1    **Supplementary Information**

2

3    **ClusterMap for multi-scale clustering analysis of spatial gene expression**

4

5    **Yichun He[1,2,8], Xin Tang[1,2,8], Jiahao Huang[2], Jingyi Ren[2,3], Haowen Zhou[2], Kevin Chen[4],**

6    **Albert Liu[2,3], Hailing Shi[2,3], Zuwan Lin[4,2], Qiang Li[1], Abhishek Aditham[2,5], Johain**

7    **Ounadjela[2,6], Emanuelle I Grody[2,6], Jian Shu[2,6,7], Jia Liu[1*]& Xiao Wang[2,3*]**

8

9    [1]John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge,

10   MA, USA.

11   [2]Broad Institute of MIT and Harvard, Cambridge, MA, USA.

12   [3]Department of Chemistry, Massachusetts Institute of Technology, Cambridge, MA, USA.

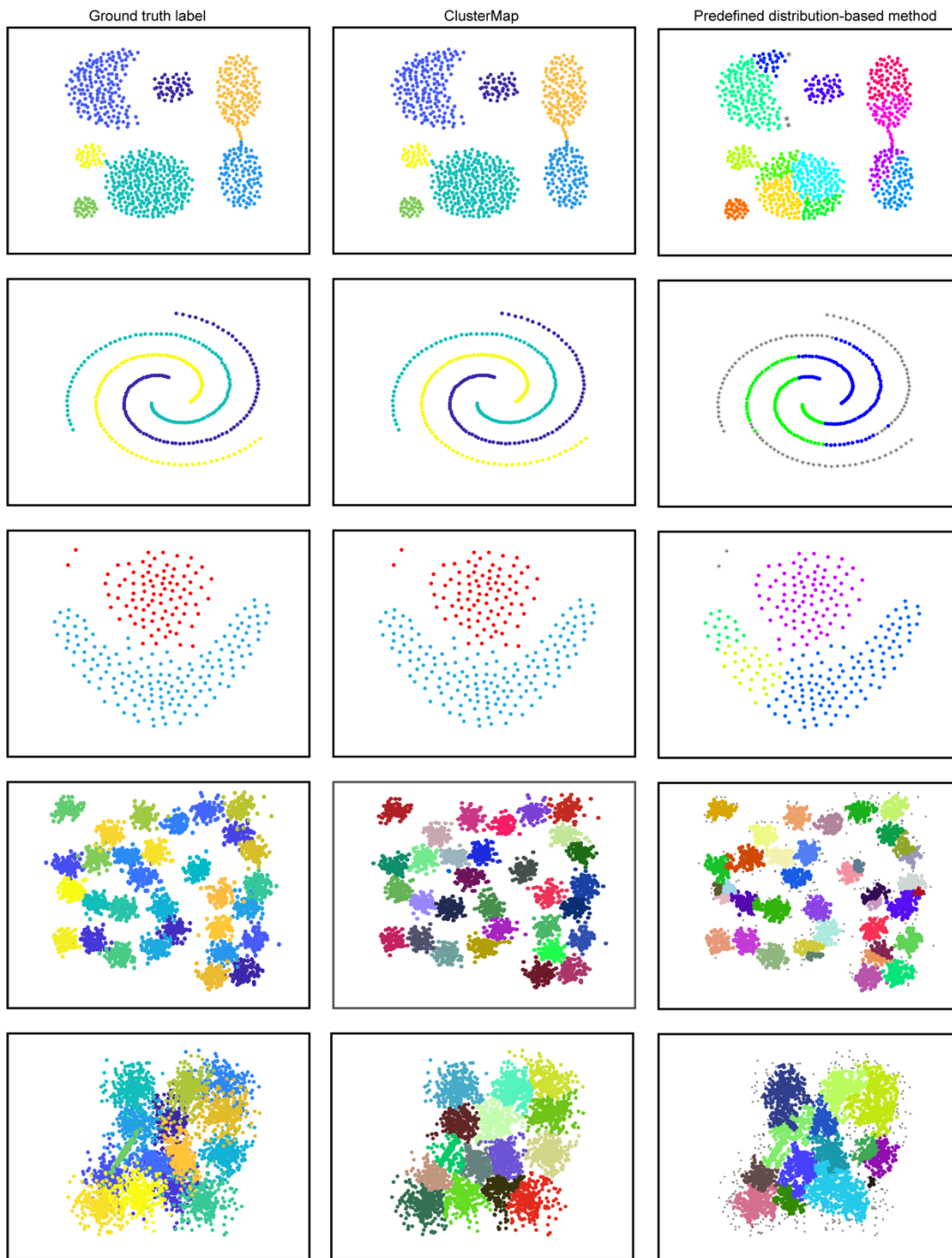13   [4]Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA, USA.

14   [5]Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA,

15   USA.

16   [6]Whitehead Institute for Biomedical Research, Cambridge, MA, USA.

17   [7]Cutaneous Biology Research Center, Massachusetts General Hospital, Harvard Medical School,

18   Boston, MA, USA.

19   [8]These authors contributed equally: Yichun He, Xin Tang.

20   [*]e-mail: wangxiao@broadinstitute.org; jia_liu@seas.harvard.edu.

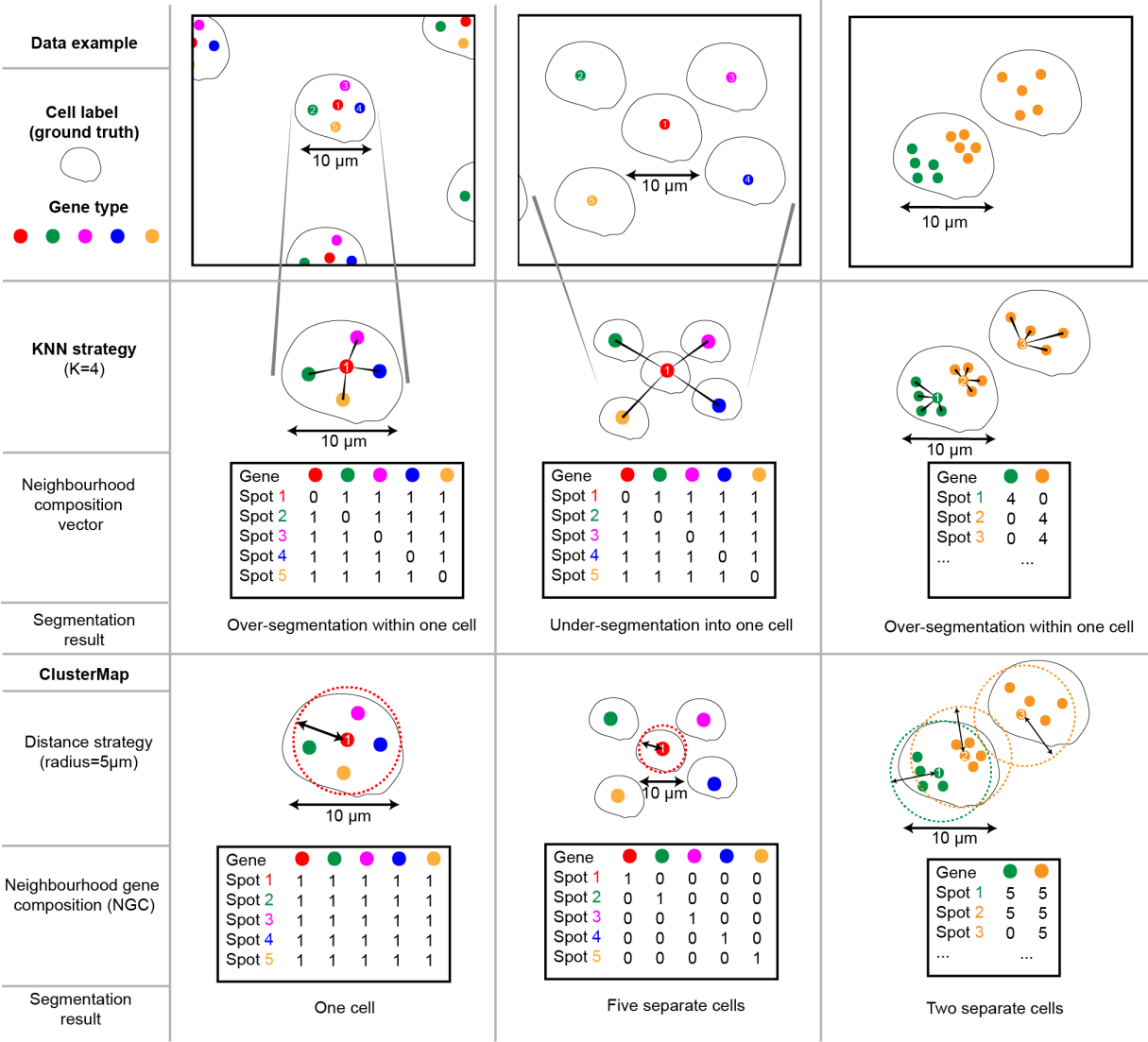| Ground truth label | ClusterMap | Predefined distribution-based method |
|---|---|---|

22    **Supplementary Figure 1**

23    **Performance comparison of ClusterMap and predefined distribution-based method in**

24    **simulated data.** Different colors represent different segmentation results. Note that the gene
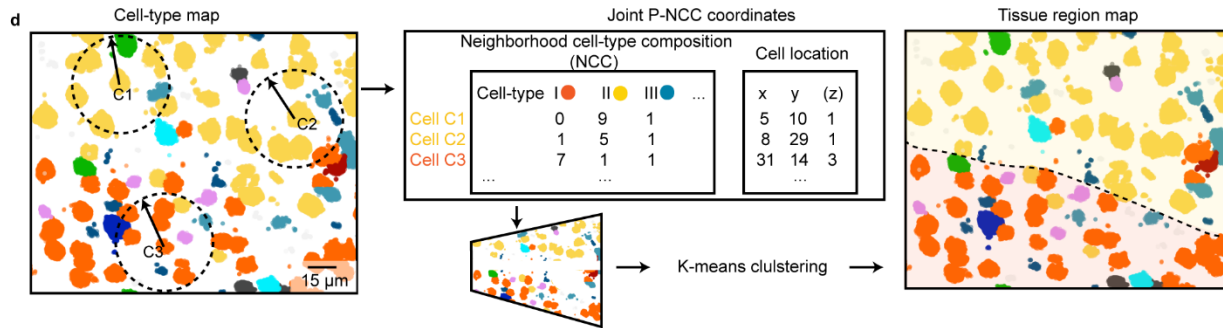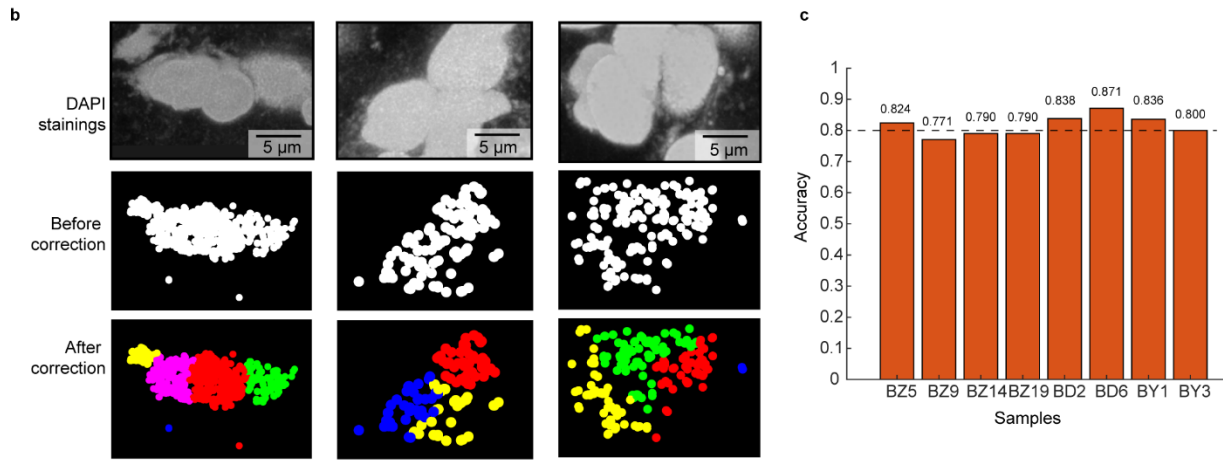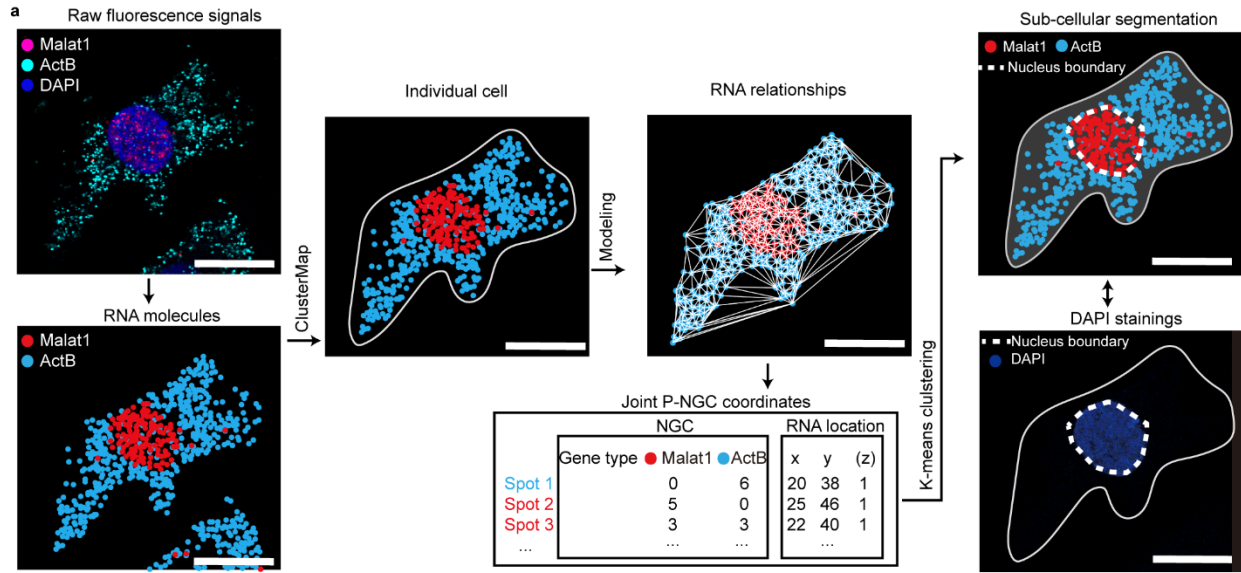
25    identity of each spot is randomly assigned from 1 to 5 as pseudo gene type. Left: ground truth;

26    Middle: ClusterMap results; Right: results using predefined distribution-based method[20].

| | Gene | 🔴 | 🟢 | 🟣 | 🔵 | 🟠 |
|---|---|---|---|---|---|---|
| | Spot 1 | 0 | 1 | 1 | 1 | 1 |
| | Spot 2 | 1 | 0 | 1 | 1 | 1 |
| | Spot 3 | 1 | 1 | 0 | 1 | 1 |
| | Spot 4 | 1 | 1 | 1 | 0 | 1 |
| | Spot 5 | 1 | 1 | 1 | 1 | 0 |

| | Gene | 🔴 | 🟢 | 🟣 | 🔵 | 🟠 |
|---|---|---|---|---|---|---|
| | Spot 1 | 0 | 1 | 1 | 1 | 1 |
| | Spot 2 | 1 | 0 | 1 | 1 | 1 |
| | Spot 3 | 1 | 1 | 0 | 1 | 1 |
| | Spot 4 | 1 | 1 | 1 | 0 | 1 |
| | Spot 5 | 1 | 1 | 1 | 1 | 0 |

| | Gene | 🟢 | 🟠 |
|---|---|---|---|
| | Spot 1 | 4 | 0 |
| | Spot 2 | 0 | 4 |
| | Spot 3 | 0 | 4 |
| | ... | | ... |

| | Gene | 🔴 | 🟢 | 🟣 | 🔵 | 🟠 |
|---|---|---|---|---|---|---|
| | Spot 1 | 1 | 1 | 1 | 1 | 1 |
| | Spot 2 | 1 | 1 | 1 | 1 | 1 |
| | Spot 3 | 1 | 1 | 1 | 1 | 1 |
| | Spot 4 | 1 | 1 | 1 | 1 | 1 |
| | Spot 5 | 1 | 1 | 1 | 1 | 1 |

| | Gene | 🔴 | 🟢 | 🟣 | 🔵 | 🟠 |
|---|---|---|---|---|---|---|
| | Spot 1 | 1 | 0 | 0 | 0 | 0 |
| | Spot 2 | 0 | 1 | 0 | 0 | 0 |
| | Spot 3 | 0 | 0 | 1 | 0 | 0 |
| | Spot 4 | 0 | 0 | 0 | 1 | 0 |
| | Spot 5 | 0 | 0 | 0 | 0 | 1 |

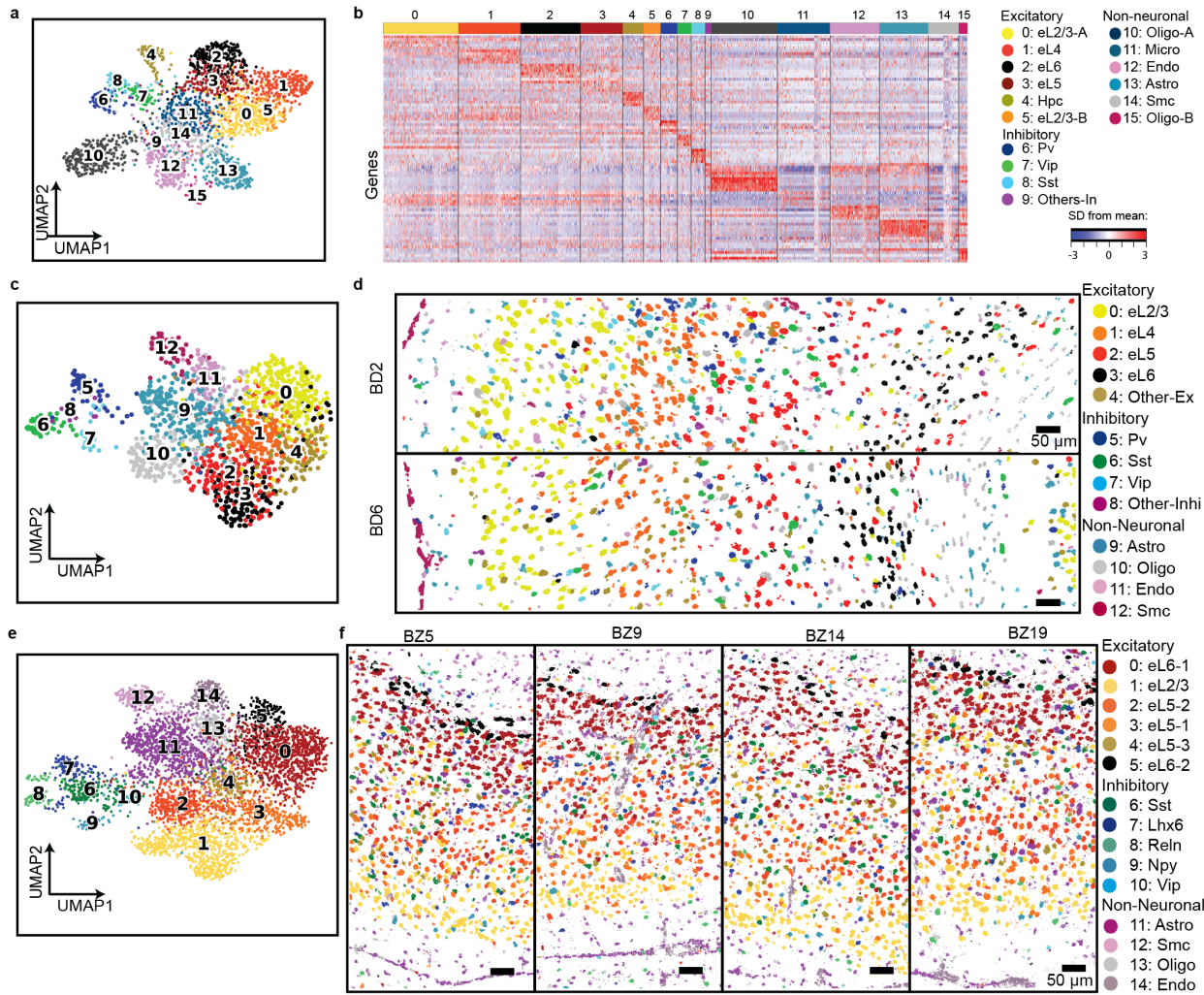| | Gene | 🟢 | 🟠 |
|---|---|---|---|
| | Spot 1 | 5 | 5 |
| | Spot 2 | 5 | 5 |
| | Spot 3 | 0 | 5 |
| | ... | | ... |

27

28

29    **Supplementary Figure 2**

30    **Comparison of RNA sampling approaches between ClusterMap using absolute physical**

31    **distance with other methods[20] using *k*-nearest neighbours (kNN) in simulated data.** Three

32    examples of RNAs with various local density demonstrating that ClusterMap preserves the local

33    physical density information while kNN does not consider the physical density of RNAs.

**a**

Raw fluorescence signals
- Malat1
- ActB
- DAPI

RNA molecules
- Malat1
- ActB

ClusterMap

Individual cell

Modeling

RNA relationships

Joint P-NGC coordinates

| | NGC | | RNA location | | |
|---|---|---|---|---|---|
| Gene type | Malat1 | ActB | x | y | (z) |
| Spot 1 | 0 | 6 | 20 | 38 | 1 |
| Spot 2 | 5 | 0 | 25 | 46 | 1 |
| Spot 3 | 3 | 3 | 22 | 40 | 1 |
| ... | | ... | | ... | |

K-means clulstering

Sub-cellular segmentation
- Malat1
- ActB
- Nucleus boundary

DAPI stainings
- Nucleus boundary
- DAPI

**b**

DAPI stainings — 5 µm

Before correction

After correction

**c**

| Sample | Accuracy |
|---|---|
| BZ5 | 0.824 |
| BZ9 | 0.771 |
| BZ14 | 0.790 |
| BZ19 | 0.790 |
| BD2 | 0.838 |
| BD6 | 0.871 |
| BY1 | 0.836 |
| BY3 | 0.800 |

**d**

Cell-type map

Joint P-NCC coordinates

Neighborhood cell-type composition (NCC)

| Cell-type | I | II | III | ... | Cell location | | |
|---|---|---|---|---|---|---|---|
| | | | | | x | y | (z) |
| Cell C1 | 0 | 9 | 1 | | 5 | 10 | 1 |
| Cell C2 | 1 | 5 | 1 | | 8 | 29 | 1 |
| Cell C3 | 7 | 1 | 1 | | 31 | 14 | 3 |
| ... | | ... | | | | ... | |

K-means clulstering

Tissue region map

36    **Supplementary Figure 3**

37    **Illustration of sub-cellular, cellular, and tissue region analyses. a**, Subcellular analysis process

38    of the panel IV in **Figure 1d** by ClusterMap. A three-channel (magenta: *Malat1*; cyan: *ActB*; blue:

39    DAPI) composite image shows raw fluorescent signals. After preprocessing mRNA molecules

40    with specific genes located, ClusterMap first performs cellular resolution and identifies individual

41    cells. Then a mesh graph that models the relationships among mRNA spots in the cell is generated

42    to compute the NGC coordinates and *K*-means clustering separate spots into two regions using

43    joint physical and NGC coordinates. 100 times K-means clustering was performed with different

44    seeds and showed the consistent same results. Finally, a convex hull is constructed from the

45    nucleus spots, denoting the nucleus boundary. The pattern of ClusterMap-constructed nucleus

46    boundary is compared with the DAPI staining. Scale bar: 20μm. **b**, Examples of the cell

47    identification in ClusterMap procedures in **Fig. 2a**. Upper: DAPI staining showing the cell nuclei.

48    Middle: mRNA spots. Lower: Clustering results. **c**, The accuracy of cell identification results from

49    eight STARmap[6] datasets compared with corresponding expert-annotated labels. BZ5, BZ9, BZ14,

50    BZ19: four STARmap[6] 166-gene sets in mouse medial prefrontal cortex (mPFC); BD2, BD6: two

51    STARmap[6] 160-gene sets in mouse V1. BY1, BY3: two STARmap 1020-gene sets in mouse V1.

52    The horizontal line is at 80% accuracy. **d**, ClusterMap constructs the tissue regions after cell-typing.

53    First, the neighborhood cell-type composition (NCC) of each cell is computed by considering a

54    sliding window over the cell-type map. Then both the NCC and physical locations of cells are

55    combined for *K*-means clustering. Cells with highly correlated neighboring cell-type composition

56    and close spatial distances are merged into a single tissue region signature.

**a**

**b**

Excitatory
0: eL2/3-A
1: eL4
2: eL6
3: eL5
4: Hpc
5: eL2/3-B
Inhibitory
6: Pv
7: Vip
8: Sst
9: Others-In

Non-neuronal
10: Oligo-A
11: Micro
12: Endo
13: Astro
14: Smc
15: Oligo-B

SD from mean:
-3   0   3

**c**

**d**

Excitatory
0: eL2/3
1: eL4
2: eL5
3: eL6
4: Other-Ex
Inhibitory
5: Pv
6: Sst
7: Vip
8: Other-Inhi
Non-Neuronal
9: Astro
10: Oligo
11: Endo
12: Smc

BD2

BD6

50 µm

**e**

**f**

BZ5      BZ9      BZ14      BZ19

Excitatory
0: eL6-1
1: eL2/3
2: eL5-2
3: eL5-1
4: eL5-3
5: eL6-2
Inhibitory
6: Sst
7: Lhx6
8: Reln
9: Npy
10: Vip
Non-Neuronal
11: Astro
12: Smc
13: Oligo
14: Endo

50 µm

59 **Supplementary Figure 4**

60 **Identification of cell types in mouse primary cortex (V1) and medial prefrontal cortex**

61 **(mPFC). a,b**, UMAP and heatmap visualization of all excitatory, inhibitory and non-neuronal cell

62 types in STARmap[6] mouse V1 1020-gene (two replicates: BY1 and BY3). The heatmap represents

63 z-scored expression matrix of all cell types, showing clustering of five most differentially

64 expressed genes per cell type. Genes shown are selected based on a false discovery rate (FDR)-

65 adjusted *p*-value threshold of 0.05 (Benjamini-Hochberg correction) and a minimum $log_{10}$ fold

66 change of 0.4, using a two-sided, unpaired t-test, for genes that are expressed in cells within each

67 cluster versus cells in any other cluster. **c**, **e**, UMAP visualization of all excitatory, inhibitory and

68 non-neuronal cell types in STARmap[6] 160-gene datasets in mouse V1 (two replicates: BD2, BD6,

69 (**c**)), and STARmap[6] 166-gene datasets in mPFC (four replicates, BZ5, BZ9, BZ14, BZ19, (**e**)). **d**,

70 **f**, Spatial organization map of cell types in BD2 and BD6 (**d**), and in BZ5, BZ9, BZ14 and BZ19
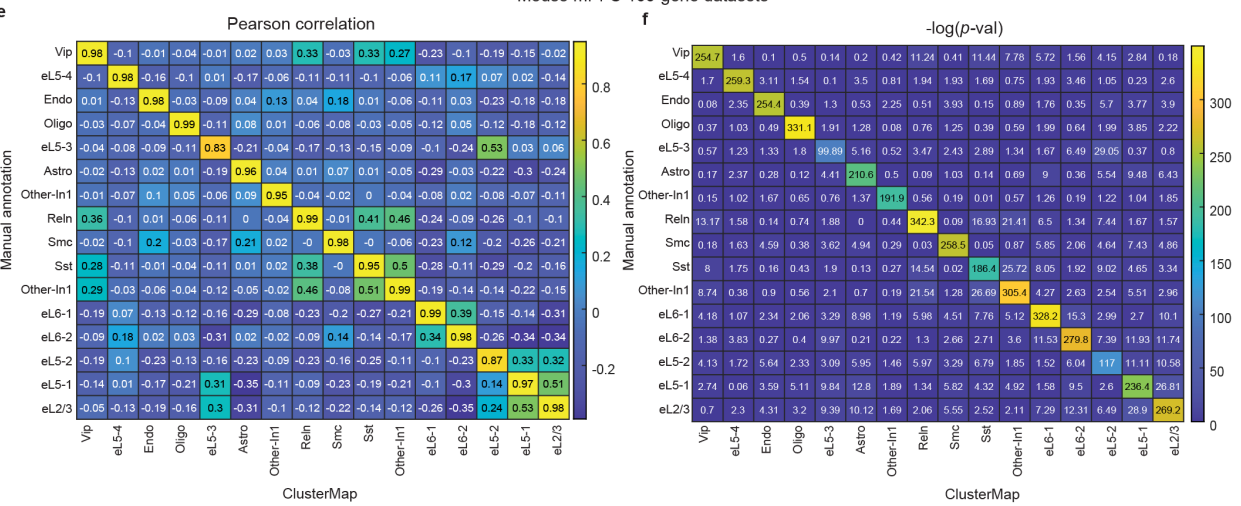
71 (**f**).

STARmap mouse V1 1020-gene datasets

**a** Pearson correlation

**b** -log(*p*-val)

STARmap mouse V1 160-gene datasets

**c** Pearson correlation

**d** -log(*p*-val)

Mouse mPFC 166-gene datasets

**e** Pearson correlation

**f** -log(*p*-val)

73    **Supplementary Figure 5**

74    **Cell-type correlation matrices comparison of ClusterMap-based and manually-segmented**

75    **cell types. a,b**, Comparison on STARmap mouse V1 1020-gene datasets. Heatmaps of Pearson

76    correlation (**a**) and -$log$($p$-value) (**b**) for null hypothesis testing. The $p$ value is based on a t statistic

77    which has $n$-2 degrees of freedom and 95% confidence interval. The single-cell gene expression

78    profiles from ClusterMap with manual annotation are compared. **c,d**, Comparison on STARmap
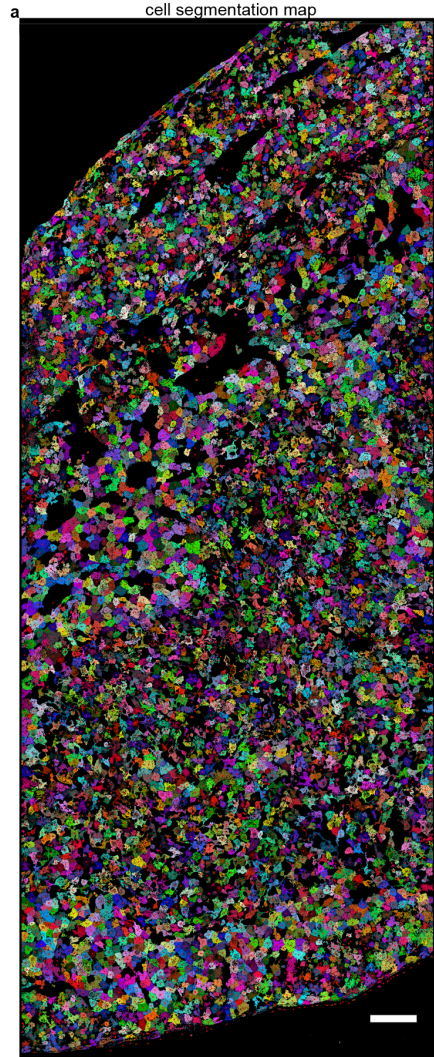
79    mouse V1 160-gene datasets. Heatmaps of correlation (**c**) and -$log$($p$-value) (**d**) comparing the

80    single-cell gene expression profiles from ClusterMap with manual annotation. **e,f**, Comparison on

81    STARmap mouse mPFC 166-gene datasets. Heatmaps of correlation (**e**) and -$log$($p$-value) (**f**)

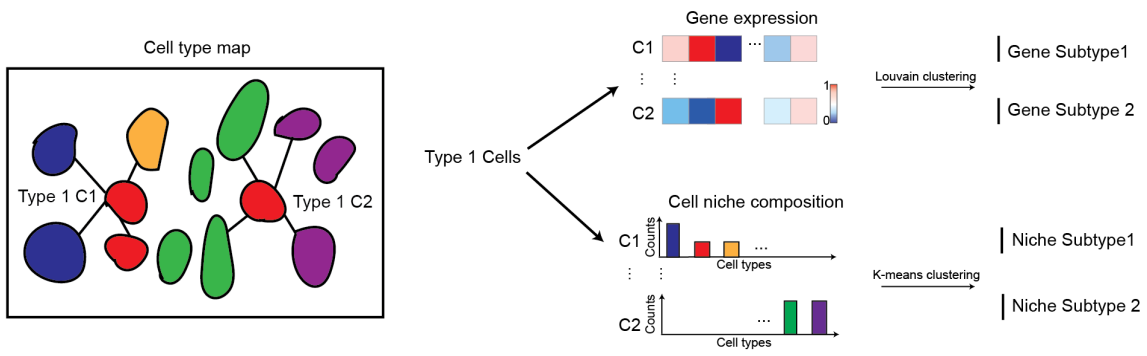82    comparing the single-cell gene expression profiles from ClusterMap with manual annotation.

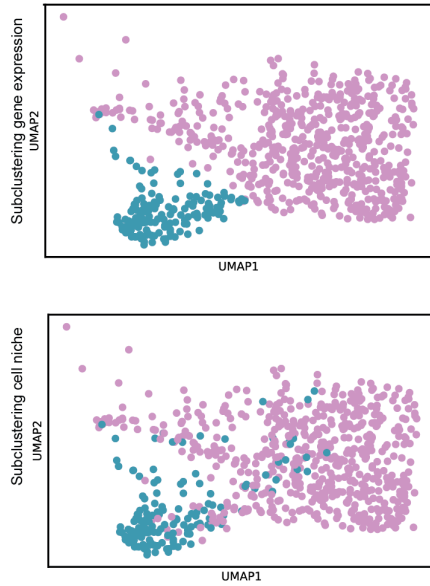83    Horizontal: ClusterMap; vertical: manual annotation.

**a** cell segmentation map

**b**

**c**

**d** UMAP from label transfer result — UMAP from cell-typing result

0: Spongiotrophoblast    5: Endothelial
1: Trophoblast Giant    6: Stromal
2: Glandular Trophoblast    7: Unknown2
3: NK    8: Maternal Decidua
4: Unknown1    9: NA

0: Trophoblast Giant -1    6: Endothelia
1: Trophoblast Giant -2    7: Trophoblast Giant -3
2: Glandular Trophoblast -1    8: Stromal
3: Spongiotrophoblast -1    9: Glandular Trophoblast -2
4: Spongiotrophoblast -2    10: NK
5: Maternal Decidua    11: Trophoblast Giant -4

84

85    **Supplementary Figure 6**

86    **Analyses of the STARmap[6] mouse placental dataset. a**, ClusterMap generates the cell

87    segmentation map of the STARmap[6] mouse placenta 903-gene dataset, including 7,224 cells. Scale

88    bar: 100 μm. **b**, Statistics of ClusterMap identified placental cells as shown in (**a**). Left: Histogram

89    of detected reads (DNA amplicons) per cell. Middle: Histogram of genes per cell. Right:

90    Correlation plot between genes per cell and reads per cell. **c**, Heatmap visualization of 12 cell types.

91    Names are in the right panel of (**d**). **d**, UMAP from label transfer results with scRNA-seq,

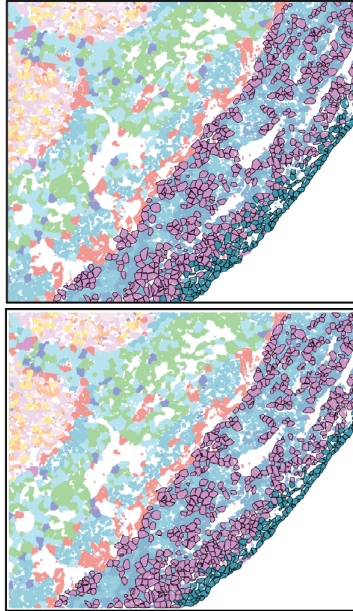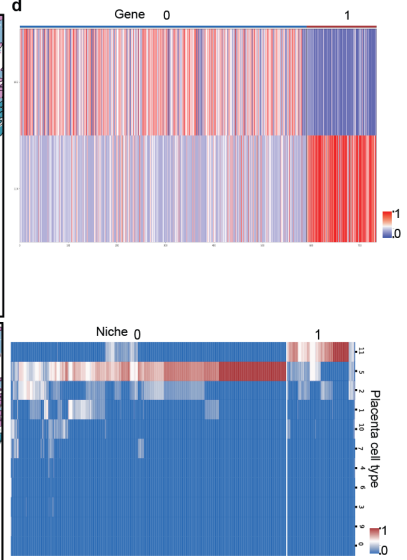92    compared with UMAP of the Louvain clustering[22] in ClusterMap.

**a**

Cell type map

Type 1 Cells

Gene expression

C1

C2

Louvain clustering → Gene Subtype1 / Gene Subtype 2

Cell niche composition

C1
Counts / Cell types

C2
Counts / Cell types

K-means clustering → Niche Subtype1 / Niche Subtype 2

**b**

Subclustering gene expression
UMAP2 / UMAP1
Gene: 0, 1

Subclustering cell niche
UMAP2 / UMAP1
Niche: 0, 1

**c**

**d**

Gene  0  1

Niche  0  1
Placenta cell type

93

94 **Supplementary Figure 7**

95 **Sub-clustering within one cell type using cell niche compositions in STARmap mouse**

96 **placenta 903-gene dataset. a**, Schematic indicating how cells in one cell type are sub-clustered

97 based on either gene expression (Louvain clustering[22]) or the cell niche compositions ($K$-means

98 clustering[19]). **b**, UMAP of gene expression sub-clustering (top) or cell niche composition sub-

99 clustering (bottom) in Maternal Decidua-1 (MD-1). **c**, Spatial subtype maps using gene expression

100 (top) or cell niche composition (bottom) in MD-1. **d**, Heatmap of sub-clustering using gene

101 expression (top) or cell niche composition sub-clustering (bottom) in MD-1. Gene markers in the

102 top heatmaps of gene expression sub-clustering are 0: *GPNMB*, 1: *CXCL14*. Row names in the
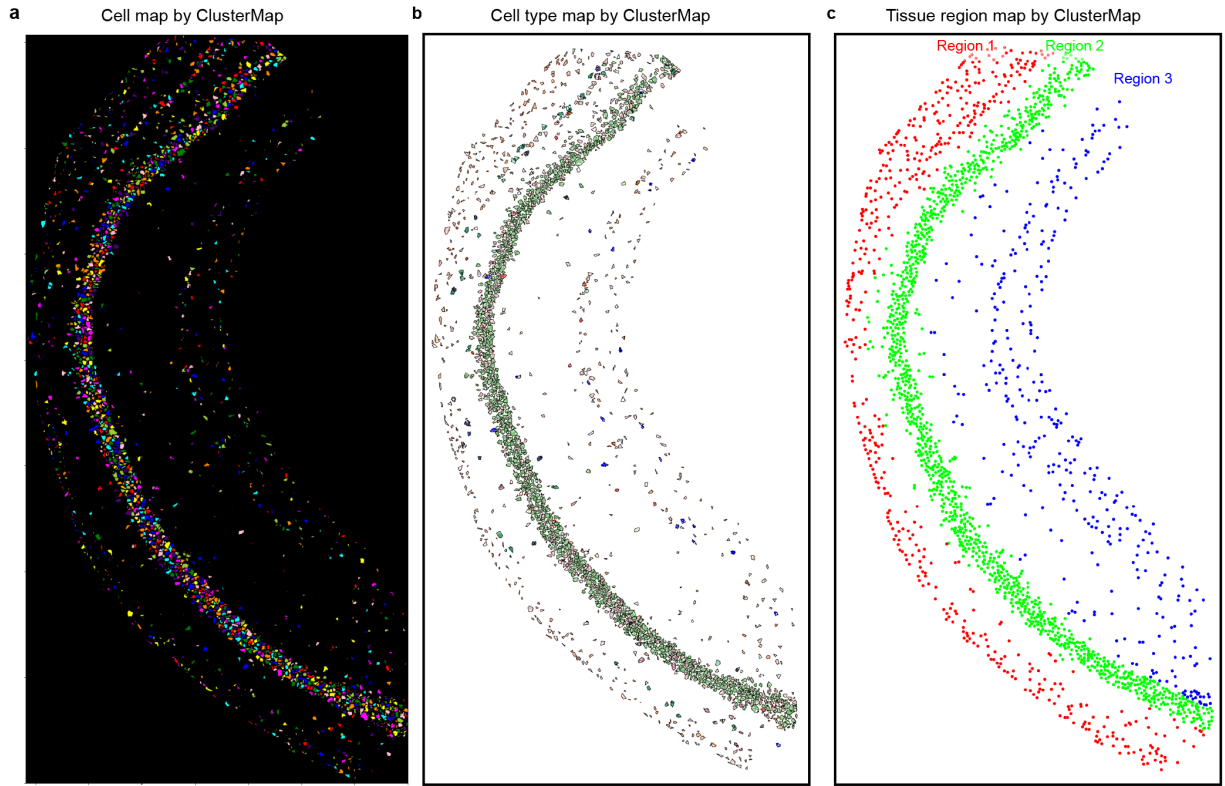
103 bottom heatmaps of cell niche composition sub-clustering are cell types in numbers annotated in
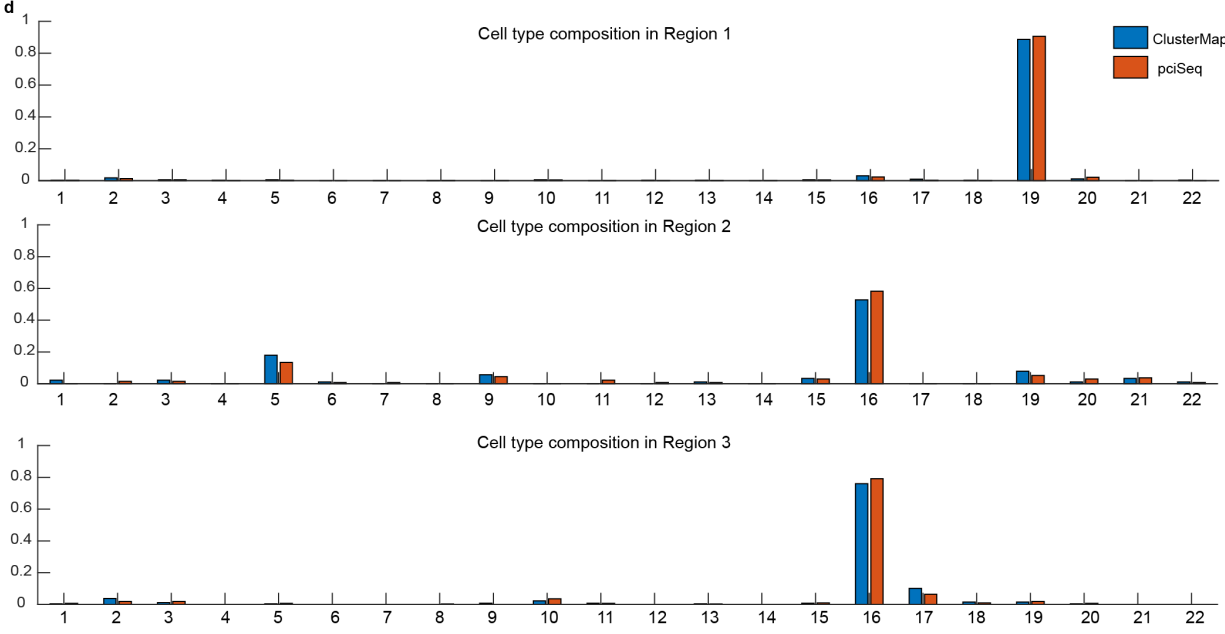
104 Figure 3.

105

**a** Cell segmentation map

**b**

UMAP 2 / UMAP 1

**c**

Gad2_Hybridization2
Slc32a1_Hybridization12
Grlnp_Hybridization10
Cnr1_Hybridization13
Vip_Hybridization6
Pthlh_Hybridization8
Crh_Hybridization10
Cpne5_Hybridization5
Lamp5_Hybridization9
Lamp5i_Hybridization9
Rorb_Hybridization7
Rorb_Hybridization7
Rorb_Hybridization7
Tbr1_Hybridization1
Syt6_Hybridization11
Kcnip_Hybridization12
Gfap_Hybridization2
Aldoc_Hybridization1
Pdgfra_Hybridization8
Bmp4_Hybridization6
Itpr2_Hybridization6
Ctps_Hybridization7
Plp1_Hybridization13
Mec1_Hybridization3
Hexb_Hybridization3
Ttr_Hybridization13
Foxj1_Hybridization1
Vtn_Hybridization12
Flt1_Hybridization2
Apln_Hybridization10
Acta2_Hybridization5

**d** Cell segmentation map

**e** UMAP 2 / UMAP 1

**f**

Gad1
Slc17a6
Sgk1
Pdgfra
Aldh1l1
Sclplg
Cd24a
Fn1
Myh11

106

107    **Supplementary Figure 8**

108    **ClusterMap analyses across different experimental methods. a**, The cell segmentation map of

109    whole osmFISH mouse somatosensory cortex (SSp) datasets. Scale bar: 100 μm. **b,c**, UMAP and

110    heatmap visualization of 31 cell types in osmFISH datasets. **d**, The 2D cell segmentation map of

111    whole MERFISH mouse preoptic area (POA) datasets. Scale bar: 200 μm. **e,f**, UMAP and heatmap

112    visualization of 9 cell types in MERFISH datasets. The number of cells increased from 6,471 to

113    8,538 for osmFISH, from 2,620 to 2,924 for pciSeq, from 6,977 to 10,320 for MERFISH. The

114    number of reads increased from 1,248,106 to 1,690,328 for osmFISH, from 31,246 to 31,750 for

115    pciSeq, from 1,927,913 to 3,065,171 for MERFISH.

**a** Cell map by ClusterMap  **b** Cell type map by ClusterMap  **c** Tissue region map by ClusterMap

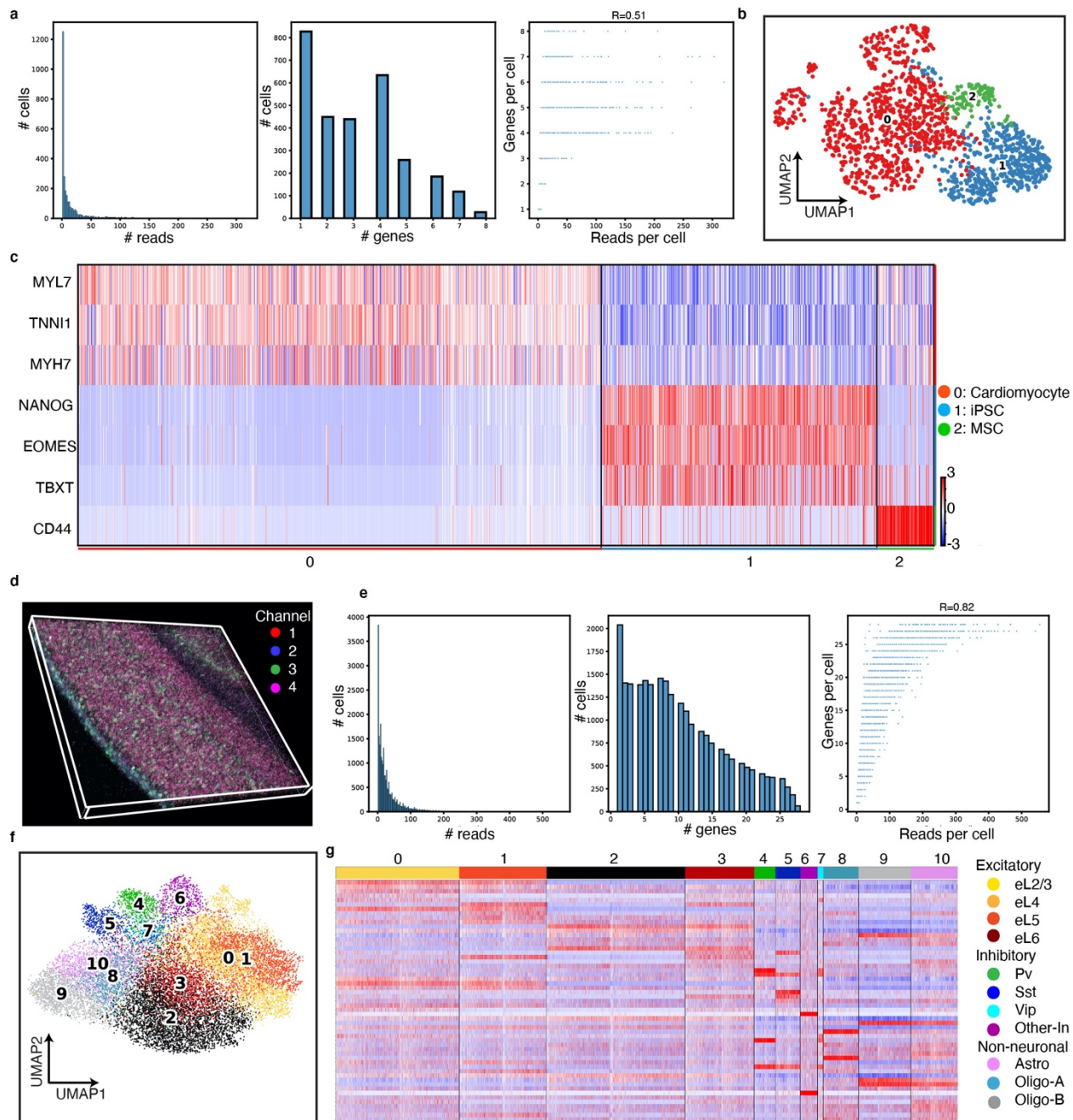Region 1   Region 2   Region 3

Legend:
- 1: Axo-axonic
- 2: Basket
- 3: Bistratified
- 4: CGE IVY
- 5: CGE NGF
- 6: Cck Calb1/Slc17a8
- 7: Cck Cxcl 14+
- 8: Cck Vip Cxcl 14-
- 9: Cxcl14 NGC
- 10: Hippo
- 11: IS1
- 12: IS2
- 13: IS3
- 14: Ivy
- 15: MGE
- 16: Nonneuron
- 17: O-Bi
- 18: O/LM
- 19: PC CA1
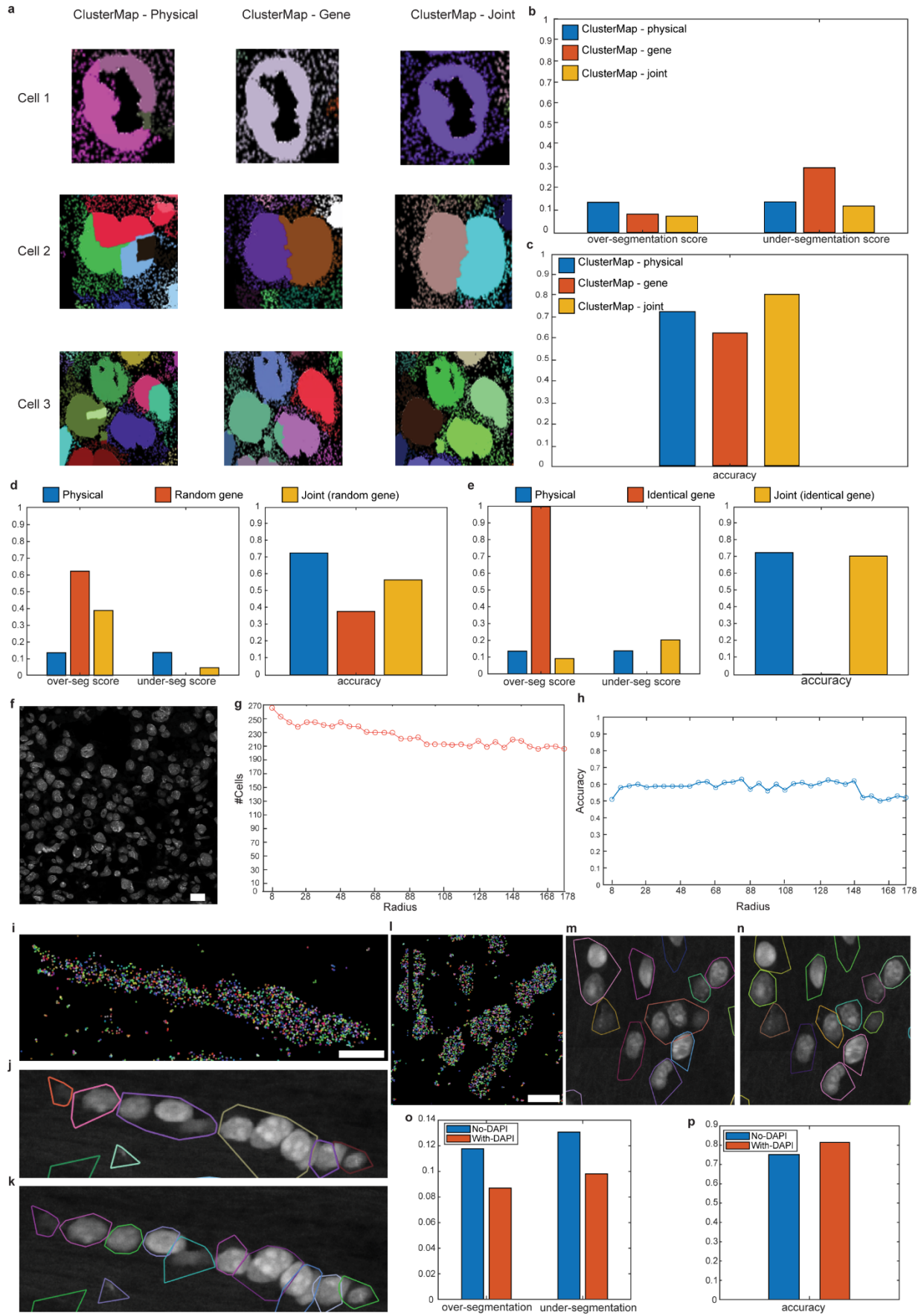- 20: PC other
- 21: Radiatum retrohip
- 22: Trilaminar
- Uncalled

**d**

Cell type composition in Region 1

Cell type composition in Region 2

Cell type composition in Region 3

ClusterMap
pciSeq

116

117     **Supplementary Figure 9**

118     **ClusterMap analyses of ISS data. a**, Cell segmentation map shows the cell segmentation results

119     by ClusterMap. Colors are randomly assigned to each cell mask. **b**, Cell type map shows the cell

120     type calling results. Colors are assigned according to their corresponding cell type categories. **c**,

121     Tissue region map shows laminar structure of hippocampus. Scale bar: 200 μm. **d**, Side-by-side

122     comparison of cell type compositions in each tissue region from ClusterMap and pciSeq of the ISS
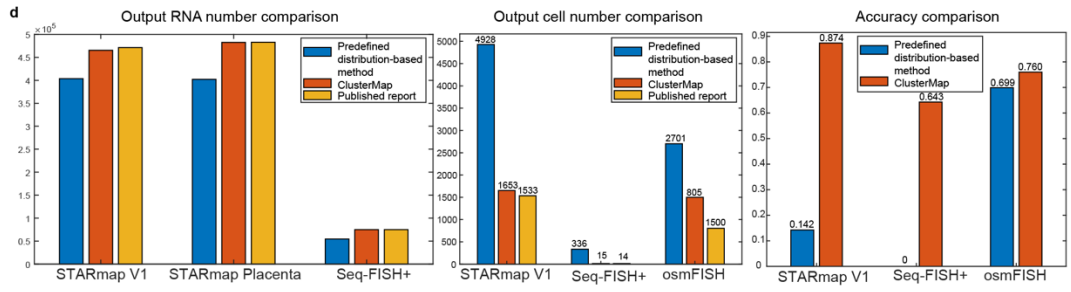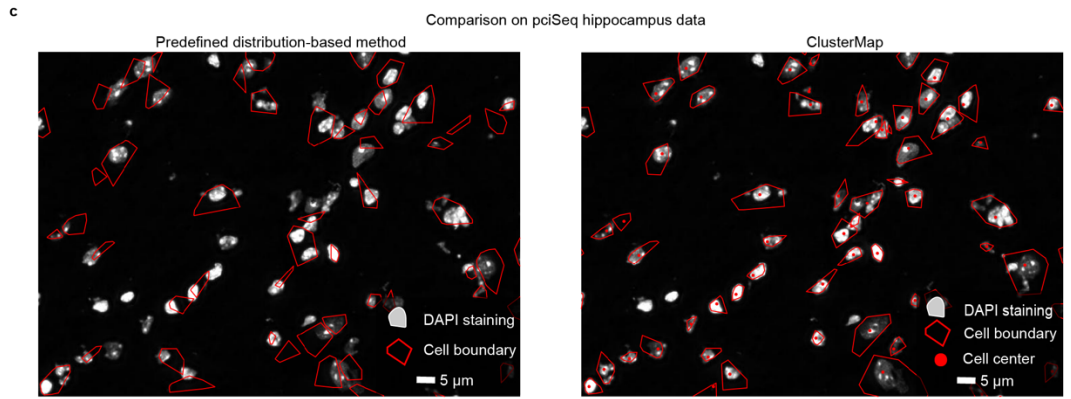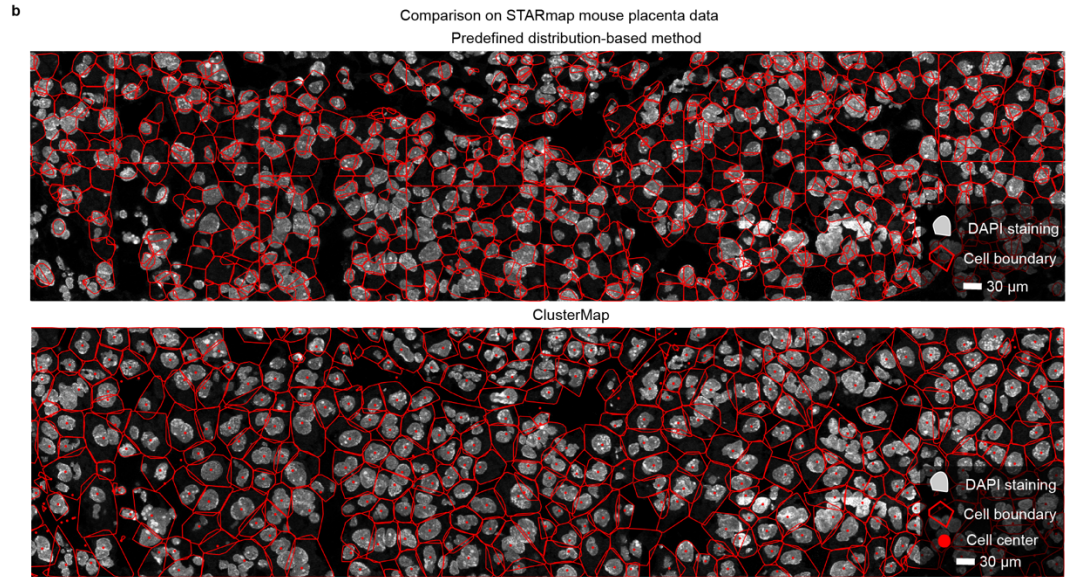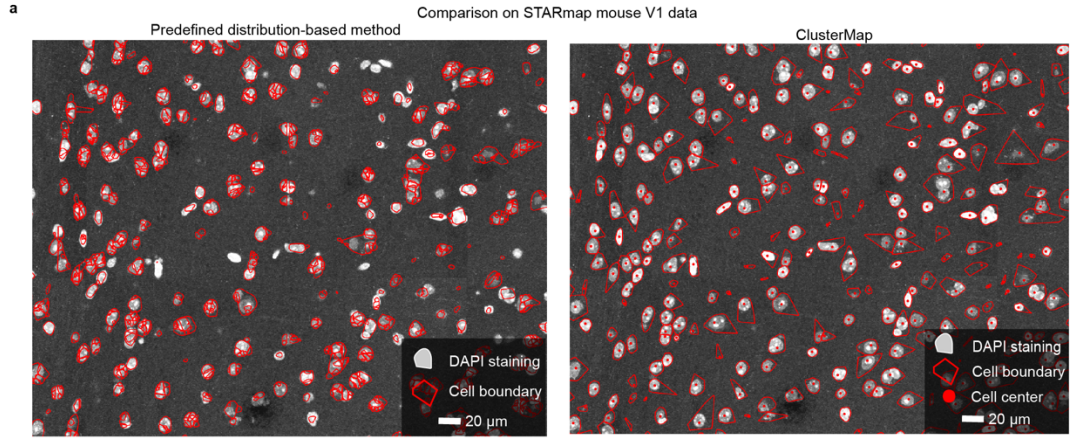
123     data.

125    **Supplementary Figure 10**

126    **ClusterMap analyses in the 3D datasets. a**, Statistics of ClusterMap identified cells in the 3D

127    STARmap[6] cardiac organoid[27] 8-gene dataset. Left: Histogram of detected reads (DNA amplicons)

128    per cell. Middle: Histogram of genes per cell. Right: Correlation plot between genes per cell and

129    reads per cell. **b**, **c**, UMAP and heatmap visualization of three cell types in the STARmap[6] cardiac

130    organoid 8-gene dataset.  The number of cells in each cell type is as follows: cardiomyocytes, 929;

131    induced pluripotent stem cells (iPSCs), 489; mesenchymal stem cells (MSCs), 101. **d**, 3D four-

132    channel composite raw fluorescent image of the first sequencing round shows spatial arrangement

133    of mRNA molecules in the STARmap[6] mouse V1 28-gene dataset. Width 184 µm, height 194 µm,

134    depth 100 µm. **e**, Statistics of ClusterMap identified cells in (**d**). Left: Histogram of detected reads

135    (DNA amplicons) per cell. Middle: Histogram of genes per cell. Right: Correlation plot between

136    genes per cell and reads per cell. **f**, **g**, UMAP and heatmap visualization of three cell types of (**d**).

138  **Supplementary Figure 11**

139  **Performance comparison of ClusterMap using physical density, gene distance, and joint**

140  **information. a**, Examples of cell segmentation using only physical density information (left), gene

141  (NGC) distance information (middle), and joint information (right). **b,c**, Bar plots demonstrating

142  the percentage of over-/under- segmented cells in ground truth cells (**b**) and overall accuracy (**c**)

143  in using physical distances information, gene (NGC) distance information, and joint information.

144  **d,e**, Bar plots demonstrating the percentage of over-/under- segmented cells in ground truth cells

145  (left) and overall accuracy (right) in using physical distances information, random (**d**) or identical

146  (**e**) gene (NGC) distance information, and joint information. **f**, Raw DAPI image of the targeted

147  mouse placenta tissue. Scale bar: 20μm. **g,h**, Line plots showing the number of cells and overall

148  accuracy to the radius. **i-m**, Two examples of the hippocampus regions in STARmap mouse V1

149  1020-gene datasets showing raw spatial transcriptomics data (**i,l**), ClusterMap results without

150  DAPI (**j,m**), and ClusterMap results with DAPI (**k,n**). Scale bar: 20 μm. **o,p**, Bar plots showing

151  the percentage of over-/under- segmented cells (**o**) and overall accuracy (**p**) from ClusterMap

152  without and with DAPI.

**a**

Comparison on STARmap mouse V1 data

Predefined distribution-based method

ClusterMap

**b**

Comparison on STARmap mouse placenta data

Predefined distribution-based method

ClusterMap

**c**

Comparison on pciSeq hippocampus data

Predefined distribution-based method

ClusterMap

**d**

Output RNA number comparison

Output cell number comparison

Accuracy comparison

153

154 **Supplementary Figure 12**

155 **Performance comparison of ClusterMap and other method across different types of *in situ***

156 **transcriptomic data. a,** Example of a region in the STARmap[6] mouse V1 1020-gene dataset with

157 DAPI signals (gray) showing ground truth cell nuclei. Red contours show cell boundaries

158 identified by predefined distribution-based method[20] (left) and ClusterMap (right), respectively.

159 **b,c**, As in (**a**) but using the STARmap[6] mouse placenta 903-gene dataset and published pciSeq[4]

160 dataset. **d**, Bar plots demonstrating the remaining RNA numbers, cell numbers and segmentation

161 accuracy for each dataset. In each bar plot, results from predefined distribution-based method[20],

162 ClusterMap, and published reports were shown in blue, red, and yellow, respectively.

163

164  **Supplementary Table**

165  **Supplementary Table 1.** Summary of the name, *in situ* sequencing protocol, number of genes,

166  number of cells, number of reads, number of cell types, corresponding figures and note of 7

167  datasets.

| Dataset | Experimental Method | # Genes | # Cells | # Reads | # Cell types | Figures | Note |
|---|---|---|---|---|---|---|---|
| STARmap mouse V1 1020-gene | STARmap | 1,020 | 1,599 | 863,426 | 16 | Fig. 1c, Fig. 2, Supplementary Figs. 3,4,5 | Source: Ref. 6. 2D analysis |
| STARmap mouse placenta 903-gene | STARmap | 903 | 7,224 | 5,090,930 | 12 | Fig. 3, Fig. 4, Supplementary Figs. 6,7 | New data. 2D analysis |
| MERFISH mouse POA | MERFISH | 140 | 10,320 | 3,065,171 | 9 | Fig. 5, Supplementary Fig. 8 | Source: Ref. 3. 3D analysis |
| pciSeq mouse hippocampus | ISS | 98 | 2,924 | 31,750 | 23 | Fig. 5, Supplementary Fig. 9 | Source: Ref. 4. 2D analysis |
| osmFISH mouse SSp | osmFISH | 33 | 8,538 | 1,690,328 | 31 | Fig. 5, Supplementary Fig. 8 | Source: Ref. 5. 2D analysis |
| STARmap cardiac organoid 8-gene | STARmap | 8 | 1,519 | 47,594 | 3 | Fig. 6, Supplementary Fig. 10 | New data, 3D analysis |
| STARmap mouse V1 28-gene | STARmap | 28 | 24,590 | 753,396 | 11 | Fig. 6, Supplementary Fig. 10 | Source: Ref. 6. 3D analysis |

168