomly add  $\text{noise\_ratio} \times \#\text{pixels}$  number  
744 of all-zero pixels.

745 Dataset 2: Use 4 cell lines as 4 clusters, for each pixel, multiply it with a random number drawn  
746 from  $U(0, \text{noise\_ratio})$ .

747 Dataset 3: Use 4 cell lines as cluster1, and the altered version of 4 cell lines as cluster2. Alteration  
748 method: for each cell, first randomly select  $\text{noise\_ratio} \times \#\text{pixel}$  pixels, then replace these pixels  
749 with samples drawn from feature-independent Gaussian fitted with original data.

750 Dataset 4: Use A549 cell line as cluster1, and use the 3 differently altered version as the other 3  
751 clusters. Alteration method: First, randomly select 2 dimensions,  $i$  and  $j$ . Then, for cluster2, multiply  
752  $\text{fold\_change}$  to the  $i$ -th dimension of all pixels of 10A cells, and the  $j$ -th dimension remains un-  
753 changed. For Cluster3, multiply  $\text{fold\_change}$  to both  $i$ -th and  $j$ -th dimension of all pixels of 10A  
754 cells. For cluster4: multiply  $\text{fold\_change}$  to the  $j$ -th dimension of all pixels of 10A cells, and the  $i$ -th  
755 dimension remains unchanged. The simulating method of dataset 4 is illustrated as (supplemen-  
756 tary Fig. 12a).

757 Dataset 5: Similar to dataset 4 but using Hela cell line.

758 Dataset 6: Similar to dataset 4 but using SK-BR-3 cell line.

759 Dataset 7: Similar to dataset 4 but using MCF 10A cell line.



760 Dataset 8: Use 10A cell line as cluster1, and use the 3 differently altered version as the other  
761 clusters. Alteration method: First, randomly select 2 dimensions,  $i$  and  $j$ , and calculate the  $mean_i$   
762 and  $variance_i$  for each cell. Second, for each cell, randomly divide pixels into two partitions of an  
763 equal number of pixels. Next, for cluster2, for each cell, replace the  $i$ -th dimension of the first  
764 partition with data drawn from  $Gaussian(fold\_chang \times mean_i, variance_i)$ , and replace the  $i$ -th  
765 dimension of the second partition with data drawn from  $Gaussian((2-fold\_change) \times mean_i,$   
766  $variance_i)$ . The  $j$ -th dimension remains unchanged. For cluster4, for each cell, replace the  $j$ -th  
767 dimension of the first partition with data drawn from  $Gaussian(fold\_change \times mean_j, variance_j)$ ,  
768 and replace the  $j$ -th dimension of the second partition with data drawn from  
769  $Gaussian((2-fold\_change) \times mean_j, variance_j)$ . The  $i$ -th dimension remains unchanged. For  
770 cluster3, the alteration for the  $i$ -th dimension is the same as cluster2, and the  $j$ -th dimension is the  
771 same with cluster4. The simulating method of dataset 8 is illustrated as (supplementary Fig. 12b).

772 Dataset 9: Use 2 differently altered versions of 10A cell line as two clusters. Alteration method:  
773 First, randomly select 2 dimensions,  $i$  and  $j$ , and calculate the  $mean_i$  and  $variance_i$  for each cell.  
774 Second, for each cell, randomly divide pixels into two partitions of an equal number of pixels. Next,  
775 for cluster1, for each cell, replace the  $i$ -th dimension of the first partition with data drawn from  
776  $Gaussian(fold\_chang \times mean_i, variance_i)$ , and replace the  $i$ -th dimension of the second partition  
777 with data drawn from  $Gaussian((2-fold\_change) \times mean_i, variance_i)$ . And replace the  $j$ -th di-  
778 mension of the first partition with data drawn from  $Gaussian(fold\_change \times mean_j, variance_j)$ ,  
779 and replace the  $j$ -th dimension of the second partition with data drawn from  
780  $Gaussian((2-fold\_change) \times mean_j, variance_j)$ . For cluster2, for each cell, replace the  $i$ -th di-  
781 mension of the first partition with data drawn from  $Gaussian(fold\_chang \times mean_i, variance_i)$ ,  
782 and replace the  $i$ -th dimension of the second partition with data drawn from  
783  $Gaussian((2-fold\_change) \times mean_i, variance_i)$ . And replace the  $j$ -th dimension of the second  
784 partition with data drawn from  $Gaussian(fold\_change \times mean_j, variance_j)$ , and replace the  $j$ -th  
785 dimension of the first partition with data drawn from  $Gaussian((2-fold\_change) \times mean_j,$   
786  $variance_j)$ . The simulating method of dataset 9 is illustrated as (supplementary Fig. 12c).

787

788 **Mixture cell datasets:** Mixture cell culture uses BrdU/IdU as ground truth label (Supplementary  
789 Fig. 13), and the BrdU/IdU stain does not affect the cell metabolic profiling (Supplementary Fig.  
790 21).

791 Dataset 10: A549 cell line stained with BrdU is mixed with Hela cell line (Supplementary Fig. 13a,  
792 b).

793 Dataset 11: NCTC1469 cell line stained with IdU is mixed with Hepa1-6 cell line (Supplementary  
794 Fig. 13c, d).

795

796 **Reporting Summary**

797 Further information on research design is available in the Nature Research Reporting Summary linked to  
798 this article.

799 **Data availability**

800 Raw SIMS data for mouse liver and lung (Fig. 1,2,3), and human liver R1 (Fig. 4) are available at  
801 Github ([https://github.com/yuanzhiyuan/SEAM/tree/master/SEAM/data/raw\\_tar](https://github.com/yuanzhiyuan/SEAM/tree/master/SEAM/data/raw_tar)). The rest of raw  
802 SIMS data and processed SIMS data are available at figshare (10.6084/m9.figshare.12622883,  
803 10.6084/m9.figshare.12622841, 10.6084/m9.figshare.12622838 and  
804 10.6084/m9.figshare.12622922). Geo-seq (Fig. 5) raw sequencing data and processed data have  
805 been deposited to NCBI GEO with accession number GSE153463.

806 **Code availability**

807 An open-source Python and MATLAB implementation of SEAM is available at GitHub ([https://](https://github.com/yuanzhiyuan/SEAM)  
808 [github.com/yuanzhiyuan/SEAM](https://github.com/yuanzhiyuan/SEAM), and <https://github.com/yuanzhiyuan/SIMS-Cut>).

809

810 **References**

- 811 1 Quail, D. F. & Joyce, J. A. Microenvironmental regulation of tumor progression and  
812 metastasis. *Nat Med* **19**, 1423-1437, doi:10.1038/nm.3394 (2013).
- 813 2 Riquelme, P. A., Drapeau, E. & Doetsch, F. Brain micro-ecologies: neural stem cell niches  
814 in the adult mammalian brain. *Philos T R Soc B* **363**, 123-137, doi:10.1098/rstb.2006.2016  
815 (2008).
- 816 3 Swain, P. S., Elowitz, M. B. & Siggia, E. D. Intrinsic and extrinsic contributions to  
817 stochasticity in gene expression. *P Natl Acad Sci USA* **99**, 12795-12800,  
818 doi:10.1073/pnas.162041399 (2002).
- 819 4 Zhang, J. W. & Li, L. H. Stem cell niche: Microenvironment and beyond. *J Biol Chem* **283**,  
820 9499-9503, doi:10.1074/jbc.R700043200 (2008).
- 821 5 Shukla, S. D. & Lim, R. W. Epigenetic effects of ethanol on the liver and gastrointestinal  
822 system. *Alcohol research: current reviews* **35**, 47 (2013).
- 823 6 Benly, P. Role of histamine in acute inflammation. *Journal of Pharmaceutical Sciences and*  
824 *Research* **7**, 373 %@ 0975-1459 (2015).
- 825 7 Stuart, T. & Satija, R. Integrative single-cell analysis. *Nat Rev Genet* **20**, 257-272,  
826 doi:10.1038/s41576-019-0093-7 (2019).
- 827 8 Pareek, V., Tian, H., Winograd, N. & Benkovic, S. J. Metabolomics and mass spectrometry  
828 imaging reveal channeled de novo purine synthesis in cells. *Science* **368**, 283-290 %@  
829 0036-8075 (2020).
- 830 9 Kennedy, D. E. *et al.* Novel specialized cell state and spatial compartments within the  
831 germinal center. *Nature Immunology*, 1-11 %@ 1529-2916 (2020).
- 832 10 Stoeckli, M., Chaurand, P., Hallahan, D. E. & Caprioli, R. M. Imaging mass spectrometry: a  
833 new technology for the analysis of protein expression in mammalian tissues. *Nat Med* **7**,  
834 493-496 %@ 1546-1170X (2001).
- 835 11 Sun, C. *et al.* Spatially resolved metabolomics to discover tumor-associated metabolic  
836 alterations. *Proc Natl Acad Sci U S A* **116**, 52-57, doi:10.1073/pnas.1808950116 (2019).
- 837 12 Hare, D. J. *et al.* Three-dimensional atlas of iron, copper, and zinc in the mouse cerebrum  
838 and brainstem. *Anal Chem* **84**, 3990-3997 %@ 0003-2700 (2012).
- 839 13 Sjövall, P., Lausmaa, J. & Johansson, B. Mass spectrometric imaging of lipids in brain tissue.  
840 *Anal Chem* **76**, 4271-4278 %@ 0003-2700 (2004).
- 841 14 Zavalin, A., Yang, J. & Caprioli, R. Laser beam filtration for high spatial resolution MALDI  
842 imaging mass spectrometry. *J Am Soc Mass Spectr* **24**, 1153-1156 %@ 1044-0305 (2013).
- 843 15 Niehaus, M., Soltwisch, J., Belov, M. E. & Dreisewerd, K. Transmission-mode MALDI-2  
844 mass spectrometry imaging of cells and tissues at subcellular resolution. *Nat Methods* **16**,  
845 925-931, doi:10.1038/s41592-019-0536-2 (2019).
- 846 16 Passarelli, M. K. *et al.* The 3D OrbiSIMS-label-free metabolic imaging with subcellular  
847 lateral resolution and high mass-resolving power. *Nat. Methods* **14**, 1175-+,  
848 doi:10.1038/Nmeth.4504 (2017).
- 849 17 Vickovic, S. *et al.* High-definition spatial transcriptomics for in situ tissue profiling. *Nat*  
850 *Methods* **16**, 987-990, doi:10.1038/s41592-019-0548-y (2019).
- 851 18 Keren, L. *et al.* A Structured Tumor-Immune Microenvironment in Triple Negative Breast  
852 Cancer Revealed by Multiplexed Ion Beam Imaging. *Cell* **174**, 1373-+,

853 doi:10.1016/j.cell.2018.08.039 (2018).

854 19 Xia, C., Fan, J., Emanuel, G., Hao, J. & Zhuang, X. Spatial transcriptome profiling by  
855 MERFISH reveals subcellular RNA compartmentalization and cell cycle-dependent gene  
856 expression. *Proc Natl Acad Sci U S A* **116**, 19490-19499, doi:10.1073/pnas.1912459116  
857 (2019).

858 20 Keren, L. *et al.* MIBI-TOF: A multiplexed imaging platform relates cellular phenotypes and  
859 tissue structure. *Sci Adv* **5**, doi:10.1126/sciadv.aax5851 (2019).

860 21 Robertson, A. R. The CIE 1976 color-difference formulae. *Color Research & Application* **2**,  
861 7-11 %@ 0361-2317 (1977).

862 22 Halpern, K. B. *et al.* Single-cell spatial reconstruction reveals global division of labour in  
863 the mammalian liver. *Nature* **542**, 352-356, doi:10.1038/nature21065 (2017).

864 23 Park, J. *et al.* Single-cell transcriptomics of the mouse kidney reveals potential cellular  
865 targets of kidney disease. *Science* **360**, 758-763, doi:10.1126/science.aar2131 (2018).

866 24 Moor, A. E. *et al.* Spatial Reconstruction of Single Enterocytes Uncovers Broad Zonation  
867 along the Intestinal Villus Axis. *Cell* **175**, 1156-+, doi:10.1016/j.cell.2018.08.063 (2018).

868 25 Becht, E. *et al.* Dimensionality reduction for visualizing single-cell data using UMAP. *Nat*  
869 *Biotechnol* **37**, 38-+, doi:10.1038/nbt.4314 (2019).

870 26 Wang, Y. J. *et al.* Multiplexed In Situ Imaging Mass Cytometry Analysis of the Human  
871 Endocrine Pancreas and Immune System in Type 1 Diabetes. *Cell Metab* **29**, 769-783 e764,  
872 doi:10.1016/j.cmet.2019.01.003 (2019).

873 27 Koller, D. & Friedman, N. *Probabilistic graphical models: principles and techniques*. (MIT  
874 press, 2009).

875 28 Hinton, G. E. & Salakhutdinov, R. R. Reducing the dimensionality of data with neural  
876 networks. *Science* **313**, 504-507, doi:10.1126/science.1127647 (2006).

877 29 Hinton, G. A Practical Guide to Training Restricted Boltzmann Machines. (2010).

878 30 Le Roux, N. & Bengio, Y. Representational power of restricted Boltzmann machines and  
879 deep belief networks. *Neural Comput* **20**, 1631-1649, doi:10.1162/neco.2008.04-07-510  
880 (2008).

881 31 Bengio, Y., Courville, A. & Vincent, P. Representation Learning: A Review and New  
882 Perspectives. *Ieee T Pattern Anal* **35**, 1798-1828, doi:10.1109/TPAMI.2013.50 (2013).

883 32 Boykov, Y., Veksler, O. & Zabih, R. Markov random fields with efficient approximations.  
884 *Proc Cvpr Ieee*, 648-655, doi:Doi 10.1109/Cvpr.1998.698673 (1998).

885 33 Boykov, Y., Veksler, O. & Zabih, R. Fast approximate energy minimization via graph cuts.  
886 *Ieee T Pattern Anal* **23**, 1222-1239, doi:10.1109/34.969114 (2001).

887 34 Doersch, C., Gupta, A. & Efros, A. A. Unsupervised Visual Representation Learning by  
888 Context Prediction. *2015 Ieee International Conference on Computer Vision (Iccv)*, 1422-  
889 1430, doi:10.1109/icc.2015.167 (2015).

890 35 Noroozi, M. & Favaro, P. Unsupervised Learning of Visual Representations by Solving  
891 Jigsaw Puzzles. *Computer Vision - Eccv 2016, Pt Vi* **9910**, 69-84, doi:10.1007/978-3-319-  
892 46466-4\_5 (2016).

893 36 Noroozi, M., Pirsiavash, H. & Favaro, P. Representation Learning by Learning to Count.  
894 *2017 Ieee International Conference on Computer Vision (Iccv)*, 5899-5907,  
895 doi:10.1109/icc.2017.628 (2017).

896 37 Hartigan, J. A. & Wong, M. A. Algorithm AS 136: A k-means clustering algorithm. *Journal*

897 *of the Royal Statistical Society. Series C (Applied Statistics)* **28**, 100-108 %@ 0035-9254  
898 (1979).

899 38 Kiselev, V. Y. *et al.* SC3: consensus clustering of single-cell RNA-seq data. *Nat Methods*  
900 **14**, 483-486, doi:10.1038/nmeth.4236 (2017).

901 39 Wang, B., Zhu, J. J., Pierson, E., Ramazzotti, D. & Batzoglou, S. Visualization and analysis  
902 of single-cell RNA-seq data by kernel-based similarity learning. *Nat. Methods* **14**, 414-+,  
903 doi:10.1038/nmeth.4207 (2017).

904 40 Maaten, L. v. d. & Hinton, G. Visualizing data using t-SNE. *J Mach Learn Res* **9**, 2579-2605  
905 (2008).

906 41 McInnes, L., Healy, J. & Astels, S. hdbSCAN: Hierarchical density based clustering. *Journal*  
907 *of Open Source Software* **2**, 205 %@ 2475-9066 (2017).

908 42 Brison, J. *et al.* ToF-SIMS imaging and depth profiling of HeLa cells treated with  
909 bromodeoxyuridine. *Surface and Interface Analysis* **43**, 354-357, doi:10.1002/sia.3415  
910 (2011).

911 43 Ramdas, A., Trillos, N. G. & Cuturi, M. On Wasserstein Two-Sample Testing and Related  
912 Families of Nonparametric Tests. *Entropy-Switz* **19**, doi:ARTN 47  
913 10.3390/e19020047 (2017).

914 44 Ben-Moshe, S. & Itzkovitz, S. Spatial heterogeneity in the mammalian liver. *Nat Rev*  
915 *Gastroenterol Hepatol*, doi:10.1038/s41575-019-0134-x (2019).

916 45 Sano, K. *et al.* Distributional Variation of P-450 Immunoreactive Hepatocytes in Human-  
917 Liver Disorders. *Hum Pathol* **20**, 1015-1020, doi:Doi 10.1016/0046-8177(89)90274-8  
918 (1989).

919 46 Brosch, M. *et al.* Epigenomic map of human liver reveals principles of zoned  
920 morphogenic and metabolic control. *Nat Commun* **9**, 4150, doi:10.1038/s41467-018-  
921 06611-5 (2018).

922 47 Shetty, S., Lalor, P. F. & Adams, D. H. Liver sinusoidal endothelial cells - gatekeepers of  
923 hepatic immunity. *Nat Rev Gastro Hepat* **15**, 555-567, doi:10.1038/s41575-018-0020-y  
924 (2018).

925 48 Ramachandran, P. *et al.* Resolving the fibrotic niche of human liver cirrhosis at single-cell  
926 level. *Nature*, doi:10.1038/s41586-019-1631-3 (2019).

927 49 Wu, X., Jiang, R., Zhang, M. Q. & Li, S. Network-based global inference of human disease  
928 genes. *Mol Syst Biol* **4**, 189, doi:10.1038/msb.2008.27 (2008).

929 50 Pillai, S. M. & Meredith, D. SLC36A4 (hPAT4) is a high affinity amino acid transporter when  
930 expressed in *Xenopus laevis* oocytes. *J Biol Chem* **286**, 2455-2460 %@ 0021-9258 (2011).

931 51 Féral, C. C. *et al.* CD98hc (SLC3A2) participates in fibronectin matrix assembly by  
932 mediating integrin signaling. *The Journal of cell biology* **178**, 701-711 %@ 1540-8140  
933 (2007).

934 52 Wang, S. *et al.* Lysosomal amino acid transporter SLC38A9 signals arginine sufficiency to  
935 mTORC1. *Science* **347**, 188-194 %@ 0036-8075 (2015).

936 53 Reimer, R. J. SLC17: a functionally diverse family of organic anion transporters. *Mol*  
937 *Aspects Med* **34**, 350-359 %@ 0098-2997 (2013).

938 54 Kempson, S. A., Zhou, Y. & Danbolt, N. C. The betaine/GABA transporter and betaine:  
939 roles in brain, kidney, and liver. *Frontiers in physiology* **5**, 159 %@ 1664-1042X (2014).

940 55 Chen, J. *et al.* Spatial transcriptomic analysis of cryosectioned tissue samples with Geo-

941 seq. *Nature Protocols* **12**, 566-580, doi:10.1038/nprot.2017.003 (2017).

942 56 Krueger, F. Trim galore: A wrapper tool around Cutadapt and FastQC to consistently apply  
943 quality and adapter trimming to FastQ files. **516**, 517 (2015).

944 57 Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes.  
945 *Nucleic Acids Res.* **47**, D766-D773, doi:10.1093/nar/gky955 (2019).

946 58 Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21,  
947 doi:10.1093/bioinformatics/bts635 (2013).

948 59 Anders, S., Pyl, P. T. & Huber, W. HTSeq-a Python framework to work with high-  
949 throughput sequencing data. *Bioinformatics* **31**, 166-169,  
950 doi:10.1093/bioinformatics/btu638 (2015).

951 60 Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion  
952 for RNA-seq data with DESeq2. *Genome Biol.* **15**, doi:ARTN 550  
953 10.1186/s13059-014-0550-8 (2014).

954 61 Zhang, Y. Y., Brady, M. & Smith, S. Segmentation of brain MR images through a hidden  
955 Markov random field model and the expectation-maximization algorithm. *IEEE T Med*  
956 *Imaging* **20**, 45-57, doi:Doi 10.1109/42.906424 (2001).

957 62 Panjwani, D. K. & Healey, G. Markov Random-Field Models for Unsupervised  
958 Segmentation of Textured Color Images. *IEEE T Pattern Anal* **17**, 939-954, doi:Doi  
959 10.1109/34.464559 (1995).

960 63 Hammersley, J. M. & Clifford, P. (1971).

961 64 Clifford, P. Markov random fields in statistics. *Disorder in physical systems: A volume in*  
962 *honour of John M. Hammersley* **19** (1990).

963 65 Besag, J. Spatial Interaction and Statistical-Analysis of Lattice Systems. *J Roy Stat Soc B*  
964 *Met* **36**, 192-236 (1974).

965 66 Panjwani, D. K. & Healey, G. Markov Random-Field Model for Unsupervised Segmentation  
966 of Textured Color Images (Vol 17, Pg 939, 1995). *IEEE T Pattern Anal* **17**, 1128-1128 (1995).

967 67 Kato, Z. & Pong, T. C. A Markov random field image segmentation model for color  
968 textured images. *Image Vision Comput* **24**, 1103-1114, doi:10.1016/j.imavis.2006.03.005  
969 (2006).

970 68 Chen, F. Q., Wu, Y., Bu, Y. D. & Zhao, G. D. Spectral Classification Using Restricted  
971 Boltzmann Machine. *Publications of the Astronomical Society of Australia* **31**,  
972 doi:10.1017/pasa.2013.38 (2014).

973 69 Boykov, Y. & Kolmogorov, V. An experimental comparison of min-cut/max-flow  
974 algorithms for energy minimization in vision. *IEEE T Pattern Anal* **26**, 1124-1137, doi:Doi  
975 10.1109/TPAMI.2004.60 (2004).

976 70 Boykov, Y. & Funka-Lea, G. Graph Cuts and Efficient N-D Image Segmentation. *Int J*  
977 *Comput Vision* **70**, 109-131, doi:10.1007/s11263-006-7934-5 (2006).

978 71 Hinton, G. E. Training products of experts by minimizing contrastive divergence. *Neural*  
979 *Computation* **14**, 1771-1800, doi:Doi 10.1162/089976602760128018 (2002).

980 72 Tieleman, T. 1064-1071.

981 73 Keyvanrad, M. A. & Homayounpour, M. M. A brief survey on deep belief networks and  
982 introducing a new object oriented toolbox (DeeBNet). *arXiv preprint arXiv:1408.3264*  
983 (2014).

984 74 Larochelle, H. & Bengio, Y. 536-543.

985 75 Hinton, G., Vinyals, O. & Dean, J. Distilling the knowledge in a neural network. *arXiv*  
986 *preprint arXiv:1503.02531* (2015).

987 76 Veselkov, K. A. *et al.* Optimized Preprocessing of Ultra-Performance Liquid  
988 Chromatography/Mass Spectrometry Urinary Metabolic Profiles for Improved Information  
989 Recovery. *Anal Chem* **83**, 5864-5872, doi:10.1021/ac201065j (2011).

990 77 Lu, Y. Unsupervised learning on neural network outputs: with application in zero-shot  
991 learning. *arXiv preprint arXiv:1506.00990* (2015).

992 78 Wu, Z., Xiong, Y., Yu, S. X. & Lin, D. 3733-3742.

993 79 Duin, R. P. W. & Pekalska, E. The dissimilarity space: Bridging structural and statistical  
994 pattern recognition. *Pattern Recogn Lett* **33**, 826-832 %@ 0167-8655 (2012).

995 80 Tax, D. M. J., Loog, M., Duin, R. P. W., Cheplygina, V. & Lee, W.-J. 222-234 (Springer).

996 81 Glorot, X., Bordes, A. & Bengio, Y. 315-323.

997 82 Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. *arXiv preprint*  
998 *arXiv:1412.6980* (2014).

999 83 Moncada, R. *et al.* Integrating microarray-based spatial transcriptomics and single-cell  
1000 RNA-seq reveals tissue architecture in pancreatic ductal adenocarcinomas. *Nat*  
1001 *Biotechnol*, doi:10.1038/s41587-019-0392-8 (2020).

1002 84 Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data  
1003 analysis. *Genome Biol* **19**, 15, doi:10.1186/s13059-017-1382-0 (2018).

1004

1005 **Acknowledgements**

1006 The authors would like to acknowledge Imaging Core Facility, Technology Center for Protein Sciences,  
1007 Tsinghua University for assistance of using LMD7000. We also thank Yalan Chen from Imaging Core  
1008 Facility for her detailed instruction on LMD7000. We thank Center of Laboratory Animal Resources, Tsing-  
1009 hua University for mice maintenance and providing CM1900 Cryostat. We thank Hongjun Li for the com-  
1010 puting resource supporting. We thank Hui Zhang for the help of ethics material preparation. We thank  
1011 Minglei Shi, Yisi Li, Zhaofeng Ye, Rui Qi and all other members of our lab for valuable comments and  
1012 discussions. We thank Minping Qian for helpful advice on algorithm development. This work was partly  
1013 supported by National Basic Research Program of China (2018YFA0801402, 2018YFB0704304,  
1014 2017YFA0505503), National Nature Science Foundation of China (31871343, 21974078, 21727813,  
1015 21621003) and foundation of BNRist (BNR2019TD01020).

1016 **Author contributions**

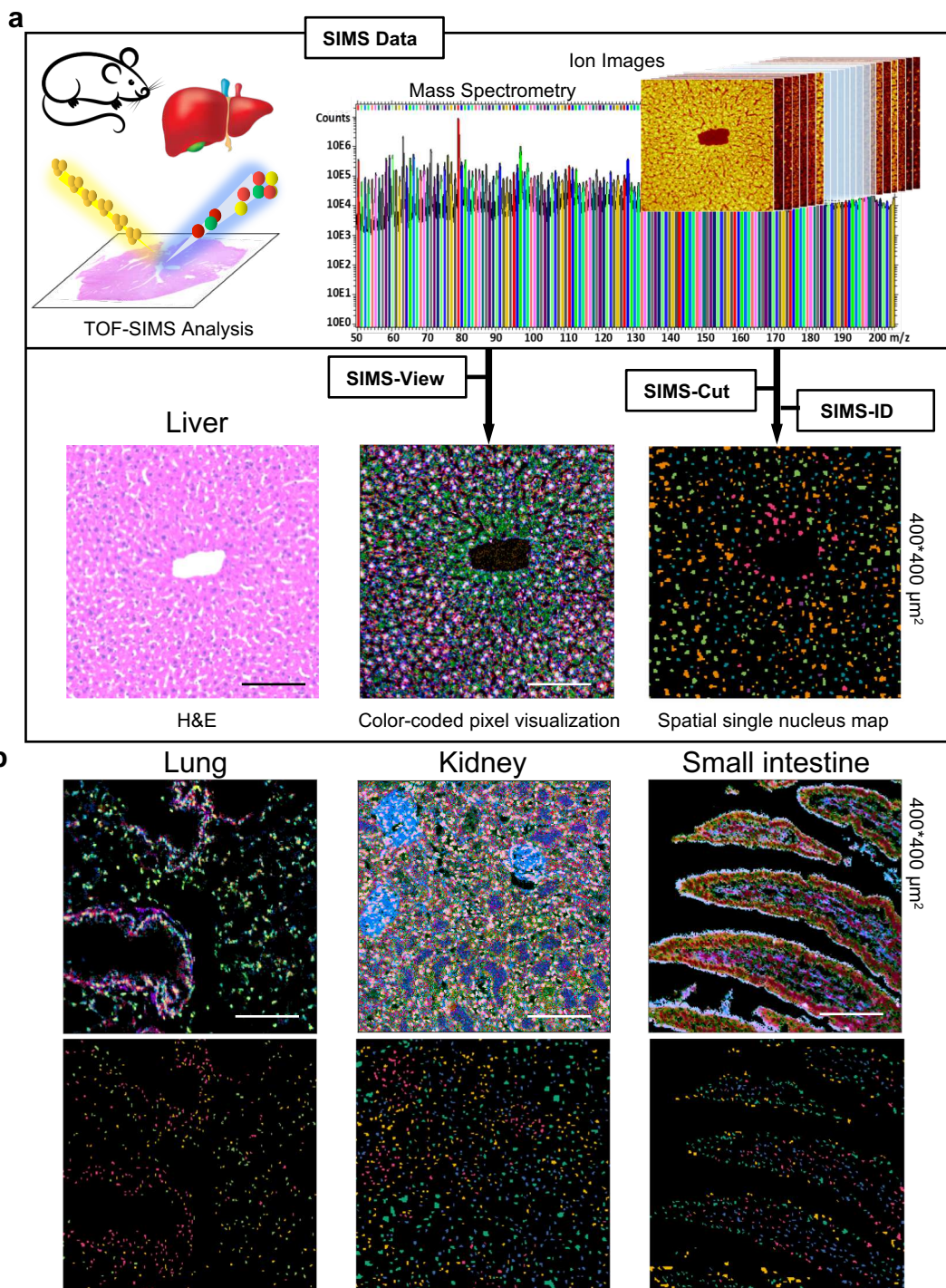
1017 Y.C, M.Q.Z and X.Z conceived and designed the project. L.C designed the IMS experiment and generated  
1018 the IMS data. Q.Z processed the mouse and human sample assisted by W.S, and generated IHC and HE  
1019 imaging data. Q.Z and L.C designed and conducted the cell culture and BrdU staining experiment. Q.Z

1020 designed and conducted the modified Geo-seq experiment. Z.Y developed and implemented the algo-  
1021 rithms under the guidance of M.Q.Z and Y.C, and assisted by Q.Z. Z.Y analyzed the SIMS data, and Q.Z  
1022 analyzed the spatial transcriptome data. Y.Z and S.Y provided the clinical samples. L.P and S.Q guided  
1023 the histological annotation. S.L gave suggestions on the application of the method. Z.Y, Q.Z, and L.C  
1024 completed the figures and manuscript with the guidance of Y.C. X.Z and M.Q.Z.

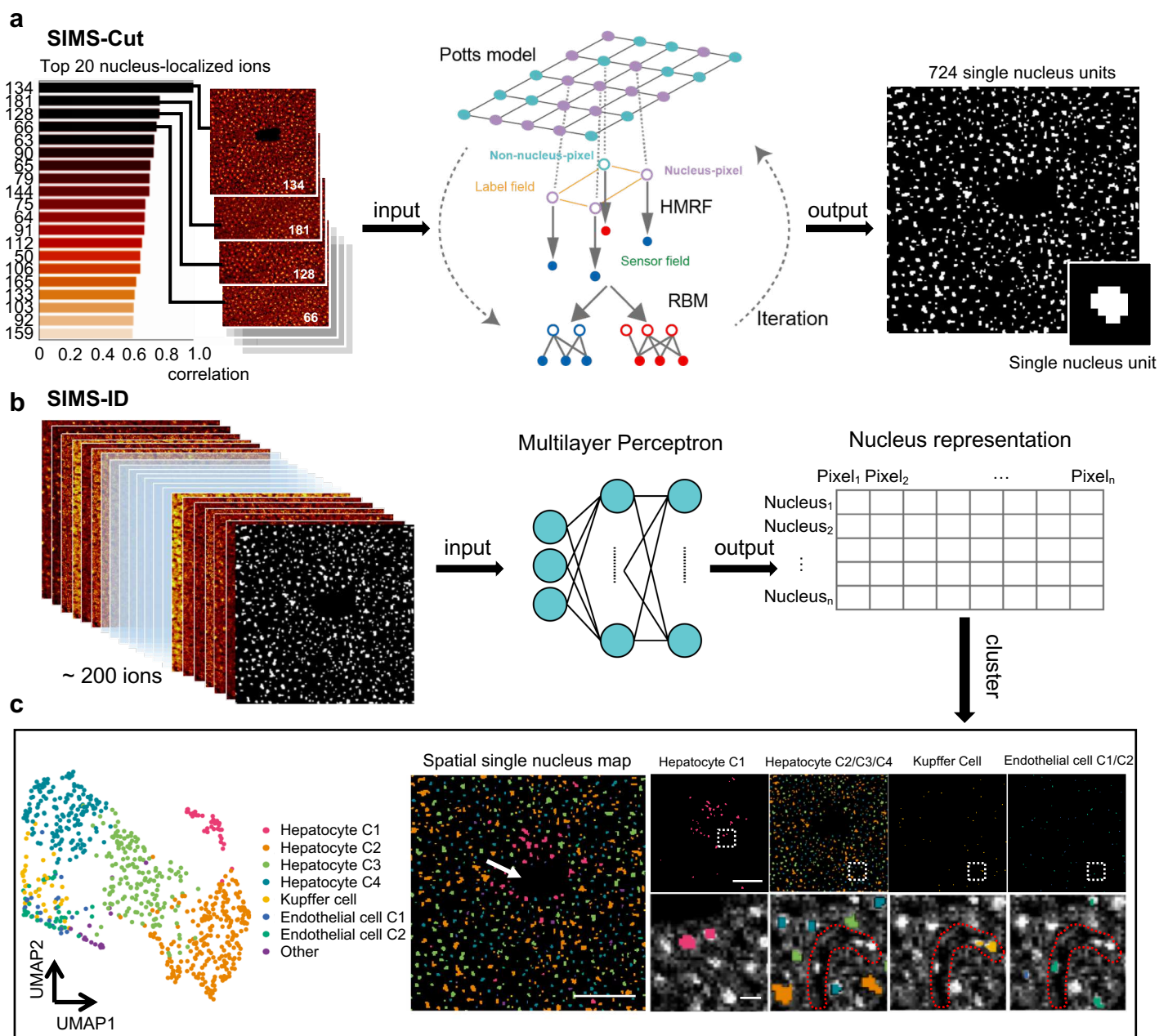
1025 **Competing interests**

1026 The authors declare no competing interests.

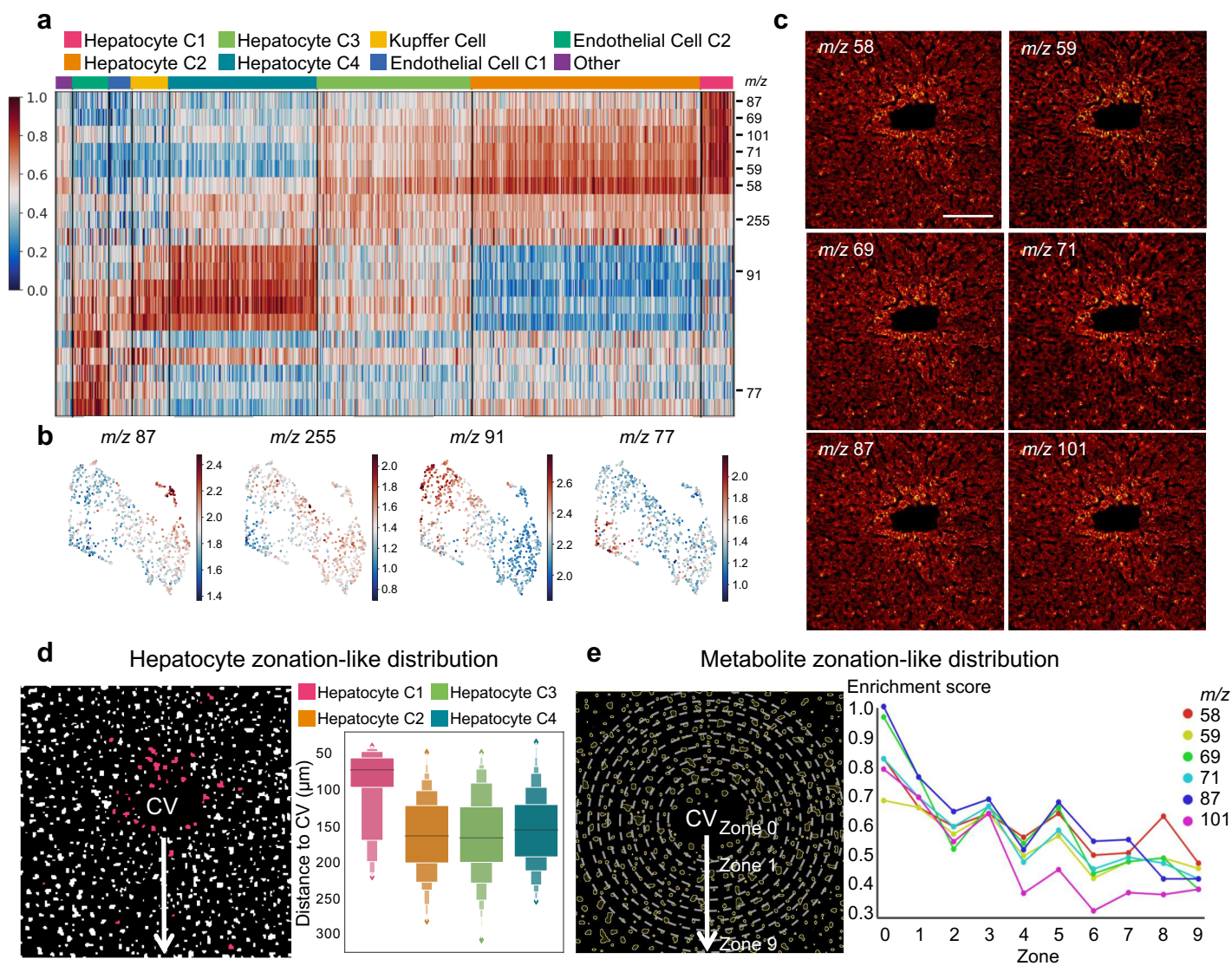




**Fig. 1 | SEAM captures spatial metabolic heterogeneity at single nucleus resolution.**  
**a**, Overview of SEAM. (Left) Tissue samples on glass slides are analyzed by TOF-SIMS to generate multiplex SIMS data containing mass spectrometry and ion images (Right). (Bottom left) H&E staining of mouse liver central vein region. (Bottom middle) Color-coded pixel visualization is obtained by SIMS-View. (Bottom right) Spatial single nucleus map is obtained by a sequential of algorithms: SIMS-Cut (segmentation), SIMS-ID (representation), and SIMS-Cluster (clustering). **b**, SEAM scales to different mouse tissues with different cell density and distribution pattern. First row is color-coded pixel visualization by SIMS-view to differentiate metabolic patterns at pixel level. Second row is spatial single nucleus map for cell type visualization at original tissue space. Scale bar 100 $\mu\text{m}$ . In Fig. 1a, Mouse illustration: Image by [OpenClipart-Vectors](#) from [Pixabay](#). Liver illustration: Image by [zachvanstone8](#) from [Pixabay](#).



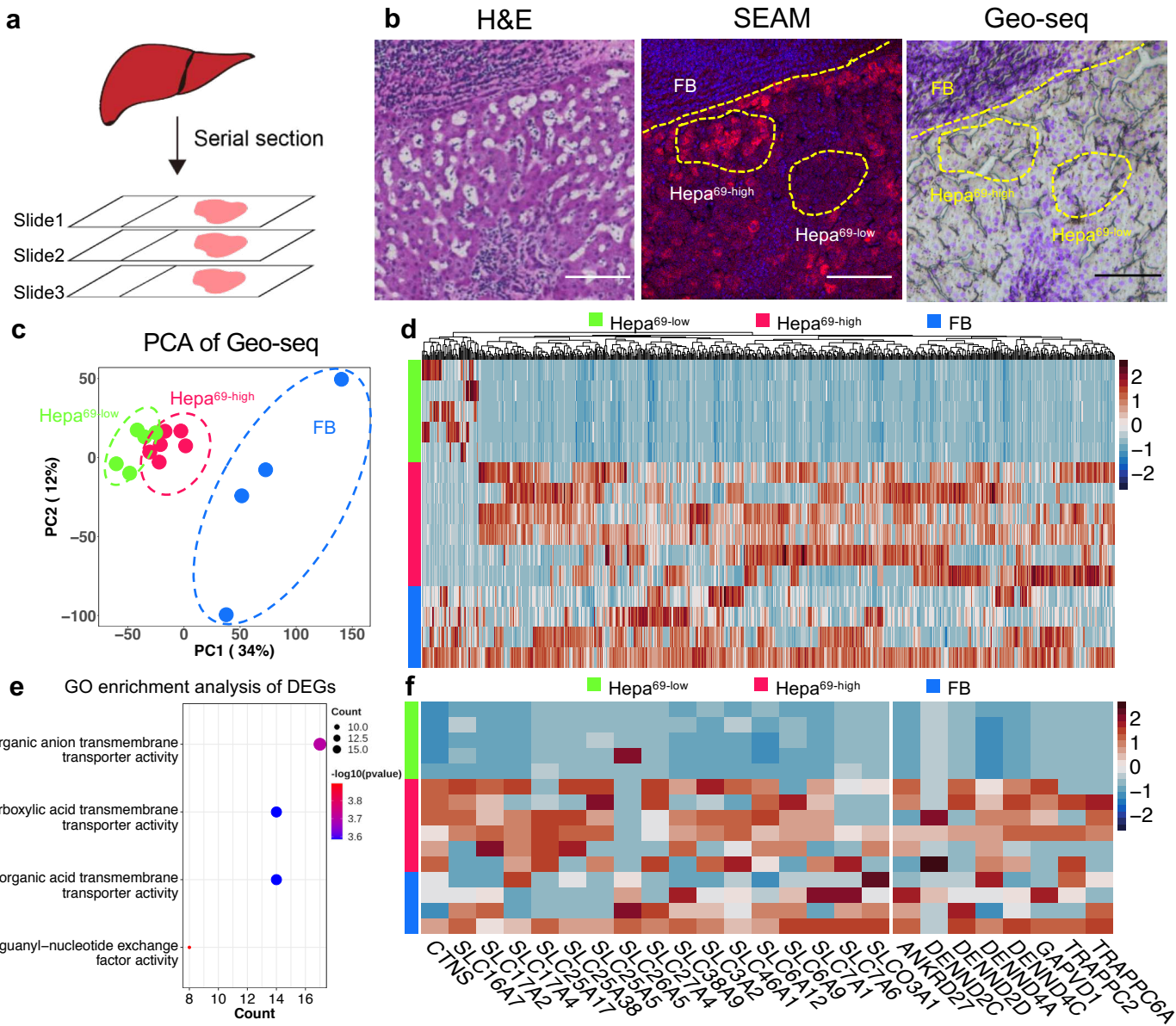
**Fig. 2 | Algorithms design and performance.** **a**, Sketch of SIMS-Cut, leveraging Potts model as prior for pixel labels and Restricted Boltzmann Machines as conditional distribution of pixel intensities. (Left) Top 20 nucleus-localized ions. (Middle) Iterative optimization between subproblems (See Methods). (Right) Cell segmentation mask. **b**, Sketch of SIMS-ID, learning vector-formed representation for each segmented cells using self-representation learning. (Left) multiplex SIMS data combined with cell segmentation mask. (Middle) A neural network for a auxiliary classification task. (Right) Single nucleus representation output. **c**, Demonstration of algorithms on central vein (CV) of wild type mouse liver. (Left) UMAP visualization of single nucleus using SIMS-ID representation, colored by SIMS-Cluster identified cell types. (Middle) Spatial single nucleus map. White arrow indicates CV. Scale bar 100 $\mu$ m. (Right top) Respective layout of cell populations. Scale bar 100 $\mu$ m. (Right bottom) Zoom in images of each population merged with grey scaled image of m/z 134. Red dotted area indicate liver sinusoid. Scale bar 10 $\mu$ m.



**Fig. 3 | SEAM detects zonation-like metabolic pattern in wild type mouse liver.** **a**, Differential metabolite analysis of mouse liver tissue in Fig. 2c. **b**, UMAP colored by abundance of representative differential metabolites. **c**, Ion images of a ion series with zonation-like distribution identified by differential analysis in **Fig. 3a**. Scale bar 100µm. **d**, Hepatocyte C1 subpopulation shows zonation-like distribution. (Left) Schematic diagram of strategy of measuring cell-to-CV distance. (Right) Hepatocyte C1 shows significantly smaller distance to CV than other clusters. **e**, Metabolite series show zonation-like distribution. (Left) Schematic diagram of strategy of measuring metabolite-to-CV distance: Concentric circles with distance of arithmetic sequence from CV partition the liver lobule into 9 zones. (Right) 6 metabolic markers of Hepatocyte C1 show gradient decrease away from CV. X-axis: zone number, Y-axis: enrichment score of each metabolites, which is the proportion of hepatocytes that highly express each metabolites in each zones.

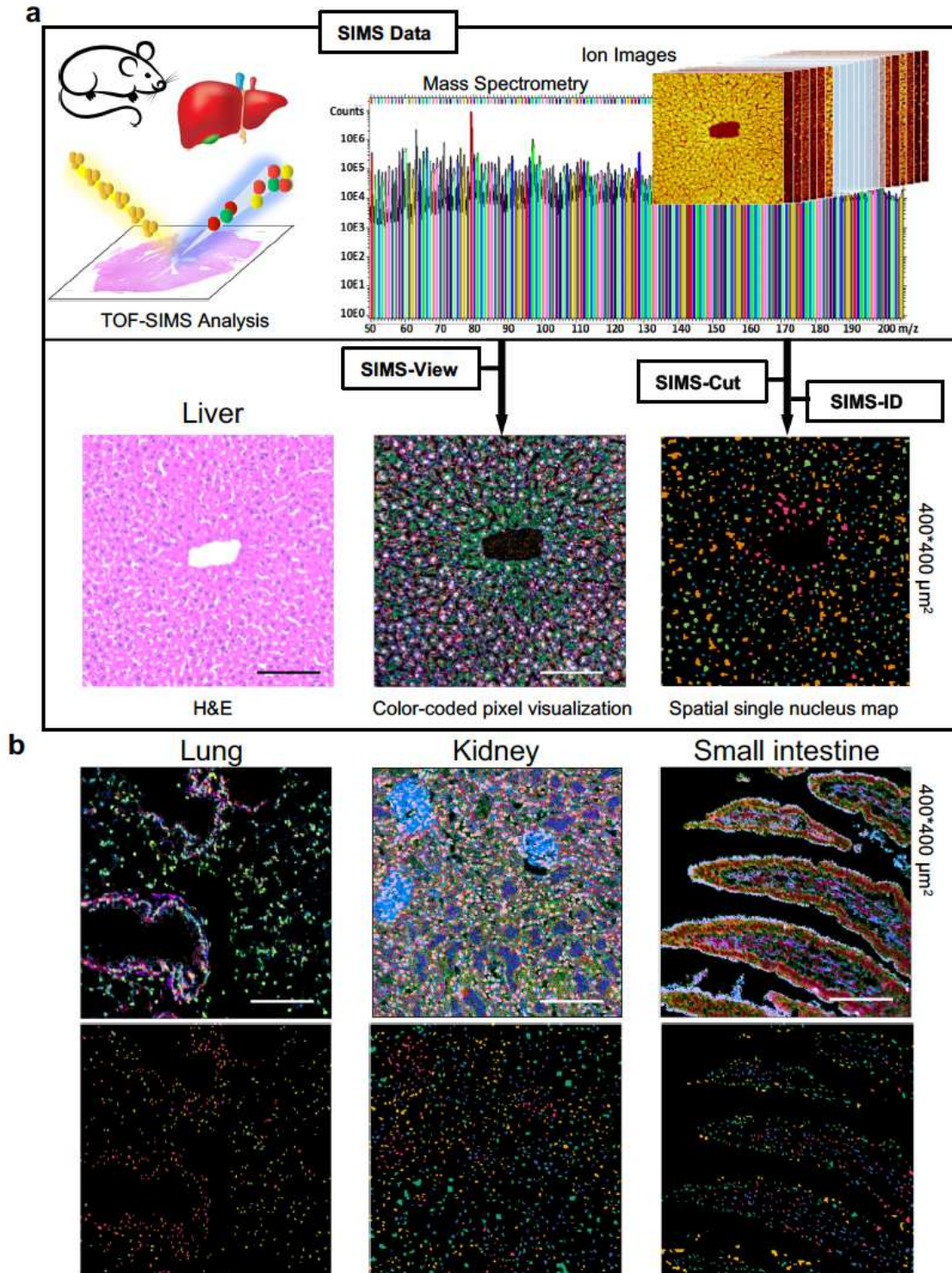


**Fig. 4 | SEAM identifies hepatocyte subtypes with differential metabolic state associated with spatial localization.** **a**, H&E staining of human liver sample post TOF-SIMS analysis. Scale bar 500 $\mu$ m. **b**, Zoom in H&E images of 4 different regions. Bottom. White arrows indicate fibrotic and inflammation niche. Scale bar 100 $\mu$ m. **c**, SEAM results of 4 regions. First column is color-coded pixel visualizations. Scale bar 100 $\mu$ m. Second column is UMAP colored by cell clusters. Third column is spatial single nucleus map. **d**, Spatial single nucleus maps of respective clusters merged with grey scaled ion image of  $m/z$  134. Scale bar 100 $\mu$ m. **e**, Differential metabolite analysis of cell clusters. **f**, (Top and middle row) Hepatocyte C1 enriched metabolites. Scale bar 100 $\mu$ m. (Bottom row left) Merged ion image of  $m/z$  69(Red) and  $m/z$  134 (Blue). (Bottom row middle and right) Spatial localization of hepatocyte C1 and C2 respectively merged with grey-scaled ion image of  $m/z$  69. **h**, Hepatocyte C1 is consistently closer to fibrotic boundary (FBD) than C2 within all 5 zones. (Left) Schematic diagram of zone definition and distance calculation. (Right) Paired boxplots of distances between C1/C2 and FBD. For Wilcoxon Rank Sum test, P-value > 0.05 is not shown on the plot. P-value  $\leq$  0.05 (\*), P-value  $\leq$  0.01 (\*\*), P-value  $\leq$  0.001 (\*\*\*) and P-value  $\leq$  0.0001 (\*\*\*\*) are shown. **i**, Normalized count of hepatocyte C1 is consistently higher than C2. (Left) Schematic diagram of normalized count ratio calculation. (Right) Normalized count ratio between C1 and C2 is a function of the distance of the outer edge (indicated by the gray line in the left part of Fig. 4i) to the FBD.



**Fig. 5 | Spatial transcriptome validated metabolism associated gene expression alteration in heterogeneous hepatocyte subtypes identified by SEAM.** **a**, Serial sections were made for cross validation among different assays. **b**, Geo-seq was performed at same location (Right) in the adjacent slide of SEAM assay (Middle,  $m/z$  134 in blue and  $m/z$  69 in red) to obtain continuous tissue spatial structure. Yellow dashed area representatively indicate the captured regions for Geo-seq. Scale bar 100 $\mu$ m. **c**, PCA plot of transcriptomic profiles from a total 15 samples of different regions. **d**, Heatmap of filtered differentially expressed genes (DEGs) between Hepa<sup>69high</sup> and Hepa<sup>69low</sup> cells. **e**, GO enrichment of DEGs. **f**, Heatmap of DEGs enriched in GO terms in **e**. Upper part is consensus 14 genes in top 3 GO terms, and lower part is 8 genes enriched in last GO term.

# Figures



**Figure 1**

SEAM captures spatial metabolic heterogeneity at single nucleus resolution. a, Overview of SEAM. (Left) Tissue samples on glass slides are analyzed by TOF-SIMS to generate multiplex SIMS data containing mass spectrometry and ion images (Right). (Bottom left) H&E staining of mouse liver central vein region.

(Bottom middle) Color-coded pixel visualization is obtained by SIMS-View. (Bottom right) Spatial single nucleus map is obtained by a sequential of algorithms: SIMS-Cut (segmentation), SIMS-ID (representation), and SIMS-Cluster (clustering). b, SEAM scales to different mouse tissues with different cell density and distribution pattern. First row is color-coded pixel visualization by SIMS-view to differentiate metabolic patterns at pixel level. Second row is spatial single nucleus map for cell type visualization at original tissue space. Scale bar 100 $\mu$ m. In Fig. 1a, Mouse illustration: Image by OpenClipart-Vectors from Pixabay. Liver illustration: Image by zachvanstone8 from Pixabay.

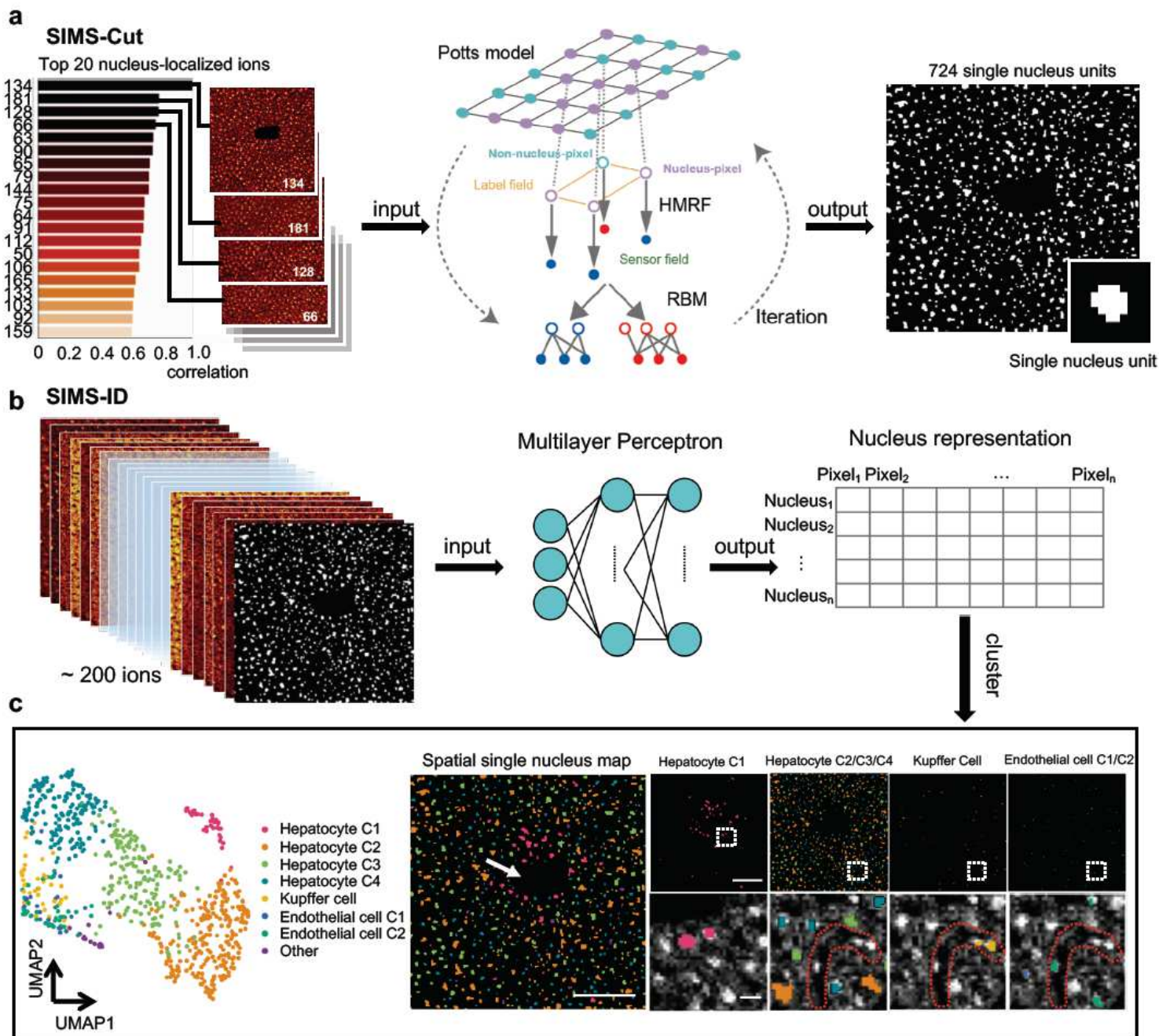


Figure 2



Algorithms design and performance. a, Sketch of SIMS-Cut, leveraging Potts model as prior for pixel labels and Restricted Boltzmann Machines as conditional distribution of pixel intensities. (Left) Top 20 nucleus-localized ions. (Middle) Iterative optimization between subproblems (See Methods). (Right) Cell segmentation mask. b, Sketch of SIMS-ID, learning vector-formed representation for each segmented cells using self-representation learning. (Left) multiplex SIMS data combined with cell segmentation mask. (Middle) A neural network for a auxiliary classification task. (Right) Single nucleus representation output. c, Demonstration of algorithms on central vein (CV) of wild type mouse liver. (Left) UMAP visualization of single nucleus using SIMS-ID representation, colored by SIMS-Cluster identified cell types. (Middle) Spatial single nucleus map. White arrow indicates CV. Scale bar 100 $\mu$ m. (Right top) Respective layout of cell populations. Scale bar 100 $\mu$ m. (Right bottom) Zoom in images of each population merged with grey scaled image of m/z 134. Red dotted area indicate liver sinusoid. Scale bar 10 $\mu$ m.

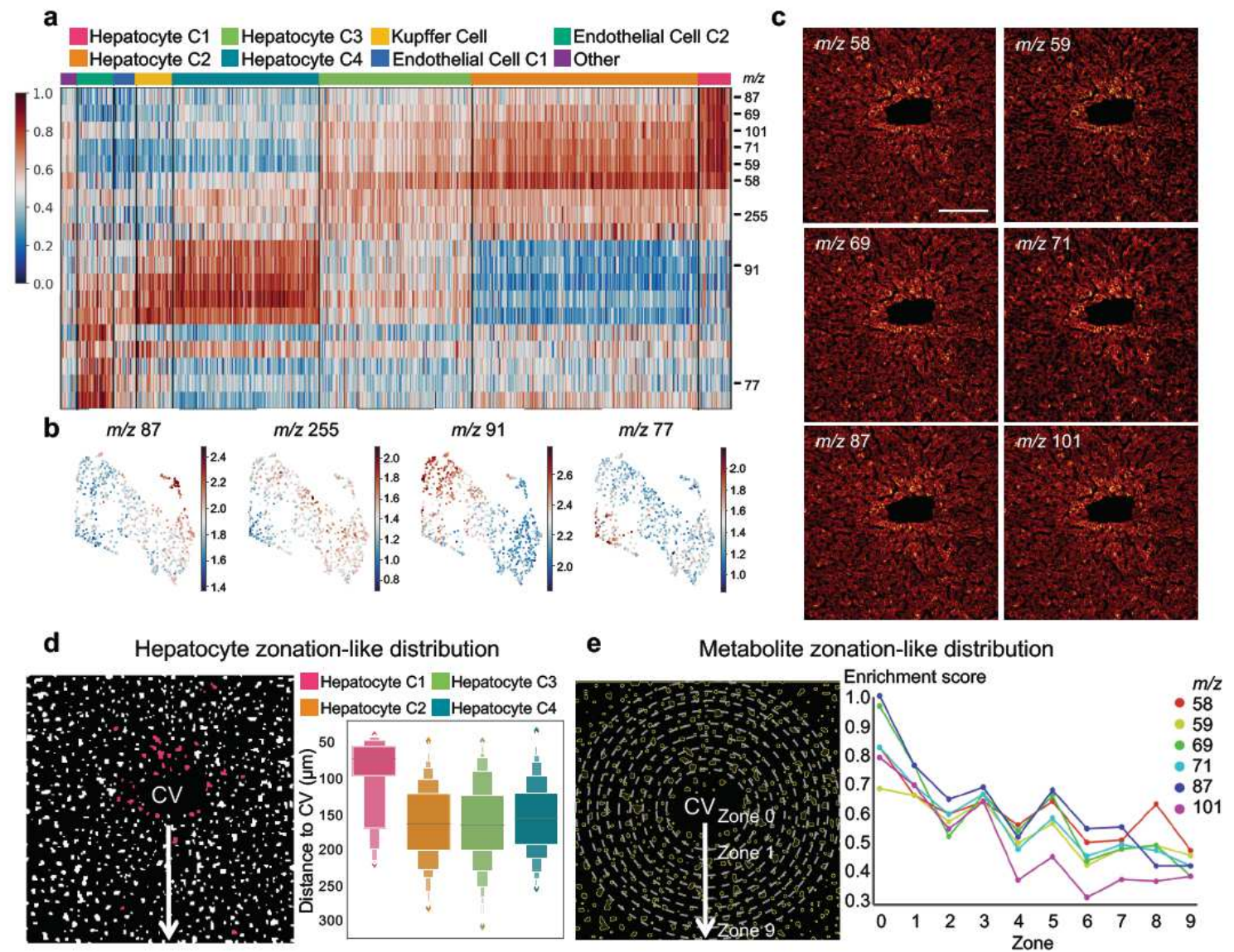
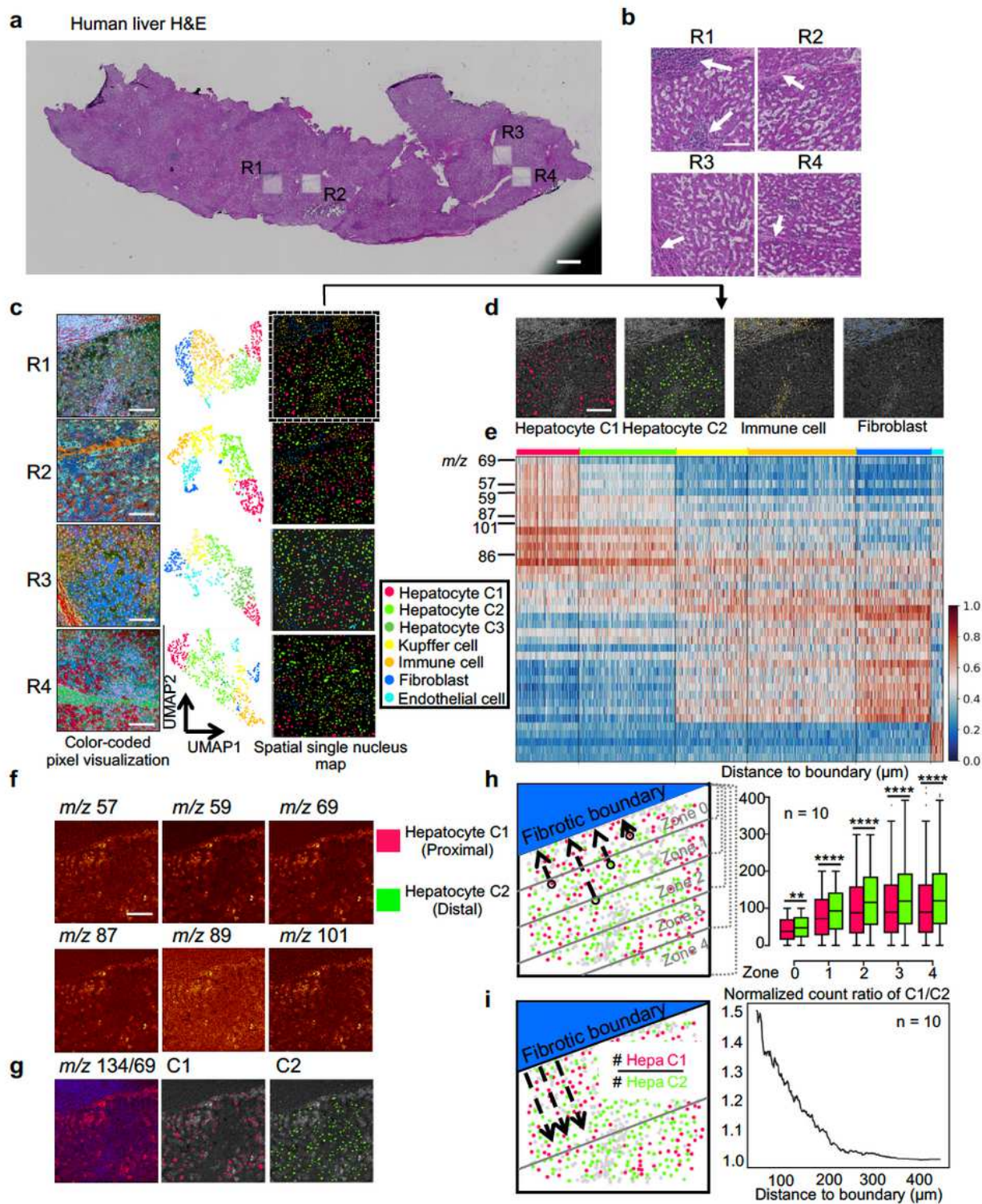


Figure 3

SEAM detects zonation-like metabolic pattern in wild type mouse liver. a, Differential metabolite analysis of mouse liver tissue in Fig. 2c. b, UMAP colored by abundance of representative differential metabolites. c, Ion images of a ion series with zonation-like distribution identified by differential analysis in Fig. 3a. Scale bar 100 $\mu$ m. d, Hepatocyte C1 subpopulation shows zonation-like distribution. (Left) Schematic diagram of strategy of measuring cell-to-CV distance. (Right) Hepatocyte C1 shows significantly smaller distance to CV than other clusters. e, Metabolite series show zonation-like distribution. (Left) Schematic diagram of strategy of measuring metabolite-to-CV distance: Concentric circles with distance of arithmetic sequence from CV partition the liver lobule into 9 zones. (Right) 6 metabolic markers of Hepatocyte C1 show gradient decrease away from CV. X-axis: zone number, Y-axis: enrichment score of each metabolites, which is the proportion of hepatocytes that highly express each metabolites in each zones.



**Figure 4**

SEAM identifies hepatocyte subtypes with differential metabolic state associated with spatial localization. a, H&E staining of human liver sample post TOF-SIMS analysis. Scale bar 500 $\mu$ m. b, Zoom in H&E images of 4 different regions. Bottom. White arrows indicate fibrotic and inflammation niche. Scale bar 100 $\mu$ m. c, SEAM results of 4 regions. First column is color-coded pixel visualizations. Scale bar 100 $\mu$ m. Second column is UMAP colored by cell clusters. Third column is spatial single nucleus map. d,

Spatial single nucleus maps of respective clusters merged with grey scaled ion image of m/z 134. Scale bar 100 $\mu$ m. e, Differential metabolite analysis of cell clusters. f, (Top and middle row) Hepatocyte C1 enriched metabolites. Scale bar 100 $\mu$ m. (Bottom row left) Merged ion image of m/z 69(Red) and m/z 134 (Blue). (Bottom row middle and right) Spatial localization of hepatocyte C1 and C2 respectively merged with greyscaled ion image of m/z 69. h, Hepatocyte C1 is consistently closer to fibrotic boundary (FBD) than C2 within all 5 zones. (Left) Schematic diagram of zone definition and distance calculation. (Right) Paired boxplots of distances between C1/C2 and FBD. For Wilcoxon Rank Sum test, P-value > 0.05 is not shown on the plot. P-value  $\leq$  0.05 (\*), P-value  $\leq$  0.01 (\*\*), P-value  $\leq$  0.001 (\*\*\*) and Pvalue  $\leq$  0.0001 (\*\*\*\*) are shown. i, Normalized count of hepatocyte C1 is consistently higher than C2. (Left) Schematic diagram of normalized count ratio calculation. (Right) Normalized count ratio between C1 and C2 is a function of the distance of the outer edge (indicated by the gray line in the left part of Fig. 4i) to the FBD.

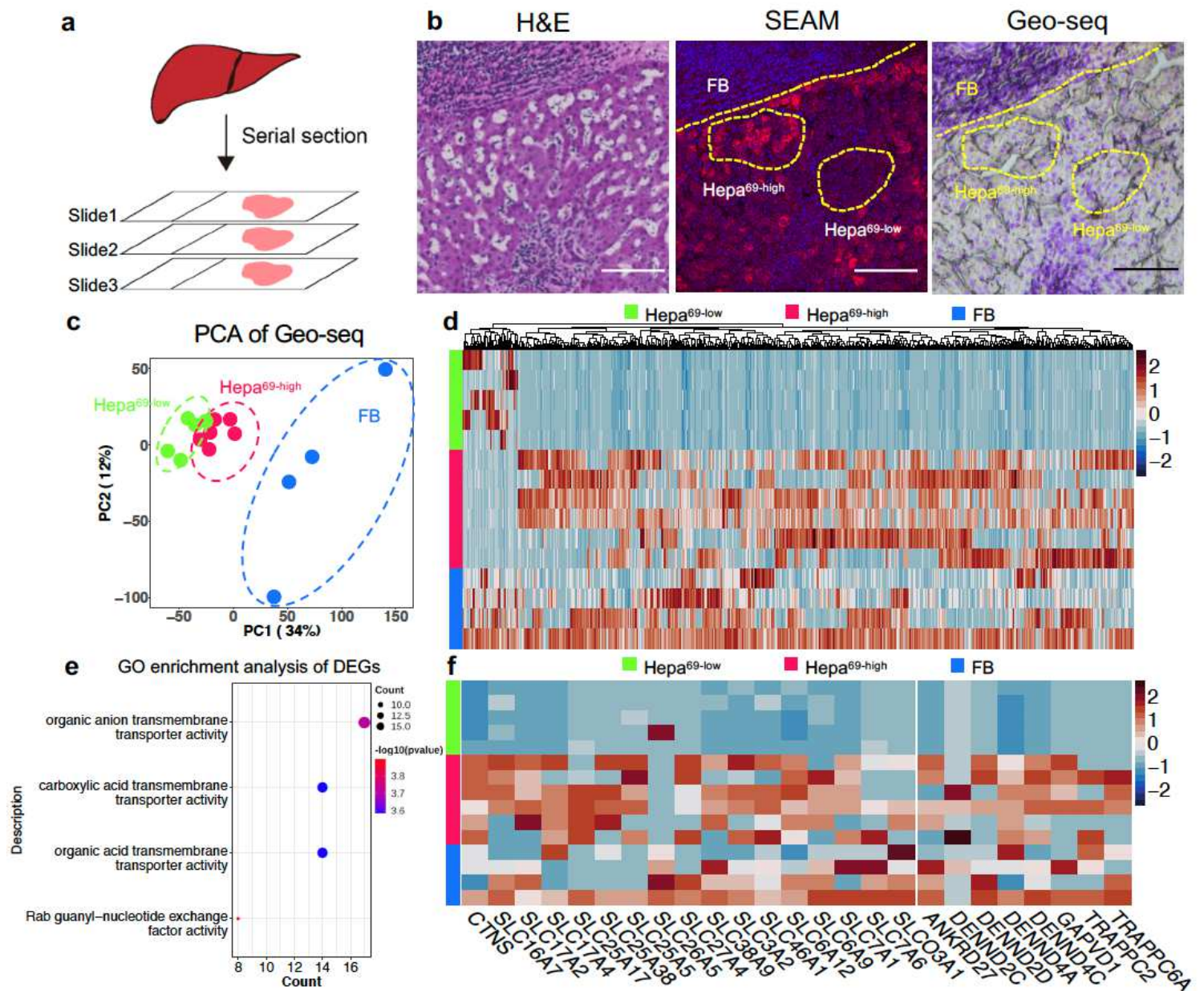


Figure 5

Spatial transcriptome validated metabolism associated gene expression alteration in heterogeneous hepatocyte subtypes identified by SEAM. a, Serial sections were made for cross validation among different assays. b, Geo-seq was performed at same location (Right) in the adjacent slide of SEAM assay (Middle, m/z 134 in blue and m/z 69 in red) to obtain continuous tissue spatial structure. Yellow dashed area representatively indicate the captured regions for Geoseq. Scale bar 100 $\mu$ m. c, PCA plot of transcriptomic profiles from a total 15 samples of different regions. d, Heatmap of filtered differentially expressed genes (DEGs) between Hepa69high and Hepa69low cells. e, GO enrichment of DEGs. f, Heatmap of DEGs enriched in GO terms in e. Upper part is consensus 14 genes in top 3 GO terms, and lower part is 8 genes enriched in last GO term.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementarySEAM.pdf](#)
- [SupplementaryTable2SEAM.xlsx](#)
- [nrreportingsummarySEAM.pdf](#)
- [nreditorialpolicychecklistSEAM.pdf](#)