1 **Cell segmentation-free inference of cell types from *in situ* transcriptomics data**

2

3 **Author names**

4 Jeongbin Park[1,2,3,†], Wonyl Choi[4,†], Sebastian Tiesmeyer[1], Brian Long[5], Lars E. Borm[6], Emma

5 Garren[5], Thuc Nghi Nguyen[5], Bosiljka Tasic[5], Simone Codeluppi[6,7], Tobias Graf[1], Matthias

6 Schlesner[8], Oliver Stegle[3,9], Roland Eils[1,10,‡,*] & Naveed Ishaque[1,‡,*]

7

8 **Affiliations**

9 [1]Digital Health Center, Berlin Institute of Health (BIH) and Charité Universitätsmedizin, Berlin,

10 Germany;

11 [2]Faculty of Biosciences, Heidelberg University, Heidelberg, Germany;

12 [3]Division of Computational Genomics and System Genetics, German Cancer Research Center

13 (DKFZ), Heidelberg, Germany;

14 [4]Department of Computer Science, Boston University, Boston, the United States of America;

15 [5]Allen Institute for Brain Science, Seattle, WA, USA;

16 [6]Division of molecular neurobiology, Department of medical biochemistry and biophysics,

17 Karolinska Institutet, Stockholm, Sweden;

18 [7]Science for life laboratory, Stockholm, Sweden;

19 [8]Bioinformatics and Omics Data Analytics, German Cancer Research Center (DKFZ),

20 Heidelberg, Germany;

21 [9]Genome Biology Unit, European Molecular Biology Laboratory (EMBL), Heidelberg, Germany;

22 [10]Health Data Science Unit, Heidelberg University Hospital, Heidelberg, Germany;

23

24 **Author List Footnotes**

25 [†]These authors contributed equally to this work.

26   ‡These authors jointly supervised the work.

27

## 28   **Contact information**

29   *Correspondence: Roland Eils (roland.eils@charite.de) and Naveed Ishaque

30   (naveed.ishaque@charite.de)

31 **Abstract**

32 Multiplexed fluorescence *in situ* hybridization techniques have enabled cell-type identification,

33 linking transcriptional heterogeneity with spatial heterogeneity of cells. However, inaccurate cell

34 segmentation reduces the efficacy of cell-type identification and tissue characterization. Here,

35 we present a novel method called Spot-based Spatial cell-type Analysis by Multidimensional

36 mRNA density estimation (SSAM), a robust cell segmentation-free computational framework for

37 identifying cell-types and tissue domains in 2D and 3D. SSAM is applicable to a variety of *in*

38 *situ* transcriptomics techniques and capable of integrating prior knowledge of cell types. We

39 apply SSAM to three mouse brain tissue images: the somatosensory cortex imaged by

40 osmFISH, the hypothalamic preoptic region by MERFISH, and the visual cortex by multiplexed

41 smFISH. We found that SSAM detects regions occupied by known cell types that were

42 previously missed and discovers new cell types.

43
44 **Keywords**

47 **Introduction**

48 The underlying transcriptional and spatial heterogeneity of cells gives rise to the plethora of

49 phenotypes observed in cell types, tissues, organs, and organisms. Recent technological

50 advances[1] have seen the profound adoption of single-cell sequencing to unravel transcriptional

51 heterogeneity in healthy and diseased tissues, and have subsequently given rise to international

52 consortia such as the Human Cell Atlas (HCA)[2]. Such efforts would not be possible without

53 computational frameworks supporting the analysis of single-cell sequencing data[3]. Linking this

54 transcriptional heterogeneity with spatial heterogeneity of cells is a critical factor in

55 understanding cell identity in the context of the tissue, for example, revealing the transcriptional

56 basis of invasive cancer regions[4] and highlighting the rich diversity of neuronal subtype

57 expression and localization[5]. Recently developed multiplexed fluorescence in-situ hybridization[6–

58 8] and *in situ* mRNA tissue sequencing techniques[9–14] have enabled the simultaneous

59 measurement of multiple mRNAs in a spatial context.

60

61 Traditionally, mRNA molecules identified by *in situ* transcriptomics are assigned to cells and

62 subsequently used for computing gene expression profiles of those cells[15–18]. Identification of

63 cells relies on cell segmentation, a procedure demarcating the interior and exterior of the cell

64 membranes, which relies on additional signals or landmarks obtained by staining nuclei[19], cell

65 membrane[20–22], or total poly-A RNA[5,6]. However, accurate cell segmentation is difficult to

66 achieve with current techniques due to tightly apposed or overlapping cells, uneven cell borders,

67 varying cell and nuclear shapes, signal intensity variation, probe fluorescence emission

68 efficiency variation, and tiling artifacts[23]. Such obstacles can result in detecting fewer cells or

69 incorrect cell borders. Subsequent analysis would then be spatially restricted to inaccurately

70 segmented cells and may mean that large portions of meaningful mRNA signals are discarded.

71 This may result in incorrect cell-type signatures, incomplete cell-type maps, or missing rare cell

72 types. Therefore, there is a need for robust cell segmentation-independent methods for

73 identifying cell-type signatures, cell-type organization, and tissue domains from

74 multidimensional mRNA expression data in complex tissues. These methods could be used for

75 datasets lacking landmarks or to validate segmentation-based approaches.

76

77 Here we introduce a novel computational framework named Spot-based Spatial cell-type

78 Analysis by Multidimensional mRNA density estimation (SSAM). In contrast to existing methods,

79 SSAM departs from the spatial restriction of approaches based on cell segmentation and

80 instead identifies cell types using mRNA signals in the image, without the need for prior cell

81 segmentation. Furthermore, instead of labelling only segmented regions, our approach assigns

82 cell-type labels to each pixel, ensuring a more complete picture of cell-type specific spatial

83 heterogeneity.

84

85 We apply SSAM to three mouse brain tissue images obtained by different techniques: the

86 somatosensory cortex (SSp) by osmFISH, the hypothalamic preoptic region (POA) by

87 MERFISH, and the visual cortex (VISp) by multiplexed smFISH. With all three datasets, we

88 demonstrate the robustness of SSAM in identifying 1) cell types *in situ*, 2) spatial distribution of

89 cell types, 3) spatial relationships between cell types, and 4) tissue domains (e.g., cortical layers)

90 based on the local composition of cell types without fine-tuning of parameters. We demonstrate

91 that SSAM 1) correctly identifies the spatial distribution of known cell types in regions missed in

92 the SSp by cell segmentation based methods for the osmFISH data ; 2) can analyze the POA

93 MERFISH 3D data using the same parameters as for the 2D SSp osmFISH data without any

94 extra adjustments of the settings; 3) identifies new and rare cell types in the VISp, multiplexed

95 smFISH data.

96

## Results

**The SSAM computational framework**

SSAM consists of 4 major steps (**Fig. 1**), namely 1) mRNA signal estimation and downsampling; 2) computation of cell-type signatures; 3) generation of a cell-type map; and 4) identification of tissue domains.

In the first step, SSAM estimates mRNA signal intensity over the tissue image (**Fig. 1A**). Firstly, for each gene, mRNA signal intensity distribution is estimated by applying a Kernel Density Estimation (KDE) with a Gaussian kernel, which is then resolved to pixels in the image. The mRNA signal intensity distribution for each gene is stacked to create a gene expression vector field, which is a multichannel image where the pixels encode the expected density of mRNA count for each gene. This essentially assigns gene expression profiles to pixels in the image.

In the second step, SSAM identifies cell-type gene expression signatures by clustering (**Fig. 1B**). Before running the clustering algorithm, SSAM downsamples gene expression vectors to reduce computational processing time. As default, SSAM performs informed downsampling by selecting pixels that are local maxima in the gene vector field (Methods). After that, both the downsampled vectors and the gene expression vector field are normalized (Methods). SSAM clusters the sampled vectors using either DBSCAN[24], HDBSCAN[25], OPTICS[26], or the Louvain community detection method implemented in Seurat[27] (Methods). The Louvain methods is the default as it has been widely utilized to analyze single cell data. After the clustering step, sampled vectors with a large distance in gene expression space to their cluster medoid are removed as outliers to ensure the quality of selected vectors (**Supplementary Fig. 1B**). The gene expression cluster centroids are used to represent the gene expression signature of a cell type.

6

123    In the third step, SSAM classifies each pixel in the image to create a "cell-type map" (**Fig. 1C,**

124    **Supplementary Fig. 2A**). SSAM includes a guided mode, which assigns pixels to a labeled set

125    of given gene expression signatures (e.g. from scRNA-seq/segmentation), as well as a *de novo*

126    mode, which assigns pixels to the cell type signatures obtained in the previous clustering step.

127    For the classification of pixels, SSAM first creates signature prototypes by averaging the

128    signatures per cell-type class of the given signatures, then it classifies all spots in the vector

129    field according to the maximum correlation to any of the signature prototypes.

130

131    In the fourth step, SSAM identifies tissue domains that have distinct cell-type composition (**Fig.**

132    **1D**). SSAM computes the cell-type compositions in a circular (or spherical) sliding window over

133    the cell-type map and clusters the cell-type composition of each window using agglomerative

134    hierarchical clustering (**Supplementary Fig. 2B**). The resultant clusters represent putative

135    tissue domains. Clusters with high mutual correlation are then merged into a single tissue

136    domain signature, and the cell-type composition of each domain is calculated.

137

138    In the following sections we apply SSAM to three multiplexed FISH datasets obtained using

139    different techniques. We reanalyze two previously published datasets, profiled by osmFISH[6] and

140    MERFISH[5], to demonstrate SSAM's strength in comparison to earlier methods. For a newly

141    generated multiplexed smFISH dataset we demonstrate that SSAM can unravel novel biological

142    insights into the spatial cellular organization of the brain.

143

144    **SSAM improves astrocyte and ventricle detection in the mouse brain somatosensory**

145    **cortex (SSp)**

146    To demonstrate the utility of SSAM, we analyzed published osmFISH data, where the

147    transcripts of 33 cell-type marker genes were localized in 2D space of the mouse brain

148    somatosensory cortex (SSp)[6] (**Fig. 2, 3, Supplementary Fig. 3, 4**). We compare results

149 obtained from SSAM against the results obtained from Poly-A segmentation from the original

150 study.

151

152 The osmFISH dataset was first analyzed using the guided mode of SSAM. Cell-type maps were

153 generated using cell-type signatures from the prior segmentation-based approach[6] and another

154 from scRNA-seq[28,29] (**Supplementary Fig. 4E**).

155

156 To quantify the similarity between the prior segmentation and the cell-type maps generated by

157 SSAM, we calculated a "matching score" for each cell type (Methods). The matching scores

158 between the segmentation from the previous study and SSAM guided by both segmentation-

159 based and scRNA-seq cell-type signatures were generally high (mean and median matching

160 score of 0.67 and 0.78 for segmentation-based, 0.60 and 0.70 for scRNA-seq-based signatures,

161 respectively), indicating a strong agreement of the two cell-type maps as visually apparent

162 (**Supplementary Table 1, 2, Supplementary Fig. 5, 6**).

163

164 Next, we continued with completely *de novo* cell-type identification. The resulting 30 cell-type

165 signatures (**Fig. 2A, B, Supplementary Fig. 7-10**) were consistent with those identified in the

166 segmentation-based clustering and scRNA-seq based cell-type signatures[6] (**Supplementary**

167 **Fig. 4C, D**), implicating the robustness of the *de novo* cell-type calling by SSAM. Each of the

168 SSAM *de novo* cell-type signature clusters were assigned the label of the closest correlating

169 segmentation-based cluster.

170

171 As with the guided mode analysis, we limit the comparison to the most comparable cell types,

172 excluding cell types with low correlation in gene expression signatures (< 0.8) (**Supplementary**

173 **Table 3, Supplementary Fig. 11**). The matching score result showed high average values

174 (mean and median of 0.76 and 0.83, respectively) and 81% of cell types had a matching score

8

175    of greater than 0.6. Comparing marker gene expression of cell types having lowest matching

176    score (< 0.3) (**Supplementary Table 3**) confirmed that the SSAM guided cell-type map is in

177    better agreement to their marker gene expression (**Supplementary Fig. 12-13**). Given the low

178    correlation of C. Plexus cell type to the corresponding osmFISH cluster, which is one of the

179    dominant cell types in the ventricle region, high-resolution investigation of Poly-A and DAPI

180    signals confirm the existence of both cell types in the ventricle area (**Fig. 2D**). Since ependymal

181    and choroid plexus cells were small and tightly packed and exhibit relatively lower DAPI and

182    poly-A signal, we concluded that the performance of the watershed algorithm was insufficient to

183    identify cells in the area. Furthermore, we statistically evaluated this for each cell type by

184    comparing the gene expression in the unique parts of the segmentation and SSAM *de novo* cell-

185    type map, to the overlapping parts (Methods). Gene expression of the unique part of SSAM *de*

186    *novo* cell-type map showed higher correlation to the overlapping regions compared to the

187    unique parts of the segmentation (**Supplementary Fig. 14**).

188

189    We then performed domain analysis on the SSAM *de novo* cell-type map. Identified domains

190    correlated well with the known cerebral cortex layers, consistent with results reported in the

191    previous study (**Fig. 3A**). Laminar distribution of cell types is established [30], and can be

192    considered as a ground truth for validating the cell type map. Cell-type assignments of

193    excitatory pyramidal cells in the cortical layers conformed closely to known localizations

194    (**Supplementary Fig. 15**). The domains identified as: layer 2/3 primarily consists of Pyramidal

195    L2-3/L5, L2-3, and L3-4 cell types; layer 4 consists of Pyramidal L4 and L3-4 cell types; layer 5

196    consists of Pyramidal L3-5 and L5 cell types; and layer 6 consists of Pyramidal L6 cell types.

197

198    In addition, cell-type composition of the domains revealed that *Mfge8* expressing astrocytes

199    (Astrocyte Mfge8) contributed 7-14 % of each of the tissue layers (**Fig. 3B**), in contrast to the

200   significantly fewer numbers of Astrocyte Mfge8 cells called in the previous study[6]. Comparison

201   of high-resolution images of DAPI and poly-A signals with Mfge8 expression densities implicates

202   that the poly-A signal was not strong enough to discriminate the presence of astrocyte Mfge8

203   cells from the background, while the DAPI images clearly supported the existence of *Mfge8*

204   expressing astrocytes at positions identified by SSAM (**Fig. 2E**). The clear DAPI signal but low

205   poly-A signal for these astrocytes Mfge8 suggested that they have a lower mRNA content

206   compared to other cells. We compared the total counts of mRNA molecules of astrocytes and

207   other cell types from mouse brain scRNA-seq data[31] and found that astrocytes exhibited

208   significantly less mRNA molecules than other cell classes (**Supplementary Fig. 4B**). Our

209   observation reveals the inadequacy of the watershed segmentation algorithm applied to poly-A

210   signal when not considering cells with a low total mRNA content. This implies that the original

211   segmentation of these cell types could be less accurate than the SSAM *de novo* cell-type map,

212   therefore also reducing the matching score for these cell types.

213

214   **SSAM confirms diversity of inhibitory and excitatory neuron cell types and their**

215   **localization in the hypothalamic preoptic region (POA) in 3D**

216   To demonstrate the performance of SSAM for three-dimensional *in situ* transcriptomics data, we

217   applied SSAM to previously published MERFISH data, where 135 transcripts were localized in

218   3D space of the hypothalamic preoptic region (POA) of a mouse brain[5] (**Fig. 4, Supplementary**

219   **Fig. 16, 18**). We compare results obtained from SSAM against the results obtained from DAPI

220   segmentation from the original study.

221

222   We applied both SSAM guided mode and *de novo* mode. For guided mode, the previously

223   known cell-type signatures obtained by segmentation and scRNA-seq were used. For both

224   guided and *de novo* modes, SSAM analysis was performed in 3D space, generating a 3D cell-

225   type map (**Fig. 4B**). The resulting cell-type maps on the x-y plane at the center of slice on the z-

226   axis (at 5μm) were visually similar to the previous study (**Supplementary Fig. 17G**). SSAM cell-

227   type signatures showed high expression of their marker genes (**Supplementary Fig. 18-21**) and

228   a high correlation to the cell-type signatures from both the segmentation-based clusters and

229   scRNA-seq clusters (**Supplementary Fig. 17E, F**). Among them, 7 inhibitory and 4 excitatory

230   neuronal cell types showed very high correlation (>0.8) to the segmentation-based neuronal

231   signatures, and also showed distinctive tissue localization patterns (**Fig. 4D, E**), similar to those

232   previously reported (**Supplementary Fig. 22**).

233

234   We then quantified the similarity of the SSAM cell-type maps with the cell segmentation by

235   Moffitt et al. The SSAM guided mode cell-type map achieved high matching scores for

236   comparable cell types (mean and median of 0.76 and 0.83 for segmentation-based, 0.88 and

237   0.94 for scRNA-seq-based signatures, respectively), with only 6 of 76 cell-types exhibiting a low

238   matching score (< 0.3) for segmentation-based case (**Supplementary Table 4, 5,**

239   **Supplementary Fig. 23, 24**). Comparing the SSAM *de novo* cell-type map also yielded high

240   matching scores (mean and median of 0.83 and 0.93, respectively) (**Supplementary Table 6,**

241   **Supplementary Fig. 25**), further validating the computational approach adopted by SSAM to

242   identify *de novo* cell-type signatures and generating cell-type maps. One of the most notable

243   differences in the SSAM cell-type map was that we found a higher density of astrocytes

244   compared to Moffit et al. A comparative analysis revealed that some astrocyte signals identified

245   by SSAM were not found in the segmentation by Moffit et al. Note that the existence of

246   astrocytes is clearly shown by the corresponding marker gene expression (**Supplementary Fig.**

247   **26**).

248

249   The generated tissue domain map identifies several domains consisting of regions consisting

250   primarily of inhibitory neurons, excitatory neurons and oligodendrocytes, as well as the ventricle

251     structure (**Supplementary Fig. 27**).

252

253     Finally, we reconstructed a three-dimensional cell-type map (**Movie 1**). While the thickness of

254     the tissue image is limited (10 μm), we demonstrate the shape and size difference of the whole

255     cell-type map and the cell-type specific maps for inhibitory neurons, excitatory neurons and

256     astrocytes (**Movies 2, 3, 4**).

257

258     Despite the difference of dimensionality between the osmFISH data (2D) and the MERFISH

259     data (3D), SSAM was able to successfully process the data and produce meaningful results.

260     More importantly, the analyzes in this section were performed with almost the same procedure

261     and parameters applied to the osmFISH data. Therefore, we set these parameters as the

262     default values to facilitate rapid and robust analysis of other multidimensional *in situ*

263     transcriptomics dataset using SSAM.

264

265     **SSAM identifies rare cell types and novel cortical sub-layering in the adult mouse visual**

266     **cortex (VISp)**

267     To further demonstrate that SSAM can be used for rapid and robust analysis of *in situ*

268     transcriptomics data, we applied SSAM to unpublished multiplexed smFISH data of the mouse

269     primary visual cortex (VISp) generated as part of the SpaceTx consortium[32] (**Fig. 5, 6,**

270     **Supplementary Fig. 28, 29**). In total, the expression of 22 genes was quantified *in situ*

271     (Methods).

272

273     Analysis of the tissue image was restricted to the manually defined VISp region

274     (**Supplementary Fig. 28D**). SSAM was performed in both guided mode and *de novo* mode

275     (**Supplementary Fig. 29A**). The guided mode of SSAM was performed using scRNA-seq data[30].

276     For the *de novo* run, the identified cell-type signature clusters were assigned the label of the

12

277    cluster in the scRNA-seq data with the highest correlation (**Fig. 5A, B**). Then, the tissue

278    domains were identified based on the *de novo* cell-type map (**Fig. 6**), with the result showing the

279    laminar structure of the VISp region. We identified two distinct layer 4 (L4) neuronal clusters.

280    Interestingly, both of them showed the highest correlation to the single L4 IT type identified via

281    scRNA-seq, but their spatial locations show a clear difference (**Fig. 5C, Supplementary Fig.**

282    **29B**). We named the cluster localizing to the superficial region of layer L4 as 'L4 IT Superficial'

283    (L4 IT 2). This finding adds context to the previously observed heterogeneity of the L4 IT cell

284    type[30], where we show that this heterogeneity determines superficial and deep localization in

285    layer 4.

286

287    The cell-type map generated by SSAM guided mode were visually similar to that of *de novo*

288    mode, except for the cell types found in the layer 2 (L2) (**Supplementary Fig. 29A**). We found

289    that the majority of cell types found in L2 were assigned to the VLMC type in SSAM guided

290    mode. We observed that this type was actually a neuronal type in L2. This cell type showed high

291    expression of *Alcam*, a marker gene of the VLMC cell type, but low expression of other genes.

292    Due to the limited number of genes profiled in the multiplexed smFISH experiment, lack of other

293    neuronal marker genes led to incorrect high correlation of this type VLMC. However, SSAM

294    properly assigned the centroid to be L2 neurons in *de novo* mode.

295

296    SSAM was also able to identify a rare cell type, Sst Chodl, which is known to be related to long-

297    range projection and sleep-active neurons[33–35]. In addition, we mapped the Sst Chodl cell-type

298    signal to between layer L5 and L6 (**Supplementary Fig. 29C**), consistent with previously

299    reported localization to L5 and L6[33]. This finding was validated against its marker gene

300    expression (**Supplementary Fig. 30-32**), and ultimately demonstrates SSAMs ability to identify

301    cell-type signatures of lowly abundant and rare cell-types.

302

13

303 **Discussion**

304 We describe a segmentation-free computational framework for processing *in situ*

305 transcriptomics data and demonstrate its performance on three different adult mouse brain

306 datasets: the somatosensory cortex (SSp) profiled by osmFISH, the hypothalamic preoptic

307 region (POA) by MERFISH, and the visual sensory cortex (VISp) by multiplexed smFISH. We

308 find that the cell-type signatures and maps generated by SSAM for both osmFISH and

309 MERFISH datasets were similar to the previously reported ones, validating the underlying

310 methodology of SSAM. Based on this, we successfully determined cell types and constructed

311 cell-type and tissue domain maps in the multiplexed smFISH mouse VISp dataset.

312

313 In the osmFISH dataset our method outperforms the original segmentation-based cell-type map

314 reconstruction due to limitations in the segmentation process. In the MERFISH dataset we show

315 that SSAM is able to identify diverse populations of cell types and that SSAM is scalable to 3D

316 image data. For the VISp multiplexed smFISH dataset, SSAM identified a rare cell type and

317 elucidated a suspected spatial heterogeneity of cell types in the cortex without segmenting a

318 single cell. Overall, the results show that SSAM is not only a robust tool to validate

319 segmentation-based methods, but also a reasonable alternative when segmentation is difficult

320 or DAPI or Poly-A images are lacking.

321

322 However, for some questions it is important to distinguish between cells to e.g. delineate growth

323 arising from increasing cell size vs cell proliferation or to investigate multinucleation in

324 cardiomyocytes or cytotrophoblast cells. In cases such as these, we recommend the use of

325 SSAM as a complementary method to segmentation-based analysis in two ways. First, the

326 output of SSAM can be compared to validate that the segmentation process did not introduce

327 artifacts. Secondly, to use the SSAM output as an input for the segmentation process to refine

328 the segmentation procedure for different domains or cell-type signals.

14

329

330    In terms of methodological parsimony, SSAM minimizes the number of assumptions, avoids

331    iterative optimization and thus offers maximal transparency, interpretability and reproducibility.

332    The lightweight nature of the algorithm typically brings a considerable runtime advantage over

333    other available packages. SSAM is written as a Python library, with some core analysis

334    functions wrapped up with external C functions to speed up the computation. The package is

335    available as an easily installable Python package, and can easily be extended with existing *in*

336    *situ* transcriptomics pipelines, e.g. starfish (https://github.com/spacetx/starfish) or Giotto[36].

337    SSAM is accompanied with a notebook outlining all the steps presented in this paper. Taken

338    together, we present a novel, flexible and robust method for fully automated cell-type and tissue

339    domain analysis that is readily applicable to numerous *in situ* transcriptomics methods.

340

341    **Materials and Methods**

342
343    **Using Kernel Density Estimation to generate the gene expression vector field**

344    We used the n-dimensional KDE algorithm to estimate the density of mRNAs in 2D and 3D. To

345    compute Gaussian KDE, we used our own implementation of the KDE algorithm for rapid

346    computation. Spatial distribution of the probability of mRNA presence $\hat{p}$ is estimated using the

347    kernel density estimation;

348

$$\hat{p}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} \kappa_h (\mathbf{x} - \mathbf{x_i})$$

349

350

351    where:

352       - $\kappa_h$: a kernel function with a bandwidth $h$

15

353      -   $N$: the number of data points

354      -   $\mathbf{x}_i$: location vector of the data point i (i.e. location of i-th mRNA)

355

356   Here we use the Gaussian kernel:

357

$$\kappa_h(\mathbf{x}) = \frac{1}{(2\pi h^2)^{d/2}} \; e^{-\frac{1}{2}\|\mathbf{x}\|^2/h^2}$$

358

359

360   where:

361      -   $h$: bandwidth of the Gaussian kernel

362      -   $d$: dimension of the space where the data points reside (2 for 2D, or 3 for 3D mRNA

363       locations)

364      -   $\|\mathbf{x}\|$: Euclidean norm (i.e. L2 norm) of vector $\mathbf{x}$

365   Note that the integration of $\hat{p}(\mathbf{x})$ all over the space is 1. Therefore the gene expression density

366   is calculated by multiplying the number of mRNAs per gene to $\hat{p}$.

367

368   **Calculation of spatial gene expression**

369   The continuous estimation of gene expression density is discretized over pixels of the tissue

370   image, which in our examples is set to a size of 1□m. The expectation value of the estimated

371   density in a unit pixel is approximated by multiplying the area of the unit pixel to the estimated

372   gene expression density at the location of the pixel. Finally, we stack the estimated gene

373   expression densities of genes to define the gene expression vector field over the image.

374

375   **Selection of local maxima**

376   Local maxima were selected based on the L1-norm of the vectors in the vector field, which is

16

377    the total size of each vector in the image. For the selection algorithm, we used scikit-image

378    Python package to select local maxima. Briefly, 1) maximum filter is applied to dilate the original

379    image, 2) the locations where the maximum filtered image equal to the original image are

380    selected. The maximum filter with size 3 was used throughout the examples presented in this

381    paper.

382
383    **Downsampling of the vector field**

384    For a scalable cell-type identification analysis, the vector field is downsampled to a smaller set

385    of vectors based on local maxima selection strategy (Supplementary discussion). SSAM applies

386    two thresholds for local maxima selection: 1) a minimum expression threshold for a single gene

387    defined as the height of a single Gaussian kernel to avoid regions with signal from only the

388    Gaussian tail (see Discussion section for details), which also corresponds to the position of the

389    observable drop in the histograms of gene expression (**Supplementary Fig. 3A, 17A, 28A**); 2)

390    a minimum total gene expression (i.e. L1-norm) threshold (**Supplementary Fig. 3B, 17B, 28B**).

391    Furthermore, we implemented an optional "input mask" feature to limit sampling of vectors to

392    regions of the image containing informative data, e.g. a mask outlining the informative tissue

393    area.

394

395    **Comparison of local maxima and random sampling strategies**

396    The two local maxima sampling methods, 1) local maxima sampling and 2) random

397    downsampling, were compared to justify our preference of local maxima sampling method for

398    the downstream analysis. The osmFISH data was used for the comparison. Firstly 11,469 local

399    maxima vectors were found in the vector field using a window size of 3, a minimal gene

400    expression and L1 norm thresholding. For comparison, the same number of vectors were

401    randomly sampled from the vector field, using the same thresholds used for local maxima

402    selection. At the locations of the vectors, both the local maxima and the random sampled

403    locations, the classified cell types on the cell-type map guided by segmentation-based

404    signatures are called. For each case, the Pearson's correlation coefficients between the vectors

405    and the signature of the cell types are calculated and plotted as a distribution (**Supplementary**

406    **Fig. 39**).

407

**Variance stabilization of local maxima vectors and the vector field**

409    Since the gene expression profiles of local maxima vectors are representative of the

410    transcriptomes of cells, we considered them to be analogous to the gene expression count

411    matrix obtained from single cell RNA sequencing (scRNA-seq) using unique molecular

412    identifiers (UMI). Therefore, we normalized the local maxima vectors of the vector field (which

413    would be representative of single cells) using *sctransform*[37], a normalization and regularization

414    algorithm for UMI count data. After that, each vector of the vector field is normalized using

415    *sctransform*, with the same parameters previously used to normalize the local maxima.

416

**Clustering of representative gene expression vectors**

418    The SSAM framework supports clustering via DBSCAN[24], HDBSCAN[25], OPTICS[26] and an

419    implementation of the Louvain algorithm equivalent to that in the R package, Seurat[27]. DBSCAN,

420    HDBSCAN and OPTICS are implemented via the scikit-learn Python library. The Louvain

421    clustering algorithm is based on the R package Seurat[27] reimplemented in Python. In short, an

422    SNN network with correlation metric is built using a python package NetworkX[38]. The weight of

423    the network is calculated by a Jaccard similarity coefficient. A weight smaller than 1/15 was set

424    to zero. Clustering was done by detecting communities in the network using a Louvain

425    community detection algorithm implemented in Python (python-louvain, https://python-

426    louvain.readthedocs.io/). It is known that the Louvain algorithm is not sensitive in detecting small

427    clusters[39], optionally DBSCAN algorithm can be applied to subcluster each Louvain cluster. This

428    sub-clustering strategy is conceptually similar to the "Polished Louvain" algorithm in Zeisel et

18

429    al[31].

430

**Diagnostic plots**

432    After unsupervised clustering of gene expression vectors, some clusters may need to be

433    manually merged or discarded. SSAM supports merging of clusters based on correlation of

434    gene expression profile, however in many cases manual inspection is needed to rule out any

435    non-trivial issues. To guide this process, SSAM generates a cluster-wise 'diagnostic plot', which

436    consists of four panels: 1) location of the clustered vectors on the tissue image, 2) the pixels

437    classified to belong the cluster signature (the cluster centroid), 3) the mean expression profile of

438    the clustered vectors, and 4) the t-SNE or UMAP embedding.

439

440    In the three datasets analyzed the clusters to be merged or removed often showed a

441    discordance between the location of sampled vectors used to determine the cluster (panel 1)

442    and the pixels classified to belong to that cluster (panel 2). In case of overclustering, i.e. when a

443    cell-type signature is split over 2 clusters, the map typically does not classify the full shape of

444    the cells but instead only fragments (panel 2), and having almost the same marker gene

445    expression of another cluster (panel 3). Such clusters can be merged. For dubious clusters that

446    should be removed, we observed that vectors usually originate from outside the tissue region or

447    from image artifacts (panel 1), or that the gene expression does not show any clear expression

448    of marker genes or similarity to expected gene expression profiles (panel 3).

449    The remaining clusters are then annotated by comparing cluster marker genes to known cell-

450    type markers. Note that in many cases, the identity of clusters can be easily assigned by

451    comparing the centroids of the clusters to the known cell-type signatures, e.g., from single cell

452    RNA sequencing. To support rapid annotation of cell types to clusters, SSAM additionally shows

453    the highest correlating known cell-type signature should this data be available in panel 3. The

454    diagnostic plots for osmFISH, MERFISH, and multiplexed smFISH data are available online in

19

455    the Jupyter notebook uploaded to zenodo (http://doi.org/10.5281/zenodo.3478502).

456

457    **Statistical evaluation of cell-type mapping**

458    The accuracy of the SSAM cell-type map was validated by comparing the published osmFISH

459    segmentation and the SSAM *de novo* cell-type map by two different methods.

460

461    Firstly, to quantitatively compare concordance of cell-type we implemented a matching score.

462    The matching score for any given cell type is defined as the number of segmented cells with at

463    least 10% of matched with the SSAM guided or *de novo* mode cell type map of the

464    corresponding cell type of the segment, divided by the total number of segments of the cell type

465    which represents the ratio of segments identified by SSAM. The threshold of 10% was

466    empirically selected to account for differences in cell location in the tissue, especially for very

467    small cells where subtle changes in cell-type labeling can drastically reduce the overlap within

468    the segmented area.

469

470    Secondly, for evaluation of discrepancies in cell-type locations compared to the original studies,

471    we compare the unique part of each segmentation and SSAM *de novo* cell-type map to the

472    parts that are overlapping in both maps in the osmFISH dataset. The gene expression vectors

473    originating from overlapping parts of the same cell types (**Supplementary Table 3**), were

474    regarded as the ground truth set. Then, two sets of unique vectors were defined: 1) the

475    segmentation-only set, the vectors from the regions occupied by segments excluding the

476    overlap, and 2) the SSAM-only set, the vectors from SSAM cell-type map only regions. The

477    distribution of the gene expression vectors in the overlapping set was then compared to the two

478    unique parts (**Supplementary Fig. 14A**). To compare the accuracy of cell-type mapping of the

479    two unique parts, Pearson's correlation coefficient is calculated between the mean expression

20

480    of the ground truth set and the vectors in each set (**Supplementary Fig. 14B**).

481

**Quantification of doublets**

483    The doublet rates were evaluated by two Python packages, DoubletDetection[40] and Scrublet[41]

484    (**Supplementary Table 8**). As the two algorithms require raw counts as input, the unnormalized

485    raw vectors at local maxima used for clustering analysis were used as input of the two

486    algorithms, as an analogy of the raw counts. For DoubletDetection, the doublet rate was

487    calculated by dividing the number of doublets reported by the number of total local maxima. The

488    doublet rate quantification by both methods was consistent, and negligible in the osmFISH and

489    multiplexed smFISH datasets (average doublet rate of <0.5% for both), and marginal for

490    MERFISH (average doublet rate of 3%).

491

**SSAM analysis of osmFISH data**

493    KDE was performed with a bandwidth of 2.5 µm. The individual gene expression threshold and

494    total gene expression threshold for selection of local maxima were 0.027 (the height of a single

495    Gaussian) and 0.04, respectively (**Supplementary Fig. 3A, 3B**). Since the selected local

496    maxima includes many locations outside of the tissue area, we further filtered local maxima

497    based on their local density approximated using the k-nearest neighbor algorithm. More

498    specifically, local maxima with a density lower than 0.002 over the closest 100 local maxima,

499    corresponding to fewer than 100 local maxima in a 126.2 µm radius, were filtered out

500    (**Supplementary Fig. 3C**). The selected local maxima vectors were passed to *sctransform* to

501    determine normalization parameters, after which the whole vector field was normalized.

502

503    In SSAM guided mode, the mRNA count matrix of both the previously segmented cells and the

504    scRNA-seq data were normalized by *sctransform*. The centroid of each of the annotated

505    clusters was used to classify cell types in the vector field, generating a cell-type map guided by

506    prior knowledge.

507

508    In SSAM *de novo* mode, the selected local maxima vectors were clustered using the Louvain

509    algorithm with a resolution of 0.15, resulting in 66 clusters (**Supplementary Fig. 4A**). Distinct

510    clusters representing the same cell types were identified and then manually merged, and

511    spurious clusters were removed, resulting in a total of 30 clusters (**Fig. 2A, 2B**). For each

512    cluster, the vectors with insufficient correlation to its cluster medoid were excluded from the

513    centroid calculation (**Supplementary Fig. 1B**). The cluster centroids were compared to that of

514    the segmentation-based (**Supplementary Fig. 4B**) and scRNA-seq cell-type signatures

515    (**Supplementary Fig. 4C**) using Pearson's correlation coefficient. The *de novo* clusters were

516    named after the highest correlating segmentation-based cluster. Note that clusters closest

517    mapped to Inhibitory IC and Inhibitory CP cell types do not only appear in the internal capsule

518    and caudoputamen, but also in the cortex. Therefore, we renamed these clusters to Inhibitory

519    Kcnip2 (since Kcnip2 was the third most expressed gene for this cluster) and Inhibitory Rest,

520    respectively. After classification of the local maxima, we quantified the doublet rates (Methods,

521    **Supplementary Table 8**).

522

523    Tissue domain analysis was performed using a sliding circular window with radius 100 µm with

524    a step of 10 µm. The cell-type proportions from each window were clustered using

525    agglomerative hierarchical clustering with 15 clusters as an initial estimate, subsequently

526    merging the clusters with correlation coefficients higher than 0.8. Spatially connected clusters

527    with a correlation coefficient higher than 0.6 were merged. The resulting domain map was

528    resized to match the size of the cell-type map, after which the cells in different domains were

529    colored.

530

22

531 **Quantification of mRNA abundance in astrocytes and other brain cell types for osmFISH**

532 **data interpretation**

533 The "L5_All.loom" loom object containing scRNA-seq expression data of half a million cells from

534 the mouse nervous system[31] was downloaded (http://mousebrain.org/downloads.html). The total

535 number of mRNA molecules per cell were extracted and aggregated by their level 2 class labels

536 (astrocytes, immune, vascular, ependymal, neuronal, peripheral glia and oligodendrocyte cells)

537 using Python. The counts were log normalized and subsequently followed a normal distribution

538 (tested using the Shapiro-Wilk test for normality, all *p-values* < 1 x 10e-4 for each class),

539 therefore a Student's t-test was applicable. For each of the two classes of interest

540 ('Astrocytes', 'Immune'), we performed independent log-space t-tests for unequal sample sizes

541 and unequal variance against each of the other classes. Both astrocyte and immune cell

542 classes have significantly lower mRNA molecule counts compared to other cell types (all *p-*

543 *values* < 1 x 10e-12). While the distribution of mRNA counts in log space followed a normal

544 distribution, the use of a Student's t-test for large numbers may be not appropriate. Hence, we

545 also describe the difference in their distributions. For both astrocyte and immune cell classes,

546 more than half of the cells of each class exhibited a lower UMI count than the lowest quartile of

547 any other cell class.

548

549 **SSAM analysis of MERFISH data**

550 KDE was performed with bandwidth 2.5 μm. Local maxima were filtered using a gene

551 expression threshold of 0.0055, and then filtered with total gene expression threshold of 0.0035

552 (**Supplementary Fig. 17A, B**). The selected local maxima vectors were passed to *sctransform*

553 to determine normalization parameters, after which the whole vector field was normalized.

554

555 In SSAM guided mode, the mRNA count matrix of both the previously segmented cells and the

556 scRNA-seq data were normalized by *sctransform*. The centroid of each of the annotated

23

557   clusters was used to classify cell types in the vector field, generating a cell-type map guided by

558   prior knowledge.

559

560   For SSAM *de novo* mode, the selected vectors were clustered using the Louvain algorithm with

561   a resolution of 0.15, resulting in 68 clusters (**Supplementary Fig. 17C**). By manual inspection of

562   gene expression and localization, overclustering was merged, and spurious clusters were

563   removed, resulting in a total of 50 clusters (**Fig. 2A, 2B**). For each cluster, the vectors that did

564   not have high correlation to its cluster medoid were excluded from the centroid calculation

565   (**Supplementary Fig. 1B**). The centroids of the clusters are compared with that of the

566   segmentation-based clustering result and scRNA-seq result using Pearson's correlation

567   coefficient (**Supplementary Fig. 17E, F**). The SSAM *de novo* clusters correlating best to

568   inhibitory and excitatory neurons were named based on the most highly expressed gene of each

569   cluster, and the non-neuronal clusters were named based on the previous study[5]. After

570   classification of the local maxima, we quantified the doublet rates (Methods, **Supplementary**

571   **Table 8**). We noticed a number of small blobs on the cell type map, which are resultant from

572   cells on a different plane in the 3D image (**Movie 2**). After classification of the local maxima, we

573   quantified the doublet rates (Methods, **Supplementary Table 8**).

574

575   Tissue domain analysis based on the cell-type map was performed using a sliding spherical

576   window with radius 100 μm with a step of 10 μm. The cell-type proportions from each window

577   were clustered using agglomerative hierarchical clustering with 20 clusters as an initial estimate,

578   subsequently merging the clusters with correlation coefficient higher than 0.8. The resulting

579   domain map was resized to match the size of the cell-type map, after which the cells in different

580   domains were colored.

581

**Comparison of localization of inhibitory and excitatory neurons**

For a number of inhibitory and excitatory neuronal subtypes identified in the posterior POA

tissue image using SSAM *de novo* mode, we identified the best matching cell types based on

Pearson correlation of their gene expression signatures (**Supplementary Fig. 17F).** We

matched the following cell types: SSAM cluster 39 (C39) called Inhibitory Coch to Moffitt cluster

I-12, C16 Inhibitory Arhgap36 to I-13, C45 Inhibitory Isr4 to I-15, C34 Inhibitory Calcr to I-14 ,

C14 Inhibitory Gda to I-23, C19 Excitatory Cbln1-Cbln2 to E-19, C42 Excitatory Omp to E-16,

C25 Excitatory Necab1-Gda to E-9, C8 Excitatory Necab1 to E-14, and C36 Excitatory Col25a1

to E-24. For these cell types we checked the tissue localizations reported in the previous studies

figures 5a, 5c, 5e, 6b, 6d, and S17[5]. Side-by-side comparison of the localization of these

neuronal cell types revealed very similar patterns of localization computed by SSAM and the

original publication (**Supplementary Fig. 22**).


**3D modelling of MERFISH cell-type maps**

Firstly, the connected components in 3D were determined using the python package connected-

components-3d (https://github.com/seung-lab/connected-components-3d). Components

comprising fewer than 100 voxels were removed. After this, the voxels filling connected

components were removed, and only the contours were used for the vertex of the 3D models.

For each vertex, the vertex normal was calculated by simple physics simulation, assuming that

the direction of a vertex normal vector is the same as the force vector when there are pulling

forces between all of the contour voxels. The surface of the objects was reconstructed using

screened Poisson reconstruction algorithm[42,43] using default parameters. The number of

vertices was reduced to 5% of the total number of vertices using the 'vtkQuadricDecimation'

function[44,45] of VTK library[46]. Finally, the objects were merged into a single file. Each scene of

the rotating movie was created using Meshlab[47].

608 **VISP multiplexed smFISH data generation**

609 Multiplexed smFISH data of the mouse primary visual cortex (VISp) was generated as part of

610 the SpaceTx consortium. Tissue processing was carried out as previously described[48], with

611 some modifications.

612

613 Silanization of coverslips (#1.5, Thorlabs CG15KH) was performed by plasma cleaning for 30

614 min in a Plasma-Prep III (SPI 11050-AB), followed by vapor deposition of 3-

615 aminopropyltriethoxysilane (APES, Sigma A3648) in a vacuum for 10 minutes. Coverslips were

616 then washed in 100% methanol for 2 x 5 minutes, allowed to dry, and stored in a dust-free

617 environment until use.

618

619 Fresh-frozen mouse brain tissue was sectioned at 10 μm onto silanized coverslips, let dry for 20

620 min at -20°C, then fixed for 15 min at 4 °C in 4% PFA in PBS. Sections were washed 3 × 10 min

621 in PBS, then permeabilized and dehydrated with chilled 100% methanol at -20°C for 10 min and

622 allowed to dry. Sections were stored at –80 °C until use. Frozen sections were rehydrated in 2X

623 SSC (Sigma 20XSSC, 15557036) for 5 min, then treated 10 min with 8% SDS (Sigma 724255)

624 in PBS at room temperature. Sections were washed 5 times in 2X SSC. Sections were then

625 incubated in hybridization buffer (10% Formamide (v/v, Sigma 4650), 10% dextran sulfate (w/v,

626 Sigma D8906), 200 μg/mL BSA (ThermoFisher AM2616), 2 mM ribonucleoside vanadyl

627 complex (New England Biolabs S1402S), 1 mg/ml tRNA (Sigma 10109541001) in 2X SSC) for 5

628 min at 37°C. Probes were diluted in hybridization buffer at a concentration of 250 nM and

629 hybridized at 37°C for 2 h. Following hybridization, sections were washed 2 × 10 min at 37°C in

630 wash buffer (2X SSC, 20% Formamide), and 1 × 10 min in wash buffer with 5 μg/ml DAPI

631 (Sigma 32670), then washed 3 times with 2X SSC. Sections were then imaged in Imaging buffer

632 (20 mM Tris-HCl pH 8, 50 mM NaCl, 0.8% glucose (Sigma G8270), 30 U/ml pyranose oxidase

633 (Sigma P4234), 50 μg/ml catalase (Abcam ab219092). Following imaging, sections were

26

634     incubated 3 × 10 min in stripping buffer (65% formamide, 2X SSC) at 30°C to remove

635     hybridization probes from the first round. Sections were then washed in 2X SSC for 3 × 5 min at

636     room temperature before repeating the hybridization procedure.

637

638     The multiplexed smFISH image data was collected and processed using methods previously

639     described[48], except that images from different rounds of hybridization were registered in (x, y)

640     based on the DAPI signal. The raw images are available on request.

641

642     **SSAM analysis of VISp multiplexed smFISH data**

643     KDE was performed with bandwidth 2.5 µm. Local maxima were filtered using a gene

644     expression threshold of 0.027, and then filtered with total gene expression threshold of 0.2

645     (**Supplementary Fig. 28A, B**). The selected local maxima vectors were passed to *sctransform*

646     to determine normalization parameters, after which the whole vector field was normalized. To

647     identify rare cell types expected to exist in this tissue, the initial clustering result by Louvain

648     algorithm was sub-clustered by DBSCAN (Method). Initially 49 clusters were obtained with a

649     resolution parameter of 0.15. By manual inspection, several over-clustered cell types, including

650     nine L2/3 IT 1, two L2/3 IT 2, six L4 IT 2, six L6 CT, and two L6 IT 2 clusters were merged, and

651     one spurious cluster was removed, resulting in 28 clusters. The centroids of the clusters are

652     compared with that of scRNA-seq result using Pearson's correlation coefficient

653     (**Supplementary Fig. 28E**). The clusters were named after the highest correlating scRNA-seq

654     cluster, except the newly found 'L4 IT Superficial' (L4 IT 2) cluster. After classification of the

655     local maxima, we quantified the doublet rates (Methods, **Supplementary Table 8**).

656

657     Tissue domains were defined using a sliding circular window with radius 100 µm with step of 10

658     µm over the cell-type map image. Cell type compositions of the windows were clustered using

659     agglomerative clustering, initially with 20 clusters. Clusters with Pearson's correlation higher

660    than 0.7 were merged to result in nine clusters. Further, two clusters were merged since they

661    were different parts of the Pia layer, resulting in a final set of seven clusters representing tissue

662    domains (**Fig. 6**).

663

664    **Plotting**

665    The python packages Matplotlib 3.1.0[49] and Seaborn 0.9.0[50] were used to draw 2D images,

666    plots, and heatmaps. We include helper functions in SSAM to easily generate plots.

667

668    **Movies**

669    Movies were generated by using Virtualdub (1.10.4-AMD64, http://www.virtualdub.org/). The

670    H.264 codec was used to compress videos.

671

672    **Software**

673    Python version 3.7.0 was used throughout. The following python packages were used:

674    *numpy, scipy, pandas, matplotlib, seaborn, scikit-learn, umap-learn, python-louvain, sparse,*

675    *scikit-image*. R package *sctransform* was used for normalization and variance stabilization of

676    the data.

677

678    **Data availability**

679    The source code of SSAM is available online at https://github.com/eilslabs/ssam. A Jupyter

680    notebook (https://github.com/eilslabs/ssam_example) outlines the commands used to download

681    and pre-process the data, and to reproduce the results and figures of this study. The Jupyter

682    notebooks also contain the extensive diagnostic plots used for parameter selection, and choice

683    of removal or merging of clusters. All large files are available online from

684    http://doi.org/10.5281/zenodo.3478502.

685

686    The osmFISH data (Codeluppi et al., 2018) used within the study is available from

687    http://linnarssonlab.org/osmFISH/availability/. The single cell RNA sequencing data of the

688    mouse somatosensory cortex[28,29] are available from http://loom.linnarssonlab.org/. The single

689    cell RNA sequencing data[31] used to compare total mRNA molecules between cell types are

690    available from http://mousebrain.org/. The high resolution poly-A and DAPI images of osmFISH

691    data (Codeluppi et al., 2018) were kindly provided by Sten Linnarsson. The MERFISH data

692    (Moffitt et al., 2018) is available from https://datadryad.org/handle/10255/dryad.192644. Mouse

693    VISp multiplexed smFISH data are available from http://doi.org/10.5281/zenodo.3478502.

694

695    **Acknowledgements**

708

709    **Author contributions**

710    JP, WC designed the concept and idea of SSAM.

711    JP, WC, RE, NI conceived the study.

712    BT, EG, TN.N, BL acquired and interpreted the multiplexed smFISH data.

713    JP, WC, ST, TN.N, NI performed data analysis.

714    LE.B, MS, BL, BT, TG, OS provided critical comments and discussions.

715    RE, NI supervised the study.

716    All authors commented on and critically revised the manuscript.

717

718    **Competing interests**

719    The authors declare no competing interests.

720 **References**

721   1.   Svensson, V., Vento-Tormo, R. & Teichmann, S. A. Exponential scaling of single-cell RNA-

722       seq in the past decade. *Nat. Protoc.* **13**, 599–604 (2018).

723   2.   Regev, A. *et al.* Science Forum: The Human Cell Atlas. *Elife* **6**, (2017).

724   3.   Luecken, M. D. & Theis, F. J. Current best practices in single-cell RNA-seq analysis: a

725       tutorial. *Mol. Syst. Biol.* **15**, e8746 (2019).

726   4.   Salmén, F. *et al.* Multidimensional transcriptomics provides detailed information about

727       immune cell distribution and identity in HER2+ breast tumors. *bioRxiv* 358937 (2018)

728       doi:10.1101/358937.

729   5.   Moffitt, J. R. *et al.* Molecular, spatial, and functional single-cell profiling of the hypothalamic

730       preoptic region. *Science* **362**, (2018).

731   6.   Codeluppi, S., Borm, L. E., Zeisel, A. & La Manno, G. Spatial organization of the

732       somatosensory cortex revealed by osmFISH. *Nature Methods* **15**, 932–935 (2018).

733   7.   Chen, K. H., Boettiger, A. N., Moffitt, J. R., Wang, S. & Zhuang, X. Spatially resolved, highly

734       multiplexed RNA profiling in single cells. *Science* **348**, aaa6090 (2015).

735   8.   Lubeck, E., Coskun, A. F., Zhiyentayev, T., Ahmad, M. & Cai, L. Single-cell in situ RNA

736       profiling by sequential hybridization. *Nature methods* vol. 11 360–361 (2014).

737   9.   Ke, R. *et al.* In situ sequencing for RNA analysis in preserved tissue and cells. *Nat.*

738       *Methods* **10**, 857–860 (2013).

739   10. Lee, J. H. *et al.* Fluorescent in situ sequencing (FISSEQ) of RNA for gene expression

740       profiling in intact cells and tissues. *Nat. Protoc.* **10**, 442–458 (2015).

741   11. Wang, X. *et al.* Three-dimensional intact-tissue sequencing of single-cell transcriptional

742       states. *Science* **361**, (2018).

743   12. Ståhl, P. L. *et al.* Visualization and analysis of gene expression in tissue sections by spatial

744       transcriptomics. *Science* **353**, 78–82 (2016).

745   13. Maniatis, S. *et al.* Spatiotemporal dynamics of molecular pathology in amyotrophic lateral

746        sclerosis. *Science* **364**, 89–93 (2019).

747   14. Vickovic, S. *et al.* High-definition spatial transcriptomics for in situ tissue profiling. *Nat.*

748        *Methods* **16**, 987–990 (2019).

749   15. Hodneland, E., Kögel, T., Frei, D. M., Gerdes, H.-H. & Lundervold, A. CellSegm - a

750        MATLAB toolbox for high-throughput 3D cell segmentation. *Source Code Biol. Med.* **8**, 16

751        (2013).

752   16. Salvi, M. *et al.* Automated Segmentation of Fluorescence Microscopy Images for 3D Cell

753        Detection in human-derived Cardiospheres. *Sci. Rep.* **9**, 6644 (2019).

754   17. Kong, J. *et al.* Automated cell segmentation with 3D fluorescence microscopy images. in

755        *2015 IEEE 12th International Symposium on Biomedical Imaging (ISBI)* 1212–1215 (2015).

756   18. Jiang, J., Kao, P.-Y., Belteton, S. A., Szymanski, D. B. & Manjunath, B. S. Accurate 3D Cell

757        Segmentation using Deep Feature and CRF Refinement. *arXiv [cs.CV]* (2019).

758   19. Shah, S., Lubeck, E., Zhou, W. & Cai, L. In Situ Transcription Profiling of Single Cells

759        Reveals Spatial Organization of Cells in the Mouse Hippocampus. *Neuron* **92**, 342–357

760        (2016).

761   20. Kishi, J. Y. *et al.* SABER amplifies FISH: enhanced multiplexed imaging of RNA and DNA

762        in cells and tissues. *Nat. Methods* **16**, 533–544 (2019).

763   21. Halpern, K. B. *et al.* Single-cell spatial reconstruction reveals global division of labour in the

764        mammalian liver. *Nature* **542**, 352–356 (2017).

765   22. Lignell, A., Kerosuo, L., Streichan, S. J., Cai, L. & Bronner, M. E. Identification of a neural

766        crest stem cell niche by Spatial Genomic Analysis. *Nat. Commun.* **8**, 1830 (2017).

767   23. Thomas, R. M. & John, J. A review on cell detection and segmentation in microscopic

768        images. 1–5 (2017).

769   24. Ester, M., Kriegel, H.-P., Sander, J. & Xu, X. A density-based algorithm for discovering

770        clusters in large spatial databases with noise. in *Proceedings of the Second International*

771           *Conference on Knowledge Discovery and Data Mining* 226–231 (AAAI Press, 1996).

772    25. McInnes, L., Healy, J. & Astels, S. hdbscan: Hierarchical density based clustering. *JOSS* **2**,

773           205 (2017).

774    26. Ankerst, M., Breunig, M. M., Kriegel, H.-P. & Sander, J. OPTICS: ordering points to identify

775           the clustering structure. *SIGMOD Rec.* **28**, 49–60 (1999).

776    27. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell

777           transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.*

778           **36**, 411–420 (2018).

779    28. Zeisel, A. *et al.* Cell types in the mouse cortex and hippocampus revealed by single-cell

780           RNA-seq. *Science* **347**, 1138–1142 (2015).

781    29. Marques, S. *et al.* Oligodendrocyte heterogeneity in the mouse juvenile and adult central

782           nervous system. *Science* **352**, 1326–1329 (2016).

783    30. Tasic, B. *et al.* Shared and distinct transcriptomic cell types across neocortical areas.

784           *Nature* **563**, 72–78 (2018).

785    31. Zeisel, A. *et al.* Molecular Architecture of the Mouse Nervous System. *Cell* **174**, 999–

786           1014.e22 (2018).

787    32. Perkel, J. M. Starfish enterprise: finding RNA patterns in single cells. *Nature* **572**, 549–551

788           (2019).

789    33. Tasic, B. *et al.* Adult mouse cortical cell taxonomy revealed by single cell transcriptomics.

790           *Nat. Neurosci.* **19**, 335–346 (2016).

791    34. Tomioka, R. *et al.* Demonstration of long-range GABAergic connections distributed

792           throughout the mouse neocortex. *Eur. J. Neurosci.* **21**, 1587–1600 (2005).

793    35. Gerashchenko, D. *et al.* Identification of a population of sleep-active cerebral cortex

794           neurons. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 10227–10232 (2008).

795    36. Dries, R. *et al.* Giotto, a pipeline for integrative analysis and visualization of single-cell

796           spatial transcriptomic data. *bioRxiv* 701680 (2019) doi:10.1101/701680.

797   37. Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-seq

798        data using regularized negative binomial regression. *bioRxiv* 576827 (2019)

799        doi:10.1101/576827.

800   38. Hagberg, A., Swart, P. & S Chult, D. Exploring network structure, dynamics, and function

801        using NetworkX. in *Proceedings of the 7th Python in Science Conference (SciPy2008)* 11–

802        15 (Los Alamos National Lab.(LANL), Los Alamos, NM (United States), 2008).

803   39. Fortunato, S. & Barthélemy, M. Resolution limit in community detection. *Proc. Natl. Acad.*

804        *Sci. U. S. A.* **104**, 36–41 (2007).

805   40. Gayoso, A., Shor, J., Carr, A. J., Sharma, R. & Pe'er, D. *JonathanShor/DoubletDetection:*

806        *HOTFIX: Correct setup.py installation.* (2019). doi:10.5281/zenodo.3376859.

807   41. Wolock, S. L., Lopez, R. & Klein, A. M. Scrublet: Computational Identification of Cell

808        Doublets in Single-Cell Transcriptomic Data. *Cell Syst* **8**, 281–291.e9 (2019).

809   42. Kazhdan, M., Bolitho, M. & Hoppe, H. Poisson surface reconstruction. in *Proceedings of the*

810        *fourth Eurographics symposium on Geometry processing* 61–70 (2006).

811   43. Kazhdan, M. & Hoppe, H. Screened poisson surface reconstruction. *ACM Trans. Graph.* **32**,

812        29 (2013).

813   44. Garland, M. & Heckbert, P. S. Surface simplification using quadric error metrics. in

814        *Proceedings of the 24th annual conference on Computer graphics and interactive*

815        *techniques* 209–216 (ACM Press/Addison-Wesley Publishing Co., 1997).

816   45. Hoppe, H. New quadric metric for simplifying meshes with appearance attributes. in

817        *Proceedings Visualization '99 (Cat. No.99CB37067)* 59–510 (1999).

818   46. Schroeder, W., Martin, K. & Lorensen, B. *The Visualization Toolkit: An Object-oriented*

819        *Approach to 3D Graphics.* (Kitware, 2006).

820   47. Cignoni, P. *et al.* Meshlab: an open-source mesh processing tool. in *Eurographics Italian*

821        *chapter conference* vol. 2008 129–136 (2008).

822   48. Hodge, R. D. *et al.* Conserved cell types with divergent features in human versus mouse

823        cortex. *Nature* **573**, 61–68 (2019).

824    49.  Caswell, T. A. *et al. matplotlib/matplotlib v3.1.0.* (Zenodo, 2019).

825        doi:10.5281/zenodo.2893252.

826    50.  Waskom, M. *et al. mwaskom/seaborn: v0.9.0 (July 2018).* (Zenodo, 2018).
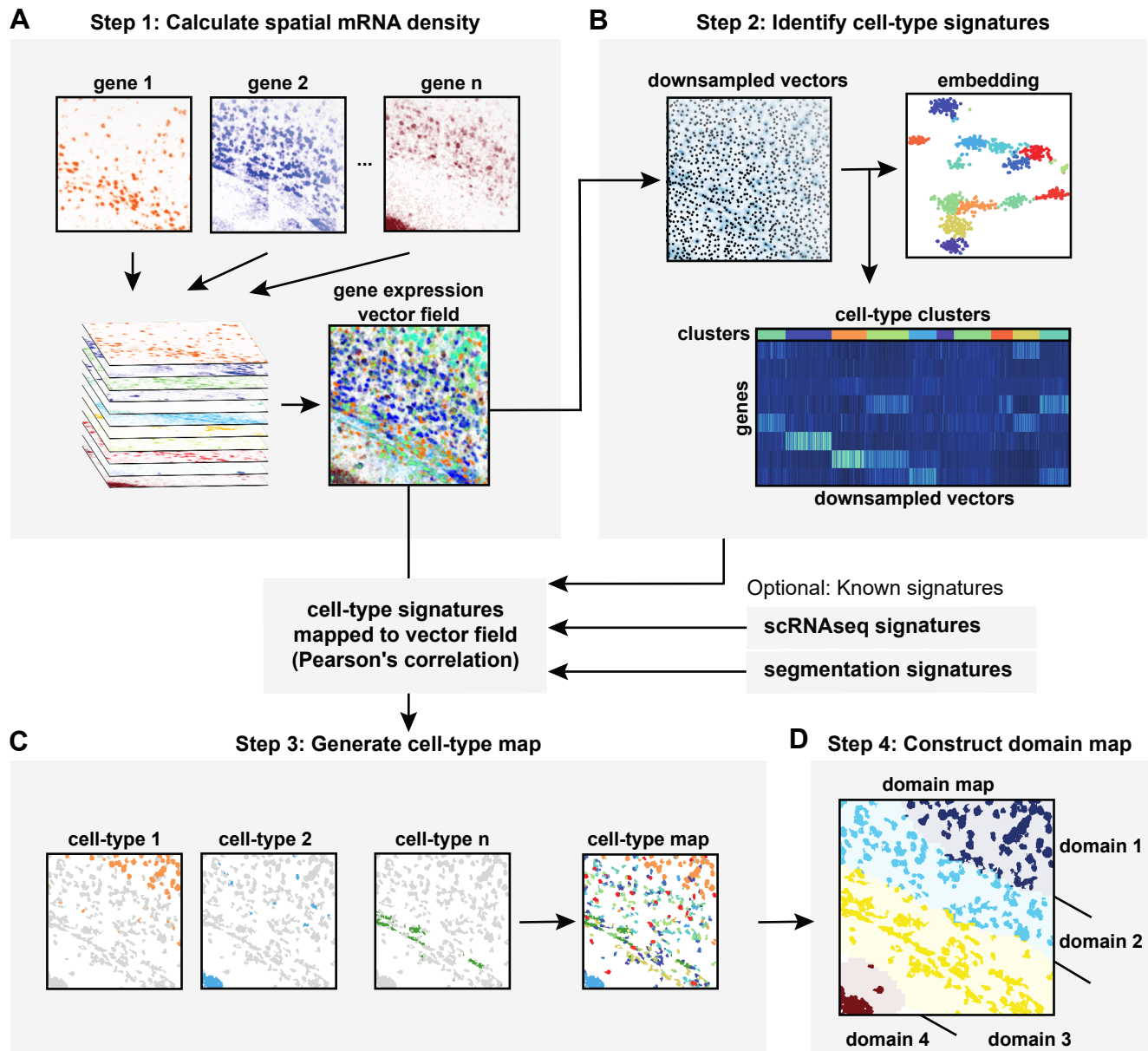
827        doi:10.5281/zenodo.1313201.

828

**Figure 1. Schematic diagram of the SSAM computational workflow for cell type and tissue domain definition based on gene expression data.**

(A) In step 1, SSAM converts mRNA locations into a vector field of gene expression values. For this, SSAM applies a Gaussian KDE to mRNA locations for each gene and projects the resulting mRNA density values to a square lattice which represents coordinates in the tissue. The mRNA density estimated per each gene are stacked to produce a "gene expression vector field" over the lattice. The gene expression vector field is analogous to a 2D/3D image where each pixel/voxel encodes the averaged gene expression of the unit area. Further details of the application of KDE can be found in Supplementary Fig. 1A; (B) In step 2, cell-type signatures are identified *de novo*. First, the gene expression profile at probable cell locations are identified as the local regions in the gene expression vector field where the signal is highest. These downsampled gene expression signals are identified and used for *de novo* cell type identification by cluster analysis. Alternatively, previously defined cell-type signatures can be used. (C) In step 3, a cell-type map is generated. For this, the cell-type signatures are mapped onto the gene expression vector field and cell types are assigned based on Pearson's correlation between each cell-type expression signature to the vector field to define cell-type distribution *in situ*. Further details about creating the cell-type map can be found in Supplementary Fig. 2A; (D) In step 4, the tissue domains are identified. The tissue domain signatures are identified using a sliding window to compute domain signatures based on the count of cell-type labels in the window. The tissue domains are defined by clustering these signatures. Further details on creating the tissue domain map can be found in Supplementary Fig. 2B.
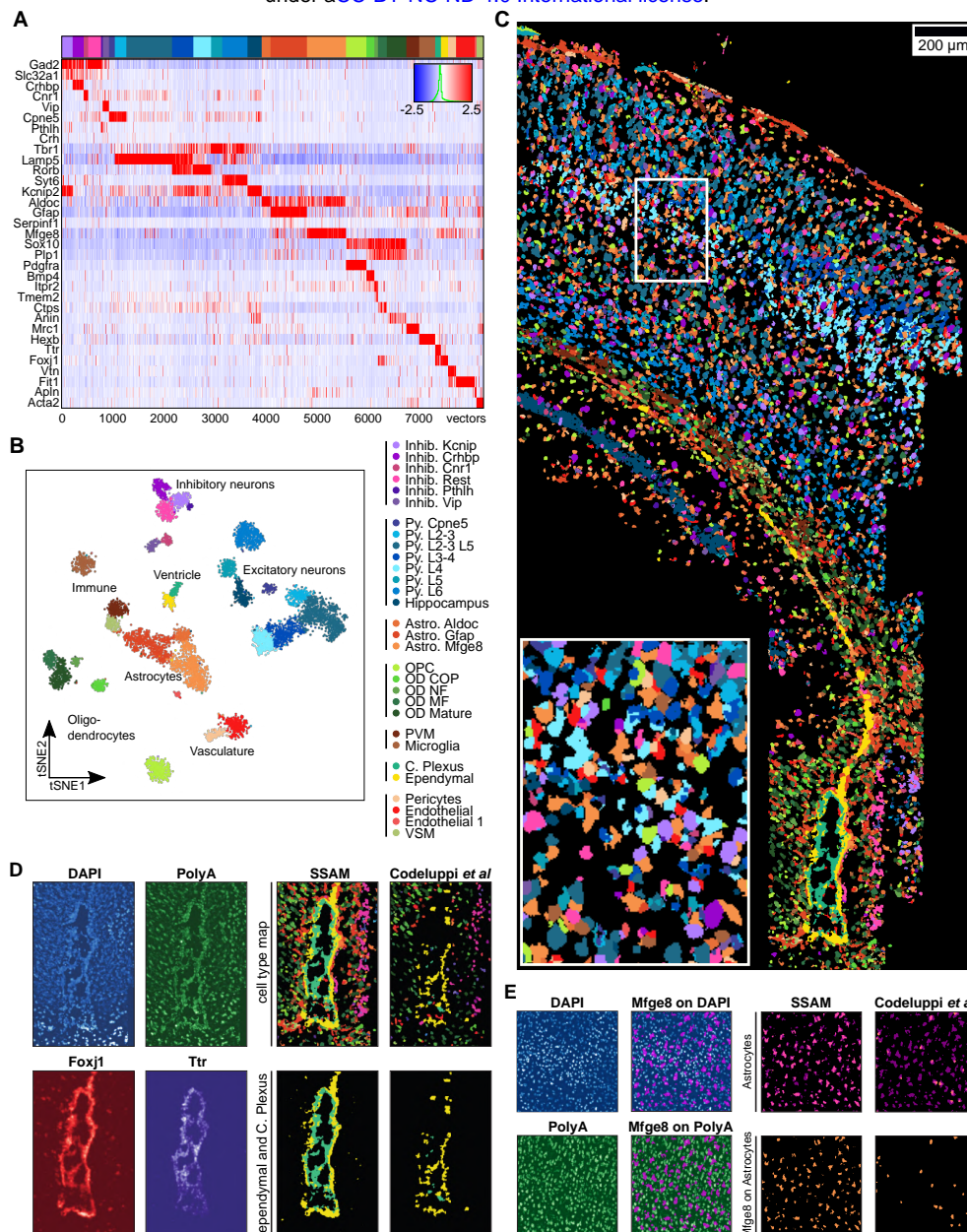
**Figure 2. SSAM improves astrocyte and ventricle detection in the mouse SSp region.**
(A) Gene expression heatmap showing cell-type specific expression of marker genes (8,252 vectors). Rows show z-score normalized gene expression and columns show the gene expression patterns of filtered local maxima vectors. The top annotation shows the cell types and coloring based on the best correlating segmentation-based cell-type signature from Codeluppi *et al*. The colors of the top annotation correspond to the cell type legend in Fig. 2B; (B) A t-SNE map of cell-type signatures with distinct expression. Cell-type clusters are visualized as a 2D t-SNE embedding of filtered local maxima vectors. Cell-type annotation and coloring are based on the best correlating segmentation-based cell-type signature from Codeluppi *et al* (Supplementary Fig. 4C,D). The cell-type legend is grouped by cell-type classes labels shown in the tSNE plot, and are based on groupings by Codeluppi *et al*.; (C) The SSAM *de novo* cell-type map showing spatial organization of the cell types signatures in the gene expression vector field. Inset shows a zoom in of the highlighted tissue region. The colors of the cell types correspond to the cell-type legend in Fig. 2B; (D) SSAM improves the reconstruction of the ventricle. The upper left 2 panels show the DAPI and Poly-A signal around the ventricle area, showing tightly packed cells (occlusion) and lower signal in the ventricle structure compared to surrounding cells. The lower left 2 panels show the KDE gene expression signature for *Foxj1* (the marker for ependymal cells) and *Ttr* (the marker for choroid plexus cells). The upper right 2 panels show the cell-type maps reconstructed by SSAM, showing a more complete reconstruction, and by Codeluppi *et al*., which misses parts of the ventricle structure. The bottom right 2 panels show the reconstructions of only the ependymal (yellow) and choroid plexus (teal) cell types by SSAM and Codeluppi *et al*.; (E) SSAM has increased sensitivity of astrocyte detection. The far left upper and lower panels show DAPI and Poly-A signal for a region in the tissue. The middle left upper and lower panels show the overlap of *Mfge8* signal (a marker for one astrocyte) with DAPI and Poly-A signals, showing that *Mfge8* signal corresponds with low Poly-A signal, but with higher DAPI signal. The top right 2 panels show the cell-type signals for *Mfge8* expressing astrocytes by SSAM and Codeluppi *et al*., showing that SSAM detect much more astrocyte cell types. The bottom right 2 panels shows the overlay of *Mfge8* signal with the cell-type calls by SSAM and Codeluppi *et al*., showing the astrocyte signals detected by SSAM correspond well with *Mfge8* signal.
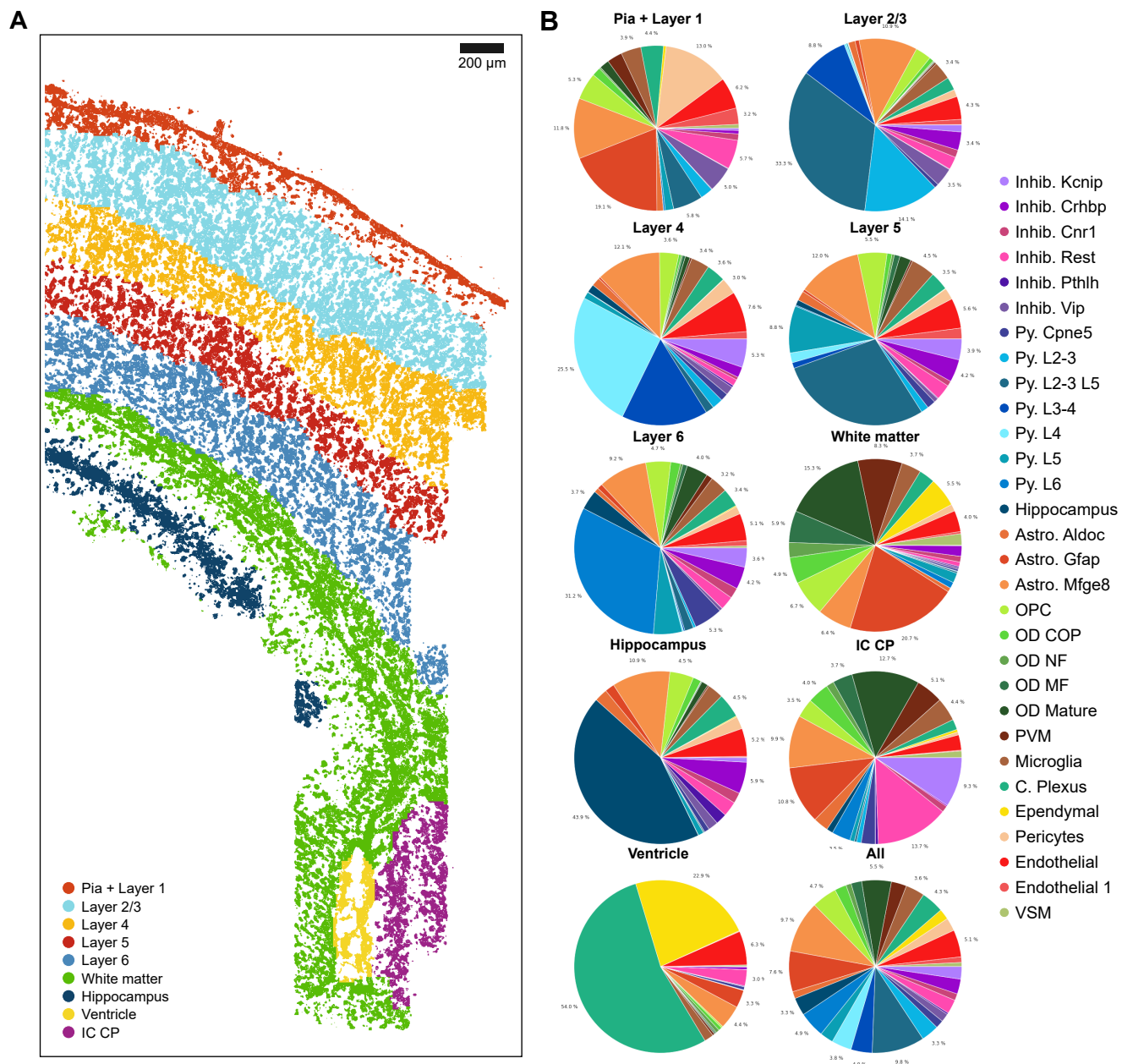
**Figure 3. SSAM identifies cortical layer tissue domains in the mouse SSp cortex.**
(A) Tissue domain map generated by SSAM. Tissue domain signatures were identified from clustering local cell-type composition over sliding 100 μm circular windows, and projected back onto the cell-type map. The reconstruction shows the various cortical layers; (B) Cell-type composition within each tissue domain. The plots show that each domain consists of 7-14% Astrocyte Mfge8 cell types, apart from the ventricle, which instead shows a majority of Choroid plexus and Ependymal cell types.
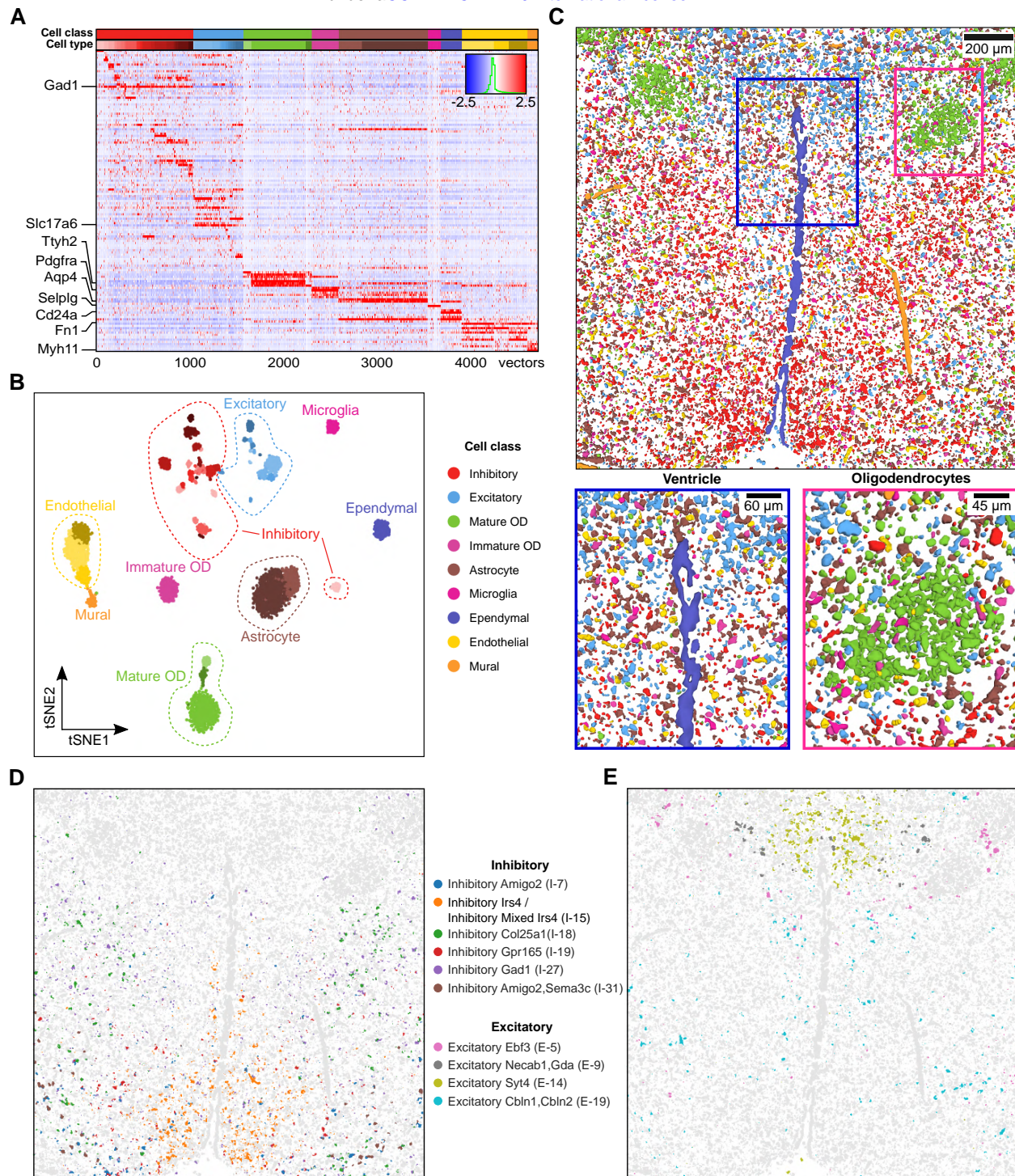
**Figure 4. SSAM 3D cell type map confirms rich diversity of heterogeneous cells in the posterior hypothalamic POA.**
(A) Gene expression heatmap showing cell-type specific expression of marker genes (4,714 vectors). Rows show z-score normalized gene expression and columns show the gene expression patterns of filtered local maxima vectors (representative of gene expression within a cell). The bottom row of the top annotation shows the cell types. Due to a rich diversity of various inhibitory and excitatory neurons captured, the cell types were grouped into classes. The top row of the top annotation shows the cell classes which are named and colored based on the best cell-type signatures and cell classes from Moffitt *et al*. The colors of the cell classes top annotation correspond to the cell-type legend in Fig. 4B. The colors of the cell types are available in Supplementary Fig. 16; (B) A tSNE map of cell-type signatures with distinct expression. Cell-type clusters are visualized as a 2D t-SNE embedding of filtered local maxima vectors. Cell-type annotation and coloring are based on the best correlating segmentation-based cell-type signature from Moffitt *et al*. The tSNE map clearly shows the distinct cluster of different inhibitory and excitatory cell-type signatures. Cell types are grouped into classes based on groupings by Moffitt *et al*.; (C) The SSAM *de novo* 3D cell-type map showing spatial organization of the cell types signatures in the gene expression vector field. Below left and right a zoom in of the highlighted tissue regions of the ventricle structure and clusters of oligodendrocyte cell types. The colors of the cell types correspond to the cell-type legend in Fig. 4B; (D) Spatial localization of various inhibitory cell-type signatures. We found a number of inhibitory cell types which both matched expression signature and tissue localization described by Moffitt *et al.* See also Supplementary Fig. 22; (E) As panel D, but for excitatory cell types.
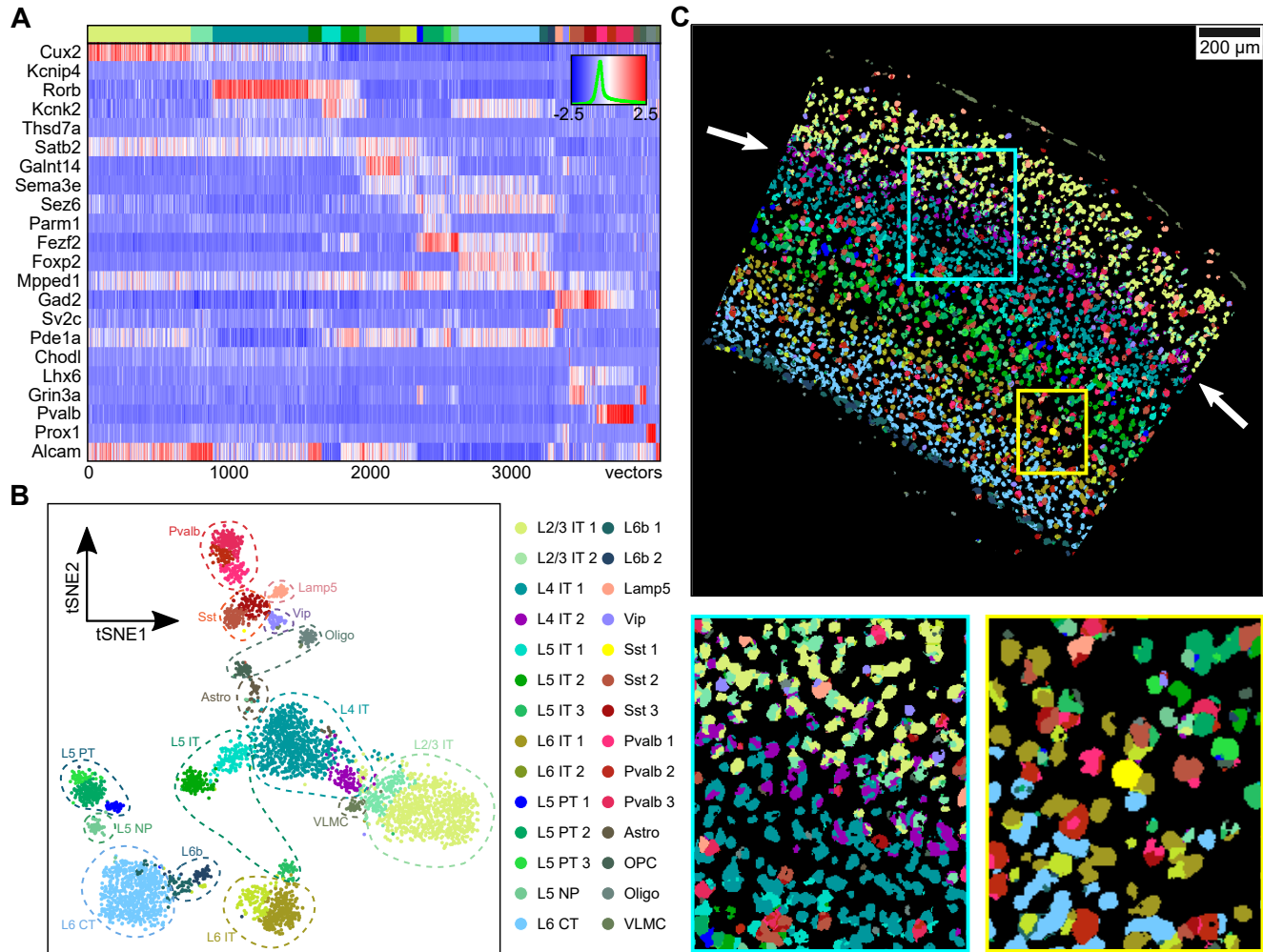
**Figure 5. SSAM identifies layer structure in VISp and confirms rare Sst Chodl cell type in the mouse VISp region.**
(A) Gene expression heatmap showing cell-type specific expression of marker genes (4,113 vectors). Rows show z-score normalized gene expression and columns show the gene expression patterns of filtered local maxima vectors. The top annotation shows the cell types and coloring based on the highest correlating single cell RNA-seq based cell-type signature from previous result (Tasic *et al.*, 2018). The colors of the top annotation correspond to the cell-type legend in Fig. 5B; (B) A tSNE map of cell-type signatures with distinct expression. Cell-type clusters are visualized as a 2D t-SNE embedding of filtered local maxima vectors, with groupings based on the supplementary table 9 of Tasic et al. 2018. Cell-type annotation and coloring are based on the best correlating segmentation-based cell-type signature from previous result (Tasic *et al.*, 2018); (C) The SSAM *de novo* cell-type map showing spatial organization of the cell types. Highlighted are the tissue regions of the cortex including novel L4 IT cell type sub-layering (main panel, purple, white arrows, lower left panel, see also Supplementary Fig. 29B), and rare Sst Chodl cell type (lower right panel, yellow, see also Supplementary Fig. 29C). The colors of the cell types correspond to the cell-type legend in Fig. 5B.
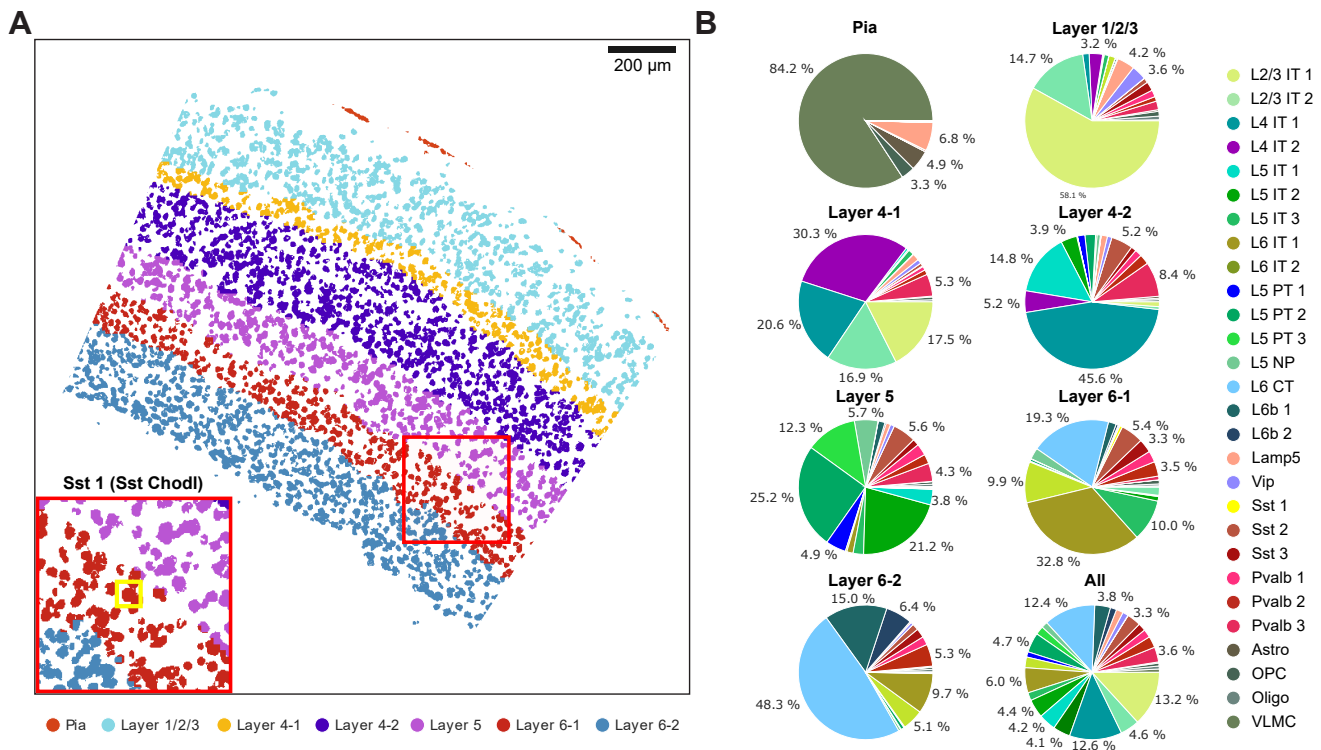
**Figure 6. Rare Sst Chodl cell type localizes to the L6-1 layer of the mouse VISp region.**
(A) Tissue domain map generated by SSAM. Tissue domain signatures were identified from clustering local cell-type composition over sliding 100 μm circular windows, and projected back onto the cell-type map. The reconstruction shows the various cortical layers within the adult mouse VISp, with very clear separation of the Pia layer, and separation of layer 4 and layer 6 into 2 sub-layers. Inset zooms into the location of the rare Sst Chodl cell type found in layer 6-1; (B) Cell-type composition within each tissue domain.