# SOAPy: a Python package to dissect spatial architecture, dynamics and communication

Heqi Wang[1#], Jiarong Li[1#], Siyu Jing[1], Ping Lin[1], Yu Li[2], Haibing Zhang[2], Yujie Chen[1], Zhen Wang[1] & Hong Li[1*]

**Affiliations**
[1]CAS Key Laboratory of Computational Biology, Shanghai Institute of Nutrition and Health, University of Chinese Academy of Sciences, Chinese Academy of Sciences, Shanghai 200031, China
[2]CAS Key Laboratory of Nutrition, Metabolism and Food Safety, Shanghai Institute of Nutrition and Health, University of Chinese Academy of Sciences, Chinese Academy of Sciences, Shanghai, China

[#]These authors contributed equally
[*]Correspondence to: lihong01@sinh.ac.cn (H. L.)

## Abstract

Advances in spatial omics technologies have brought opportunities to dissect tissue microenvironment, while also posing more requirements and challenges for computational methods. Here we developed a package SOAPy to systematically dissect spatial architecture, dynamics and communication from spatial omics data. Specifically, it provides analysis methods for multiple spatial-related tasks, including spatial domain, spatial expression tendency, spatiotemporal expression pattern, cellular co-localization, multi-cellular niches, and ligand-receptor-mediated and spatial-constrained cell communication. Applying SOAPy on different spatial omics technologies and diverse biological fields has demonstrated its power on elucidation of biological questions about tumors, embryonic development, and normal physiological structures. Overall, SOAPy is a universal tool for spatial omics analysis, providing a foundation for continued investigation of the microenvironment.

## Keywords

33  spatial omics, Python package, microenvironment, expression pattern, multi-cellular

34  niche, cell communication

35

## Introduction

37  Spatially resolved transcriptomics has been crowned Method of the Year 2020 by

38  Nature Methods[1]. Since then, more and more experimental methods for measuring

39  expression levels of genes, proteins or metabolites in a spatial context have been

40  developed. These technologies include barcode-based and imaging-based ones, which

41  differ in resolution, accuracy and throughout[2,3]. The most widely used 10X Visium

42  spatial transcriptomics measures thousands of genes in each 55μm spot that typically

43  contains 1-10 cells[4]. And imaging-based methods reach more microscopic resolution,

44  such as MIBI-TOF[5] and PhenoCycler-Fusion[6], both detecting dozens of proteins at

45  subcellular resolution. Additionally, spatial multi-omics technologies that

46  simultaneously measure multiple molecular types are emerging, e.g NanoString

47  GeoMx DSP for 18000 RNAs and 140 proteins in the region of interest (usually >100

48  cells)[7].

49  With the development of experimental methods, corresponding analysis pipelines

50  have been designed for pre-processing raw data from specific experimental platforms,

51  such as Space Ranger for 10X Visium and MCMICRO for multiplexed tissue

52  imaging[8]. Methods adapted from single-cell RNA sequencing (scRNA-seq) data

53  analysis could be used to perform standard dimensional reduction, clustering, cell

54  type annotation and marker selection for spatial-omics data[9] that do not require spatial

55  information. And for low resolution spatial technologies, various deconvolution

56  methods have been developed to impute cell-type composition from the mixture of

57  cells.

58  After these pre-processing, downstream analyses are largely independent of

59  experimental technologies, focus on the key feature of spatial omics: space. For

60  example, identifying spatial variable genes[10–12], detecting spatial domains[13], inferring

61  genes or cell-subtypes associated with spatial localization, and so on[3]. Earlier

62 algorithms were often designed for one specific task, tools that fit in with various

63 analysis tasks are becoming popular. A pioneer work Giotto not only builds a data

64 pre-processing pipeline similar to scRNA-seq data analysis[14], but also provides

65 modules for spatial pattern detection, cell neighborhood analysis, and interactive

66 visualization. Squidpy provides scalable analysis framework for both spatial

67 neighborhood graph and image, along with an interactive visualization tool[15]. stlearn

68 is another integrated package for spatial transcriptomic analysis, which adds the

69 functions of spatial trajectories and pseudotime analysis[16]. Investigating the spatial

70 organization of tissue microenvironment are important applications of spatial omics,

71 which may gain new insights in various biological fields. However, the related

72 analysis methods are scattered or lacking, a package for integrative analysis of

73 microenvironmental spatial organization is in an urgent need.

74 To address this problem, we present a package SOAPy (Spatial Omics Analysis in

75 Python) to jointly perform multiple tasks for dissecting spatial organization, including

76 spatial domain, spatial expression tendency, spatiotemporal expression pattern,

77 co-localization of paired cell types, multi-cellular niches, and cell-cell communication.

78 SOAPy improves on previous tools in three main areas (**Table S1**): (1) Providing

79 several alternative methods for most tasks to be suitable for complex and diverse

80 biological tissues and various analysis requirements. (2) Offering a factor

81 decomposition strategy for high-order spatial data to discover the major modes of

82 variations in spatial, time, sample or others. (3) Proposing a new method to combine

83 ligand-receptor expression and spatial locations to better infer short-range and

84 long-range cell communications. We also applied SOAPy to a wide range of public

85 datasets to demonstrate its general applicability and interpretability. SOAPy will be

86 one of the fundamental packages for spatial omics analysis in Python.

87

## Results

89 ### Overview of the SOAPy package

90      SOAPy is composed of four modules: **Data Preprocessing**, **Molecular Spatial**

91    **Dynamics** containing *Spatial Tendency* and *Spatiotemporal Pattern* analysis, **Cellular**

92    **Spatial Architecture** for analyzing *Spatial Proximity* and *Spatial Composition*, and

93    **Spatial Communication** that combines spatial distance, expression level and

94    interaction mechanism of ligand-receptors to infer cell interactions (**Figure 1**). In

95    addition, SOAPy provides rich visualization capabilities for all of the analysis

96    methods mentioned above.

97      The flexible **Data Preprocessing** module makes SOAPy suitable for various spatial

98    data, fitting with different modalities and different resolutions. To demonstrate the

99    utility of SOAPy, eight public datasets obtained from five state-of-the-art

100    technologies were analyzed (**Table S2**). These datasets involve multiple scenarios

101    with different molecular modalities (protein vs RNA), throughput (dozens to

102    genome-wide), spatial resolution (0.1 ~ 55μm), and in physiological and pathological

103    states.

104

### Spatial domain analysis recapitulates anatomic and pathological structures

106    Cells are not randomly distributed in tissues. They are self-organized into specific

107    structures to perform tissue functions. While in disease states, cells form abnormal

108    structures. The *Spatial Domain* analysis provides unsupervised (STAGATE) and

109    supervised (AUCell-LMI) methods to detect these structures (called spatial domains)

110    based on gene expression profiles and spatial locations[13,17,18].

111    We first tested STAGATE on Slide-seq V2 data for mouse olfactory bulb and 10x

112    Visium spatial transcriptomic data for human breast cancer[19]. Spatial domains

113    identified by STAGATE are highly consistent with the manual-labelled structures . It

114    successfully distinguishes truth anatomical structures (**Figure S1a**), malignant and

115    non-malignant tissues (**Figure S1b**, ARI=0.513), and more sophisticated pathological

116    stages (**Figure S1c**, ARI=0.580). Then we tested AUCell-LMI for finding local

117    structures with known signature genes, such as tertiary lymphoid structure (TLS)[20].

118    The results showed that supervised AUCell-LMI based on known TLS signature

119    could more accurate and more convenient identified the TLS region than

120    unsupervised STAGATE (**Figure S1d, e**). Taken together, Spatial domain analysis in

121    SOAPy could extract the interesting anatomic or pathological structures for

122    downstream analysis.

123

**Spatial tendency analysis finds genes associated with spatial structures**

125    The aim of *Spatial Tendency* analysis is to assess whether expression features were

126    influenced by spatial proximity to the region of interest (ROI). Expression features

127    could be gene expression, pathway activity, cell proportion and so on. The ROI is

128    defined by manual annotation or automatically detected by the *Spatial Domain*

129    analysis. Two kinds of methods, statistical test and regression model, are available for

130    tendency estimation in the *Spatial Tendency* module (Methods).

131    We used 10X Visium data of mouse dorsolateral prefrontal cortex (DLPFC)[21] as an

132    example to validate the feasibility of spatial tendency estimation (**Figure 2a**). The

133    sample is consisted of the grey matter of DLPFC (including six cortical layers) and

134    white matter (**Figure S2a**). To find genes whose expression changes along with the

135    distance to the white matter, three strategies were used and compared[22] (**Figure S2b,**

136    **c**): 1) cortical layers were divided into two regions and applied Wilcoxon test to

137    identify differential expressed genes; 2) cortical layers were separated to five

138    continuous zones for Spearman correlation test; 3) a polynomial regression model was

139    fitted between gene expressions and distances to the white matter. Some genes

140    identified by Wilcoxon test and Spearman correlation only express in few spots,

141    which may be the results of data sparsity instead of real biological differences (**Figure**

142    **S2e**). The regression model describes the continuous spatial variation of expression,

143    therefore it could find more complex spatial patterns than other methods[23], such as

144    nonlinear "low-high-low" spatial pattern (**Figure S2f**). Next, we analyzed the

145    expression patterns of 2857 significant (FDR < 0.05, range >0.3) genes identified by

146    polynomial regression. K-means clustering grouped them into 10 clusters (**Figure 2b**).

147    The gene clusters were compared with previously reported cortical layer specific

148    genes[24,25] (**Figure 2c**), showing high consistence. C3 is specifically highly expressed

149    near white matter regions; the expression peaks of C5, C8, C2, and C7 are at layer 6,

150    5, 4, 2, respectively (**Figure 2d**).

151        Considering that there are no predetermined structures in some scenarios, we added

152    three published methods (SpatialDE[10], SPARK[12], and SPARKX[11]) which identify

153    spatial variable genes (SVGs) but do not need a given ROI. Comparing these SVGs

154    methods with the above mentioned tendency estimation found shared and specific

155    genes among methods (**Figure S2d**). SVG methods were more inclined to show the

156    local differential expression of genes rather than the relationship with distance

157    (**Figure S2g**). Users can select sutiable methods based on their requirements.

158

159    **Tensor decomposition reveals the spatiotemporal patterns of gene expression**

160        With advances in omics techniques, spatial-resolved and time-series molecule

161    profilings are becoming available. One of the challenges is how to study the roles of

162    spatial effects and temporal effects simultaneously in biological questions. The

163    *Spatiotemporal Pattern* function in SOAPy employs tensor decomposition to extract

164    components from the three-order expression tensor ("Time-Space-Gene"), revealing

165    hidden patterns and reducing the complexity of data explanation.

166        Here, we used the mouse embryo development dataset from GeoMx Digital Spatial

167    Profiling (DSP)[7]. Limited by the availability of expression profiles, four time points

168    (E9, E11, E13, E15) and eight subtissues (Heart wall, Heart valve, heart trabecula,

169    Lung epithelium, Lung mesenchyme, Midgut epithelium, Midgut mesenchyme, and

170    Midgut neuron) from three organs were included in our analysis (**Figure 3a,b**).

171    Canonical Polyadic (CP) decomposition[26] was used to factorize the expression tensor

172    with 1000 high variable genes (a 4*8*1000 tensor) into seven factors, each of which

173    is the outer product of three vectors that contain the loadings for describing the

174    relative contribution of time, subtissues and genes (**Figure 3c**). We observed three

175    empirical spatiotemporal patterns based on the loadings of time and subtissues: pure

176    temporal variation (F1, F2), pure spatial variation (F3, F4), spatial and temporal

177    variation occur together (F5, F6, F7). We also conducted functional enrichment

178    analysis based on the loadings of genes for each factor (**Table S3**) and visualized the

179    typical genes in images (**Figure 3d**).

180 Genes in F1 (e.g. *Hbb-bh1*) highly express in heart and lung sub-tissues at E9, and

181 then gradually decrease in the later stages. Their expression pattern is consistent with

182 the enriched function "regulation of vasculature development". F1 indicates

183 co-development of heart and lung in the early embryo, which is consistent with

184 previous studies[27]. The expression of F2 genes (e.g. *Epcam*) increases significantly

185 since E11 in most sub-tissues of three organs, especially in the lungs. Expression of

186 F3 and F4 genes is stable along the developmental time. F4 genes highly express in

187 heart wall and heart trabecula, and their functions are enriched in cardiac cell

188 development as expected. Both F5 and F7 genes are enriched in midgut development.

189 F5 (e.g. *Psd*) slightly decreases from E11 to E15, while F7 (e.g. *Ndrg1*) increases

190 obviously from E11 to E15. F6 genes are specifically highly expressed in the heart

191 valve between E13-E15. In summary, the *Spatiotemporal Pattern* function in SOAPy

192 could reveal spatiotemporal specificity during development and other biological

193 processes.

194

195 **Spatial proximity analysis characterizes co-localization patterns between cell types**

196 Spatial architecture of cells is important for understanding the organization rules

197 from single cells to tissues[28–30]. SOAPy first constructs a cell/spot network

198 fromspatial locations; then implements two scenarios for deciphering spatial

199 architecture: *Spatial Proximity* analysis (including neighborhood and infiltration)

200 determines whethe two cell types or cell states within an image are significant

201 proximal; *Spatial Composition* analysis identifies multi-cellular niches that composed

202 by cell types with specific proportion.

203 We applied this analysis to a dataset of 41 triple-negative breast cancer (TNBC)

204 patients[5], which used multiplexed ion beam imaging by time-of-flight (MIBI-TOF) to

205 simultaneously quantify expression of 36 proteins in-situ at sub-cellular resolution.

206 Totally 211,649 cells were annotated to eight types (epithelial cell, endothelial cell,

207 mesenchymal cell, B, CD4 T, CD8 T, macrophage and other) based on the expression

208 of known protein markers.

209    First, *Spatial Neighborhood* analysis was performed to identify significantly

210    adjacent cell types compared to random perturbation[29]. **Figure 4a** illustrates the

211    neighborhood score of all samples for all cell type pairs, with positive or negative

212    scores corresponding to co-localization or avoidance. Different immune cells types

213    such as B, CD4 T, CD8 T and macrophage have significant co-localization in many

214    patients, which may relate with the formation of inflammatory foci (**Figure 4b**).

215    Endothelial and mesenchymal cells also prefer to co-locate together (**Figure 4c**).

216    Colocalization pattern of malignant epithelial cells and non-parenchymal cells were

217    highly heterogeneous across patients. Taking malignant epithelial cells and

218    mesenchymal cells as an example, samples with less than 200 mesenchymal cells

219    were filtered, others are subjected to *Spatial Infiltration* analysis. Samples with higher

220    and lower infiltration scores indicate mixed (e.g. sample 28) and compartmentalized

221    (e.g. sample 29) patterns between malignant epithelial cells and mesenchymal cells

222    respectively (**Figure 4 d-f**).

223

### Spatial composition analysis discovers multi-cellular niches

225    For *Spatial Composition* analysis of the TNBC dataset, the cell-cell network that

226    connected centroids of the cells within 100 pixels was built to capture the composition

227    pattern of more surrounding cells. Niche of each cell was presented by the proportion

228    of cell types of its surrounding cells, called I-niche. I-niches of 211,649 cells from 41

229    TNBC patients were clustered into 30 niche clusters, named C-niches (**Figure 5a,**

230    **Figure S3a**). The major cell types of the top two C-niches (C-niche13, C-niche18) are

231    mainly composed of malignant epithelial cells, and the percentages of other cell types

232    are less than 15%, showing the characteristics of tumor cell aggregation (**Figure 5b**).

233    Additionally, epithelial cells also form C-niches with other cell types. For example,

234    C-niche25 is composed of 38% epithelial cells, 31% mesenchymal cells, and 9%

235    macrophages; C-niche27 is composed of 23% epithelial, 28% endothelial, 10%

236    mesenchymal cells and 10% macrophages; C-niche15 is composed of 30% epithelial,

237    23% CD4 T, 13% CD8 T cells and 11% macrophages, suggesting different local

238    microenvironment exists among tumors (**Figure 5b**). We also observed four B cell

239  dominated C-niches (C-niche10, C-niche17, C-niche28, C-niche4) that may be related

240  to tertiary lymphoid structures. For example, sample 1 contains C-niche 10, 17, and

241  28 (**Figure 5c**). Around 80% of cells are B cells in C-niche10; C-niche17 majorly

242  consists of 52% B cells, 13% CD8 T cells, 10% CD4 T cells, and 11% epithelial cells;

243  C-niche28 majorly consists of 30% B cells, 10% CD8 T cells, and 37% epithelial

244  cells.

245  In order to investigate the combinational effects of non-parenchymal cell types and

246  niches on patient heterogeneity, the "Niche-CellType-Sample" tensor (30*7*41) was

247  factorized to four factors (**Methods**). All samples were clustered into five groups

248  according to the sample loadings in different factors (**Figure 5d**). Sample groups A, B,

249  C, and E have the highest loadings in factors 3, 2, 1, and 4, respectively. By checking

250  the loadings of cell types and niches in the major factors (**Figure S3b,c**), group B

251  corresponds to the above mentioned B cell enriched samples; group C is characterized

252  by niches with high proportion of mesenchymal cells; group E has niches consisted of

253  T cells and macrophages.

254  Furthermore, survival analysis was performed to explore the clinical indications of

255  niches. Eight c-niches were significantly related to survival time (P < 0.05, **Figure**

256  **S4**). For example, patients with a higher proportion of c-niche15 had a longer survival

257  time (**Figure 5e**). There also exists survival differences among the patient groups

258  identified by the "Niche-CellType-Sample" tensor decomposition, such as longer

259  survival time for group C patients that that of group D (**Figure 5f**). Taken together,

260  spatial composition analysis could find multi-cellular niches and yield insight into

261  how cells are organized into tissues.

262

263  **Ligand-receptor-mediated and spatial-constrained cell-cell communications**

264  The above spatial architecture analysis disregards interacting molecules and context,

265  while expression-based methods like CellphoneDB[31] and CellChat[32] infer cell-cell

266  communications by the expression of ligands and receptors (LRs) disregarding spatial

267  proximity. SOAPy develops a new method that simultaneously utilizes spatial

268  location and gene expression to calculate interaction scores (affinity and strength) and

269  then outputs significant LR interactions (**Figure 6a, Methods**). It can not only infers

270  short-range cell communication that relies on contact LRs to directly deliver signaling

271  between adjacent cells; but also infer long-range cell communication that does not

272  require cell–cell contact, rather depending on the diffusion of signaling molecules

273  from one cell to another after secretion[33,34].

274  The *Spatial Communication* module was applied to an ovarian cancer dataset

275  generated by the MERSCOPE platform, measuring 500 genes and 71,381 cells

276  (**Figure 6b**). Cells were classified and annotated into ten types or subtypes by Leiden

277  clustering algorithm. The spatial locations of epithelial cells C3 are very special,

278  which clearly separated with most of other cells. Therefore, our method did not find

279  significant contact LRs between epithelial cells C3 and other cell types. However,

280  CellChat, one of the most popular LR communication inference packages using

281  scRNA-Seq data, reported many LR interactions due to lack of spatial constrain

282  (**Tabls S4**), indicating lower false positives of our method.

283  We used endothelial cell as an example to present its short-range and long-range

284  communication partners. Fibroblasts and macrophages are located closest to

285  endothelial cell, while epithelial cell C3 and C4 are far away from endothelial cell

286  (**Figure 6c**). Consistently, fibroblasts have the largest number of contact LRs with

287  endothelial cells recognized by our algorithm, while there is no contact LRs for

288  distant cell types such as epithelial cells C2, C3, C4 and C5 (**Figure 6d**). For cell

289  types that are not spatially close to endothelial cells, *Spatial Communication* module

290  could infer secreted LRs that mediate long-range cell communications. The average

291  distance from epithelial cells C2 to the closet endothelial cells is significantly larger

292  than the average distance from fibroblasts to the closet endothelial cells (P <

293  3.9e-312). There are no contact LRs between epithelial cells C2 and endothelial cells

294  but 6 secreted LRs were identified (**Figure 6 d, e**).

295  Totally, we found 19 contact LRs and 66 secreted LRs that may play key roles in

296  short-range and long-range communication between endothelial cells and others

297  (**Figure 6f**). For example, COL1A1 (type I collagen) and its receptor ITGA1/ITGB1

298  (integrin α/β) highly express on spatial adjacent fibroblasts and endothelial cells, their

299    affinity and strength scores are significantly higher than random scores (**Figure 6g**).

300    Previous studies have reported that binding collagen to integrin may activate

301    downstream signaling pathways contributing to cancer progression[35]. VEGFB-FLT1

302    is an interesting LR pair for long-range communication between epithelial and

303    endothelial cells (**Figure 6h**). Epithelial cells C2 release ligand VEGFB, and

304    endothelial cells high express FLT1 (also known as VEGFR1). Their interaction may

305    promote tumor angiogenesis and are potential drug targets for anticancer therapy[36]. In

306    summary, SOAPy provide a new way to study spatial-constrained cell-cell

307    interactions and more accurately identify the related ligand-receptor pairs.

308

## Discussion

310    Tissue microenvironment is critical for understanding homeostasis, development,

311    regeneration and disease. Single-cell and spatial resolved omics are the most

312    promising technologies to investigate microenvironment. Tools for systematically

313    dissecting microenvironment and discover biologically important genes or spatial

314    cellular architecture are still falling behind, SOAPy just fill this gap. SOAPy contains

315    easy-to-use analysis modules for interpreting complex spatial microenvironments,

316    such as the spatial distribution patterns of genes and cells, dynamic changes along

317    with space and time, and cell-cell communications et al. In this article, we

318    demonstrated all SOAPy modules with various types of spatial omics data, and

319    provides complete tutorials to help users get started quickly.

320    The spatial distribution of genes or cells is associated with many elements, such as

321    time, interaction of cells, pathological foci, sample heterogeneity and so on. In the

322    face of these multi-dimensional data, how to extract important and meaningful

323    features is a key task. SOAPy utilizes tensor decomposition to discover the major

324    modes of variations from multi-dimensional data. The cases of mouse embryo

325    development and breast cancer showed that tensor decomposition in SOAPy is

326    powerful for interpret complex biological data. Another significant advantage of

327    SOAPy is the innovative *Spatial Communication* module. It combines spatial distance,

328   expression level and interaction mechanism of ligand-receptors to infer cell-cell

329   communication. The case of ovarian cancer showed that SOAPy could markedly

330   reduce false positives of interacting ligand-receptors compared to existing methods.

331       These advantages makes SOAPy differ from existing spatial data analysis tools.

332   Future extensions of SOAPy could be the integration of multi-modal spatial data to

333   delineate microenvironment, adaptation of methods from geoscience, network science,

334   or artificial intelligence to better extract biological meaningful spatial patterns. We

335   anticipate that SOAPy will be widely used by researchers to discover biological

336   insights from spatial omics data.

337

## Methods
338

### Data preprocessing
339

#### Data Import
340

341       The *Data Import* function converts data from different spatial omics technology to

342   a   unified   data   structure   that   contains   expression   profiles   of   molecules

343   (genes/proteins/metabolites) and location of cells/spots. Barcode-based data formats

344   can be read directly by passing in tables representing expression matrix and spatial

345   coordinate information. An image and a cell segmentation mask are provided for

346   imaging-based data, and the representation and coordinate matrix is extracted through

347   the tutorials on our website. We used the Scanpy toolkit[37] and generate Anndata data.

#### Spatial network construction
348

349       The *Spatial network* function provides four ways to build a neighborhood network

350   of cells/spots (Figure 1a). 1) Regular network; 2) KNN network that connects each

351   site with its K nearest neighbors; 3) Radius network that all cells/spots within the

352   given distance are connected; 4) Neighbor network based on Voronoi Diagram.

353

### Spatial domain identification
354

#### Unsupervised spatial domain identification: STAGATE
355

356       STAGATE is a graph attention autoencoder for spatial domain identification[13]. It

357    firstly integrates gene expression profiles and spatial location information to learn

358    low-dimensional latent embedding, and then assigns spatial domains by Louvain

359    clustering.

360    <u>**Supervised spatial domain identification: AUCell-LMI**</u>

361    To detect domains whose signature genes are already known, the score of signature

362    genes for each cell/spot is calculated by AUCell[38,39], and then local Moran index[17]

363    (LMI) is used to estimate the degree of spatial aggregation. LMI of cell/spot $i$ is

364    defined as:

$$I_i = \frac{x_i - \bar{x}}{s^2} * \sum_{j \in n_i} w_{ij}(x_j - \bar{x}) \#(1)$$

365    Where $x_i$ is the AUCell score of cell/spot $i$, $\bar{x} = \frac{1}{n}\sum_{i=1}^{n} x_i$, $j$ is any neighbor

366    cells/spots of $i$ based on K nearest neighbors, $w_{ij}$ is the spatial weight between $i$

367    and $j$. The P-value is calculated by permutation test and adjusted by

368    Benjamini-Hochberg method[40] to get the false discovery rate (FDR).

369    LMI of all cells/spots are illustrated by Moran scatterplot (Figure S1e). Each point

370    represents one cell/spot, the horizontal axis shows the normalized AUCell score, and

371    the vertical axis indicates the "spatial lag" which is calculated by spatial weighted

372    normalized score of neighboring sites. Sites with positive AUCell scores, positive

373    spatial lags, and low FDR were picked out as the targeted spatial domain.

374

375    **Spatial tendency analysis**

376    <u>**Definition of ROI and distance**</u>

377    Given a region of interest (ROI), the first step is to generate a binary mask file

378    (**Figure S2a**). Users can manually select ROI using tools like ImageJ to generate a

379    mask file, or get interesting cells/spots via SOAPy *Spatial domain* analysis and then

380    use SOAPy to create a mask file: Discrete cells/spots are converted to continuously

381    connected regions using a series of digital image processing steps in OpenCV library,

382    such as dilation, corrosion, removal of small connected components, and removal of

383    holes.

384    Next, the shortest distance from each cell/spot to the ROI boundary (contour) is

385    calculated. When an ROI contains multiple connected components, the closest

386    connected component is selected to calculate the distance[23].

$$d(i, C) = min_{p \in C} Enc(i, p) \#(2)$$

387    where $i$ is a cell/spot, $C$ is the boundary of ROI, and $p$ is any pixel on the

388    boundary. $Enc(\ )$ is a function of Euclidean distance. Distance with positive or

389    negative signs are used respectively to distinguish cells/spots located outside or inside

390    the ROI boundary. Then we can study the tendency of molecule expression along with

391    distance.

**<u>Identification of expression features with spatial tendency</u>**

393    SOAPy provides two statistical testing methods (**Figure S2b**): 1) wilcoxon rank

394    sum test to compare the molecule expression of cells/spots between two regions; 2)

395    spearman correlation between median expression and the rank of continuous zones.

396    To resolve more complex spatial tendency (e.g., nonlinear) or analyze ROIs without

397    prior hypothesis, SOAPy provides a parameter regression method (polynomial

398    regression model) and a non-parametric regression method (locally weighted liner

399    regression, LOESS).

400    Polynomial regression assumes that the output variable can be represented by the

401    sum of powers of the input variable.

$$Y = a_0 + \sum_{k=1}^{n} a_k \, d^k \#(3)$$

402    Where $d$ is the distance to the ROI; Y is the vector of molecule expression; $n$ is

403    the degree of the polynomial; $a_0$ is intercept; $a_k$ are slope coefficients. P-value is

404    calculated by F-test.

405    LOESS is a locally weighted polynomial regression method. Its core concept is to

406    fit weighted linear regression models with each data point using its surrounding data

407    points within the predefined window size and connect the centers of the regression

408    lines. $R^2$ (coefficient of determination) and residual standard deviation are estimated

409    to measure the goodness of fit.

410    Parameters used in both of the regression models could be customized and adjusted

411 based on the biological scenario and goodness of fit. To summarize the spatial
412 tendency of all molecules, the estimated expression values are fed into the K-means
413 clustering algorithm to obtain gene clusters with similar spatial expression tendency.

414

### Spatial architecture analysis

#### Spatial neighborhood analysis

417 For each paired cell types, a neighborhood score ($NS$) between cell type 1 ($ct1$) and
418 cell type 2 ($ct2$) is calculated as follows[29]:

$$NS_{ct1,ct2} = \frac{N_{ct1,ct2}}{N_{ct1,other} + N_{ct2,other}} \#(4)$$

419 where $N_{ct1,ct2}$ is the number of direct connections between $ct1$ and $ct2$, $N_{ct1,other}$
420 is the number of direct connections between $ct1$ and all other cell types.
421 Background distribution is generated from 1000 random permutations that fix the
422 numbers of $ct1$ and $ct2$ and randomly change their locations. P-value is the
423 proportion of permutations whose $NS$ is larger or smaller than the observed one,
424 which corresponds to either avoidance or interaction between $ct1$ and $ct2$.

#### Spatial infiltration analysis

426 An infiltration score ($IS$) is defined to present the degree of non-parenchymal
427 (immune or stromal) cells infiltration into malignant tissues:

$$IS_{m,np} = \frac{N_{m,np}}{min(N_{m,m}, N_{np,np})} \#(5)$$

428 where $N_{m,np}$ is the number of direct connections between malignant cells and
429 non-parenchymal cells. Sample with too few non-parenchymal cells are regarded as
430 cold tumor. Otherwise, larger infiltration score indicates more non-parenchymal cells
431 are mixed into malignant tissues, while smaller infiltration score suggests
432 non-parenchymal cells are more possible to be compartmentalized with malignant
433 tissues.

#### Spatial composition analysis

435 Given an index cell, niche is defined as the proportion of cell types for its
436 surrounding cells[41]. Taken all cells in one or more images, clustering algorithms like

437    K-means divides their niches into different clusters, called C-niches.

438

## Spatial-constrained cell-cell communication inference

440    Ligand-receptor (LR) pairs were obtained from the CellChat[32] package, in which

441    LR pairs were classified into contact and secreted based on their action mechanism.

442    We hypothesized that the contact LR pairs mediate short-range cell communications

443    while secreted LR pairs could mediate long-range cell communications. Therefore,

444    SOAPy infers cell communications based on the types of LR pairs and spatial

445    distance among cells (presented by a cell network). For short-range communication,

446    direct neighbors on Voronoi Diagram are connected to build a cell network. For

447    long-range communication, all cells within the given distance are connected to build a

448    cell network. Once the cell network is built, $Affinity$ and $Strength$ scores are

449    calculated for LRs on any two cell types. The LR pairs with $Affinity\ Pvalue <$

450    0.05 and $Strength > 4.0$ are considered to be significant. Paired cell types are

451    ranked based on the number of significant LRs.

452

### Cell-level ligand-receptor affinity score

454    The interaction of LR is variable among cells/spots at different spatial locations,

455    therefore we first define a cell-level ligand-receptor affinity score. Suppose a cell/spot

456    $i$ is a sender of ligand, cells/spots that have connection with $i$ and express the

457    matched receptor are receivers, the $Affinity\ score$ of ligand-receptor at location $i$

458    is defined as:

$$Affinity\ score_{l-r,i} = \sum_{j \in n_i} \frac{l_i * r_j}{1 + d_{i,j}}\ ,\ \ i\ as\ a\ ligand\ sender \#(6)$$

459    where $j$ is the cell/spot that connect to $i$ in the cell network; $l$ and $r$ are

460    expression levels of the ligand and receptor; $d$ is 0 for contact LR pairs or Euclidean

461    distance between $i$ and $j$ for secreted LR pairs. Similarly, when the cell/spot $i$ is a

462    receptor receiver, the $Affinity\ score$ of receptor-ligand at location $i$ is defined as:

$$Affinity\ score_{r-l,i} = \sum_{j \in n_i} \frac{r_i * l_j}{1 + d_{i,j}} \ , \ \ i\ as\ a\ receptor\ receiver \#(7)$$

463 The $Affinity\ Pvalue$ is obtained by random permutation:

$$Affinity\ Pvalue = \frac{\#m\{A^{(m)} \leq A^0, m = 1,2,\cdots,M\}}{M} \#(8)$$

464 $M$ is the total number of randomizations, $A^{(m)}$ is the $Affinity\ score$ under the

465 $m$-th randomization. Each randomization redistributes the expression values of the LR,

466 but keeps topology of the cell network. The affinity scores are calculated for all

467 cells/spots, and the P-values are used to find a subset of cells/spots at which the LR

468 exist interaction.

469

### CellType-level communication score

471     Suppose $ct1$ and $ct2$ are cell types that express ligands and receptors,

472 respectively. The $Affinity\ score$ between the ligand of $ct1$ and the receptor of

473 $ct2$ is the sum of cell-level scores:

$$Affinity\ score_{l,r,ct1,ct2} = \sum_{i \in ct1} \sum_{j \in n_i, ct2} \frac{l_i * r_j}{1 + d_{i,j}} \#(9)$$

474 $Affinity\ Pvalue$ is also calculated by random permutation, which randomly assign

475 a pseudo expression value to each cell/spot based on cell-type specific expression

476 distribution.

477     $Affinity$ reflects whether spatial connected $ct1$ and $ct2$ relatively more highly

478 express the LR genes. However, if the expression of ligand or receptor is too low in

479 $ct1$ or $ct2$ compared to other cell types, it is difficult to say that the LR is important

480 for cell communications; Additionally, If $ct1$ and $ct2$ are connected by too few

481 edges in the cell network, their communication may be false positive even affinity is

482 significant. To address these problems, another index 'strength' is added.

483 $Strength_{l,r,ct1,ct2}$ consists of two components: one is the relative expression level of

484 LR pairs on $ct1$ and $ct2$, and the other indicates the enrichment of real spatial

485 connections between $ct1$ and $ct2$. The detailed definition is as follows:

$$Strength_{l,r,ct1,ct2} = \left(\frac{\overline{exp}_{l,ct1}}{\overline{exp}_{l,all}} * \frac{\overline{exp}_{r,ct2}}{\overline{exp}_{r,all}}\right) * \left(\frac{2E}{1+E}\right) \quad \#(10)$$

$$E = \frac{edge_{ct1,ct2}}{\overline{edge}_{ct1,ct2}} \#(11)$$

486  where $\overline{exp}_{l,ct1}$ and $\overline{exp}_{l,all}$ are the average expression of ligand in $ct1$ and in all cells;

487  $edge_{ct1,ct2}$ and $\overline{edge}_{ct1,ct2}$ are the real and expected number of connections between

488  $ct1$ and $ct2$; $E$ is the ratio of real and expected numbers. To constrain the range of

489  $E$ and make the result more stable, a Hill function transforms $E$ into a range of $(0, 2)$

490  and keeps the transformed $E$ is still 1 when the number of real and expected

491  connections are equal.

492

493  **Tensor decomposition**

494  To discover the major modes of variation in the high-order spatial data, such as the

495  "Time-Space-Gene" tensor or "Niche-CellType-Sample" tensor, SOAPy provides

496  interface functions to conveniently build tensors from AnnData objects and then

497  decomposes tensors into several latent factors or components.

498  SOAPy implements two tensor decomposition methods, CANDECOMP

499  /PARAFAC (CP) and Tucker decomposition[26,42]. Moreover, SOAPy supports

500  non-negative constraints to make the factors more interpretable. Take non-negative

501  CP[43] as an example, an n-order tensor X is expressed as the weighted sum of R

502  (user-defined number of factors) rank-one tensors:

$$X \approx \sum_{r=1}^{R} \lambda_r a_r^{(1)} \circ a_r^{(2)} \circ \dots \circ a_r^{(n)} \#(12)$$

503  where $\lambda$ is the weight of each factor; $a_r^{(k)}$ is the non-negative loading values of k-th

504  variable in the r-th factor, indicating the relative contribution of variables to factors.

505  Each factor is the outer product of the loading vectors.

506

507  **Availability**

508    All data and code that produced the findings of the study, including all main and

509    supplemental figures, are available at https://github.com/LiHongCSBLab/SOAPy.

510

## Acknowledgements

512    We acknowledge Andrew E. Teschendorff (from Shanghai Institute of Nutrition

513    and Health, Chinese Academy of Sciences) for his advice on our manuscript. We

514    thank Bihan Shen, Xi Yan (from Shanghai Institute of Nutrition and Health, Chinese

515    Academy of Sciences) and Biao Liu (from Center for Excellence in Molecular Cell

516    Sciences, Chinese Academy of Sciences), for their help on programing and result

517    interpretation.

518

## Funding

524

## References

526    1.    Method of the Year 2020: spatially resolved transcriptomics. *Nat. Methods* **18**, 1–1
527    (2021).
528    2.    Rao, A., Barkley, D., França, G. S. & Yanai, I. Exploring tissue architecture using spatial
529    transcriptomics. *Nature* **596**, 211–220 (2021).
530    3.    Moses, L. & Pachter, L. Museum of spatial transcriptomics. *Nat. Methods* **19**, 534–546
531    (2022).
532    4.    Salmén, F. *et al.* Barcoded solid-phase RNA capture for Spatial Transcriptomics profiling
533    in mammalian tissue sections. *Nat. Protoc.* **13**, 2501–2534 (2018).
534    5.    Keren, L. *et al.* A Structured Tumor-Immune Microenvironment in Triple Negative
535    Breast Cancer Revealed by Multiplexed Ion Beam Imaging. *Cell* **174**, 1373-1387.e19 (2018).
536    6.    Keren, L. *et al.* MIBI-TOF: A multiplexed imaging platform relates cellular phenotypes
537    and tissue structure. *Sci. Adv.* **5**, eaax5851 (2019).
538    7.    Merritt, C. R. *et al.* Multiplex digital spatial profiling of proteins and RNA in fixed tissue.
539    *Nat. Biotechnol.* **38**, 586–599 (2020).
540    8.    Schapiro, D. *et al.* MCMICRO: a scalable, modular image-processing pipeline for
541    multiplexed tissue imaging. *Nat. Methods* **19**, 311–315 (2022).

9.  Cang, Z. *et al.* Screening cell–cell communication in spatial transcriptomics via collective optimal transport. *Nat. Methods* **20**, 218–228 (2023).

10. Svensson, V., Teichmann, S. A. & Stegle, O. SpatialDE: identification of spatially variable genes. *Nat. Methods* **15**, 343–346 (2018).

11. Zhu, J., Sun, S. & Zhou, X. SPARK-X: non-parametric modeling enables scalable and robust detection of spatial expression patterns for large spatial transcriptomic studies. *Genome Biol.* **22**, 184 (2021).

12. Sun, S., Zhu, J. & Zhou, X. Statistical analysis of spatial expression patterns for spatially resolved transcriptomic studies. *Nat. Methods* **17**, 193–200 (2020).

13. Dong, K. & Zhang, S. Deciphering spatial domains from spatially resolved transcriptomics with an adaptive graph attention auto-encoder. *Nat. Commun.* **13**, 1739 (2022).

14. Dries, R. *et al.* Giotto: a toolbox for integrative analysis and visualization of spatial expression data. *Genome Biol.* **22**, 78 (2021).

15. Palla, G. *et al.* Squidpy: a scalable framework for spatial omics analysis. *Nat. Methods* **19**, 171–178 (2022).

16. Pham, D. *et al. stLearn: integrating spatial location, tissue morphology and gene expression to find cell types, cell-cell interactions and spatial trajectories within undissociated tissues.* http://biorxiv.org/lookup/doi/10.1101/2020.05.31.125658 (2020) doi:10.1101/2020.05.31.125658.

17. Anselin, L. Local Indicators of Spatial Association-LISA. *Geogr. Anal.* **27**, 93–115 (2010).

18. Jong, P., Sprenger, C. & Veen, F. On Extreme Values of Moran's I and Geary's c. *Geogr. Anal.* **16**, 17–24 (2010).

19. Stickels, R. R. *et al.* Highly sensitive spatial transcriptomics at near-cellular resolution with Slide-seqV2. *Nat. Biotechnol.* **39**, 313–319 (2021).

20. Meylan, M. *et al.* Tertiary lymphoid structures generate and propagate anti-tumor antibody-producing plasma cells in renal cell cancer. *Immunity* **55**, 527-541.e5 (2022).

21. Pardo, B. *et al.* spatialLIBD: an R/Bioconductor package to visualize spatially-resolved transcriptomics data. *BMC Genomics* **23**, 434 (2022).

22. Bardou, P., Mariette, J., Escudié, F., Djemiel, C. & Klopp, C. jvenn: an interactive Venn diagram viewer. *BMC Bioinformatics* **15**, 293 (2014).

23. Hildebrandt, F. *et al.* Spatial Transcriptomics to define transcriptional patterns of zonation and structural components in the mouse liver. *Nat. Commun.* **12**, 7046 (2021).

24. He, Z. *et al.* Comprehensive transcriptome analysis of neocortical layers in humans, chimpanzees and macaques. *Nat. Neurosci.* **20**, 886–895 (2017).

25. Maynard, K. R. *et al.* Transcriptome-scale spatial gene expression in the human dorsolateral prefrontal cortex. *Nat. Neurosci.* **24**, 425–436 (2021).

26. Kolda, T. G. & Bader, B. W. Tensor Decompositions and Applications. *SIAM Rev.* **51**, 455–500 (2009).

27. Peng, T. *et al.* Coordination of heart and lung co-development by a multipotent cardiopulmonary progenitor. *Nature* **500**, 589–592 (2013).

28. Schapiro, D. *et al.* histoCAT: analysis of cell phenotypes and interactions in multiplex image cytometry data. *Nat. Methods* **14**, 873–876 (2017).

586   29.  Bäckdahl, J. *et al.* Spatial mapping reveals human adipocyte subpopulations with distinct
587       sensitivities to insulin. *Cell Metab.* **33**, 1869-1882.e6 (2021).

588   30.  Yuan, Z. *et al.* SOTIP is a versatile method for microenvironment modeling with spatial
589       omics data. *Nat. Commun.* **13**, 7330 (2022).

590   31.  Efremova, M., Vento-Tormo, M., Teichmann, S. A. & Vento-Tormo, R. CellPhoneDB:
591       inferring cell–cell communication from combined expression of multi-subunit
592       ligand–receptor complexes. *Nat. Protoc.* **15**, 1484–1506 (2020).

593   32.  Jin, S. *et al.* Inference and analysis of cell-cell communication using CellChat. *Nat.*
594       *Commun.* **12**, 1088 (2021).

595   33.  Armingol, E., Officer, A., Harismendy, O. & Lewis, N. E. Deciphering cell–cell
596       interactions and communication from gene expression. *Nat. Rev. Genet.* **22**, 71–88 (2021).

597   34.  Cheng, J., Yan, L., Nie, Q. & Sun, X. *Modeling and inference of spatial intercellular*
598       *communications and multilayer signaling regulations using stMLnet.*
599       http://biorxiv.org/lookup/doi/10.1101/2022.06.27.497696       (2022)
600       doi:10.1101/2022.06.27.497696.

601   35.  Xu, S. *et al.* The role of collagen in cancer: from bench to bedside. *J. Transl. Med.* **17**,
602       309 (2019).

603   36.  Fischer, C., Mazzone, M., Jonckx, B. & Carmeliet, P. FLT1 and its ligands VEGFB and
604       PlGF: drug targets for anti-angiogenic therapy? *Nat. Rev. Cancer* **8**, 942–956 (2008).

605   37.  Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression
606       data analysis. *Genome Biol.* **19**, 15 (2018).

607   38.  Van De Sande, B. *et al.* A scalable SCENIC workflow for single-cell gene regulatory
608       network analysis. *Nat. Protoc.* **15**, 2247–2276 (2020).

609   39.  Fang, Z., Liu, X. & Peltz, G. GSEApy: a comprehensive package for performing gene set
610       enrichment analysis in Python. *Bioinformatics* **39**, btac757 (2023).

611   40.  Haynes, W. Benjamini–Hochberg Method. in *Encyclopedia of Systems Biology* (eds.
612       Dubitzky, W., Wolkenhauer, O., Cho, K.-H. & Yokota, H.) 78–78 (Springer New York, 2013).
613       doi:10.1007/978-1-4419-9863-7_1215.

614   41.  Goltsev, Y. *et al.* Deep Profiling of Mouse Splenic Architecture with CODEX
615       Multiplexed Imaging. *Cell* **174**, 968-981.e15 (2018).

616   42.  Zhou, G., Cichocki, A., Zhao, Q. & Xie, S. Efficient Nonnegative Tucker
617       Decompositions: Algorithms and Uniqueness. *IEEE Trans. Image Process.* **24**, 4990–5003
618       (2015).

619   43.  Shashua, A. & Hazan, T. Non-negative tensor factorization with applications to statistics
620       and computer vision. in *Proceedings of the 22nd international conference on Machine*
621       *learning  - ICML '05* 792–799 (ACM Press, 2005). doi:10.1145/1102351.1102451.

622

623

624

625

626

627   **Figures**

**Figure 1.** Schematic diagram of SOAPy. **a,** "Data Preprocessing" module that imports data, generates cell network and identifies spatial domains. Data from different spatial omics technologies are converted to a unified data structure. Cell network can be built by any of the four methods. Spatial domains are inferred by unsupervised learning from expression and morphological data, or supervised classification based on the expression of signature genes. **b,** "Molecular Spatial Dynamics" module. Spatial tendency analysis finds genes or cells whose expression change with spatial distance to the given region. **c,** Spatiotemporal Pattern analysis performs a tensor decomposition to discover the major modes of variation in space and time. **d,** "Cellular Spatial Architecture" module. Neighborhood and infiltration analysis find spatial proximal cell types. Spatial composition reveals conserved niches in which surrounding cells of the index cell are consisted of specific proportion of cell types. **e,** Innovative "Spatial Communication" module that combine spatial distance, expression level and action mechanism of ligand-receptors (LRs) to infer cell interactions. The contact and secreted LRs are considered for short-range and long-range cell communications, respectively. Results at cell/spot level indicate the heterogeneous interaction among different spatial locations, they are further integrated to cell type-level to report significant LRs for any two cell types.

**Figure 2.** Spatial tendency analysis finds genes associated with spatial structures. **a,** HE image of a human dorsolateral prefrontal cortex (DLPFC) sample. Regions of white matter (WM) and six neuronal layers (L6 to L1) are labeled on the image. **b,** Regression curves between gene expression and the distance to WM. Polynomial regression models were fitted to identify genes whose expression varied along with the distance to WM boundary. These genes were grouped into 10 clusters by K-means clustering algorithm. Each curve present a cluster of genes with similar spatial expression tendency. Zero at the horizontal axis indicates the outer boundary of WM. **c,** Association between gene clusters and previously reported layer specific genes. Each row corresponds to a prior gene-list that specifically expresses in one neuronal layer[24]. Each red unit indicates the cluster of genes (column) is enriched in the prior

658    gene-list (row). **d**, Spatial distributions and fitted curves of the representative genes.

659

660    **Figure 3.** Tensor decomposition reveals the spatiotemporal patterns of gene expression

661    during mouse embryo development. **a,** The spatiotemporal dataset of mouse

662    development is represented by a three-order tensor (4 time points * 8 sub-tissues *

663    1000 highly variable genes), and then it's decomposed into seven latent factors. **b,**

664    Representative spatial locations of sub-tissues at four time points. Each spot in the

665    subtissues represents an ROI. **c,** Loading vectors of space and time for each factor

666    obtained by tensor decomposition. Higher loading values indicates larger contribution

667    of sub-tissues or time points to the expression variation of this factor. **d,** Spatial

668    expression of example genes for each factor. The contours of heart, lung and midgut

669    are colored by red, blue and green curves. ROIs of gene expression are presented by

670    cyan points. The darker the cyan color, the higher the gene expression level.

671

672    **Figure 4.** Spatial proximity analysis characterizes cellular co-localization patterns.

673    The triple negative breast cancer (TNBC) dataset contains 41 samples and 7 cell types.

674    **a,** Heatmap showing the neighborhood scores of any two cell types in all TNBC

675    samples. **b,** A representative sample with strong co-localization among immune cells.

676    **c.** A representative sample with strong co-localization between endothelial and

677    mesenchymal cells**. d,** The red bars show the number of mesenchymal cells and the

678    blue bars show the infiltration score of mesenchymal cells into malignant epithelial

679    cells. **e,** A representative sample with low infiltration score, suggesting

680    compartmentalization between mesenchymal cells and tumor tissues. **f,** A

681    representative sample with high infiltration score, suggesting mixture of mesenchymal

682    cells into malignant epithelial cells.

683

684    **Figure 5.** Spatial composition analysis discovers multi-cellular niches in TNBC

685    samples. **a.** Heatmap on the left shows the composition of neighbor cells in each

686    C-niche. The right barplot shows the number of cells belonging to each C-niche. **b,**

687    Representative samples of different C-niches, characterizing tumor cell aggregation

688 and different local microenvironment of tumors. **c**, The left image shows an example

689 sample that has B cell dominated C-niches (the region of red box). Cells are colored

690 by C-niches. 'other' are low-frequent c-niches whose proportion is less than 2%.

691 Right images are amplified views of three representative C-niches. Black or gray cell

692 contours indicate cells belonging to or not belonging to the C niche. The fill colors of

693 cells represent cell types involved in the definition of the C-niche. **d**, Heatmap

694 showing the loading values and clusters of samples. The three-order

695 'Niche-CellType-Sample' tensor was decomposed to four latent factors (**Figure S3b,**

696 **c**). Samples are clustered into five groups according to their loading vectors. **e,**

697 Survival curves stratified by the proportion of C-niche-15. **f**, Comparison of survival

698 curves between two groups of patients.

699

700 **Figure 6**. Ligand-receptor-mediated and spatial-constrained cell-cell communications.

701 **a**, The brief flow chart of our method. Short-range interaction is mediated by contact

702 LRs on neighbor cells, long-range interaction is mediated by secreted LRs on cells

703 within the given radius. Two new metrics, affinity and strength, are defined to

704 estimate the probability of LR interactions in any two cell types. Only when both

705 metrics are high, the LR is significant to mediate the interactions of these two cell

706 types. **b**, MERSCOPE data from an ovarian cancer sample. **c,** Barplot showing the

707 shortest distance from other cell to the closest endothelial cell. **d, e**. Short-range and

708 long-range cell communication networks between endothelial cells and other cell

709 types. Edges in d and e are the number of contact and secreted LRs. Edge width is the

710 number of significant ligand-receptor pairs (affinity P-value < 0.05, strength > 4). **f,**

711 Dot plot with ligand-receptor interactions corresponding to d and e. Each row

712 indicates a ligand-receptor pair, with the first and the second genes representing a

713 ligand and a receptor, respectively. Dot size indicates P-value of affinity. Color

714 indicates the strength score. **g**, An example of contact LR that mediates the

715 communication between spatially colocalized fibroblast and endothelial cells.

716 COL1A1 is the ligand on sender fibroblast cells, and ITGA1/ITGB1 is the receptor on

717 receiver endothelial cells. Expression was scaled to the range of 0-1 by normalization.

718 **h**, An example of secreted LR, corresponding to the communication between spatially

719 separate epithelial and endothelial cells. VEGFB is the ligand on sender

720    epithelial-hypoxia cells, and FLT1 is the receptor on receiver endothelial cells.

721

## Supplementary Information

723    **Figure S1.** Spatial domain analysis recapitulates anatomic and pathological structures.

724    **a**, Anatomical structure of mouse olfactory bulb (Slide-seq V2 data) and domains

725    identified by STAGATE. **b-c**, Expert-annotated pathological regions of a breast

726    cancer sample (10x Visium), and the estimated 2-class and 19-class domains based on

727    the results of by STAGATE. **d**, Expert-annotated tertiary lymphoid structure (TLS) on

728    a kidney cancer sample (10x Visium), and the estimated TLS by the AUCell-LMI

729    method. **e**, Moran scatterplot. The x-axis is the Z-transformed AUC, which presents

730    the activity for the signature genes of TLS. The y-axis is the spatial weighted

731    normalized AUC scores of neighboring locations. Hotspot presented by red points

732    (FDR < 0.05, x > 0, y > 0) is regarded as tertiary lymphoid structure.

733

734    **Figure S2.** Spatial tendency analysis. **a**, Steps of image per-processing to generate a

735    binary mask file for the given region of interest (ROI). **b**, Illustration of three spatial

736    tendency analysis strategies: wilcoxon test, spearman correlation, and regression. **c**,

737    Venn diagram shows the overlap of top 1000 genes (FDR q-value < 0.05) obtained

738    from three spatial tendency analysis strategies. There are 380 overlapped genes, 352,

739    209 and 227 genes uniquely identified by a method (**Figure S2c**). **d**, Intersection plot

740    showing the agreement for seven methods. Four methods estimate the tendency of

741    gene expressions changing with the distance to a given region: wilcoxon test,

742    spearman correlation, polynomial regression and LOESS regression. Other three

743    methods identify spatially variable genes (SVGs) whose expressions depend on their

744    spatial locations: SPARK, SPARKX and SpatialDE. The top-ranked genes with equal

745    number obtained from each method were compared. Genes of LOESS are ranked by

746    R-square, and genes of the remaining methods are ranked by FDR values. **e-g**,

747    Representative genes identified by different kinds of methods. **e**: MIAT that is

748    significant by Wilcoxon test and Spearman correlation analysis but not significant by

749    regression methods; **f**: PVALB that is significant by regression methods; **g**:

750    Expression of TFF1 is spatially variable but do not show tendency of change with the

751    distance to WM.

752

753    **Figure S3**. **a**, Heatmap showing the proportion of niches in all TNBC samples. **b-c**,

754    Loadings of cell types and niches obtained from the decomposition of

755    "Niche-CellType-Sample" tensor.

756

757    **Figure S4**. Results of survival analysis.

758

759    **Table S1**. Comparison with existing tools of spatial omics data analysis.

760

761    **Table S2**. Examples datasets that were used in this study.

762

763    **Table S3**. Enriched functional terms by gene set enrichment analysis. Genes were

764    pre-ranked based on the loading values of each factor obtained from tensor

765    ("Time-Space-Gene") decompositon.

766

767    **Table S4.** Predicted LR interactions between spatial-separated epithelial cells C3 and

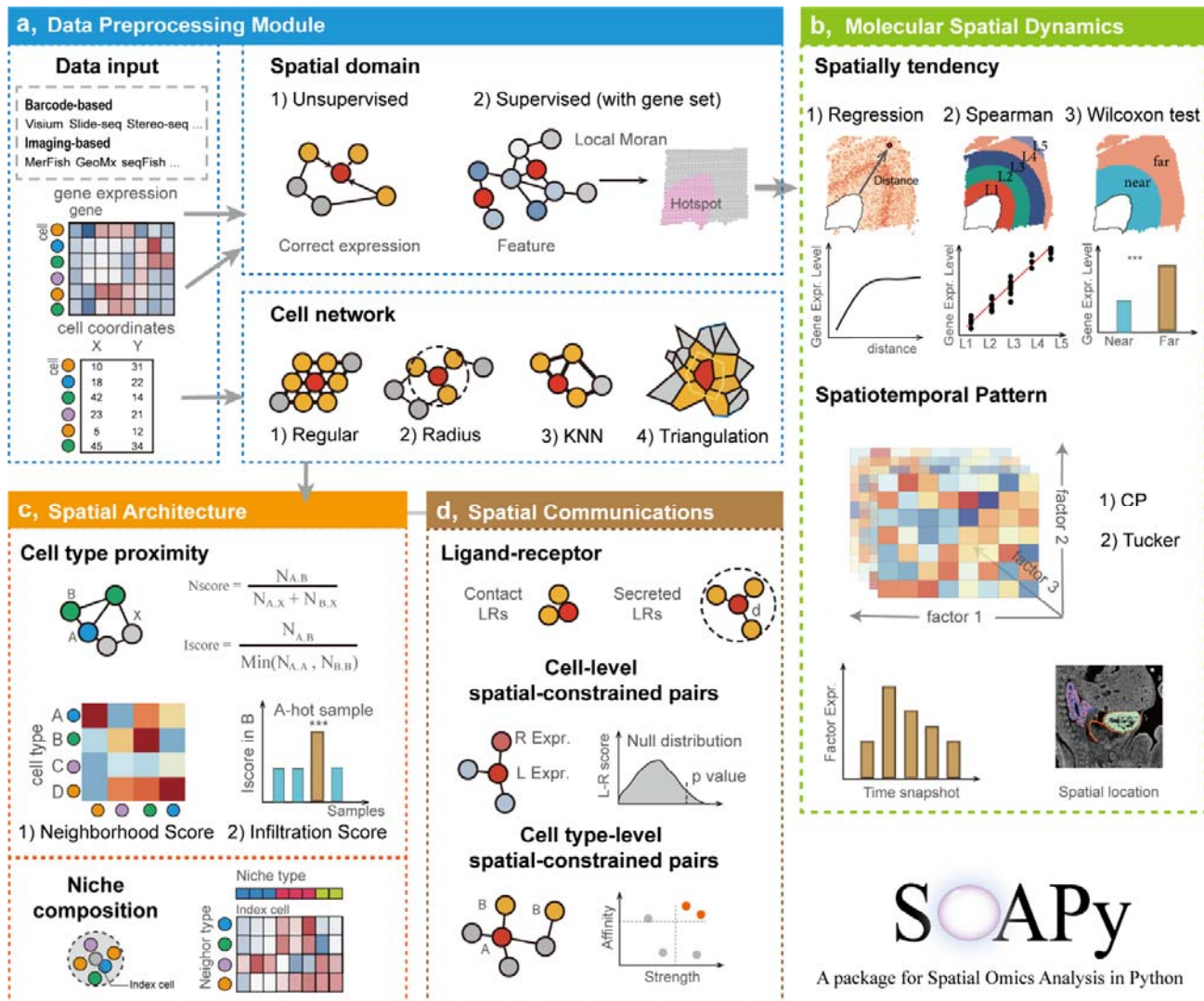768    other cell types by CellChat and SOAPy.
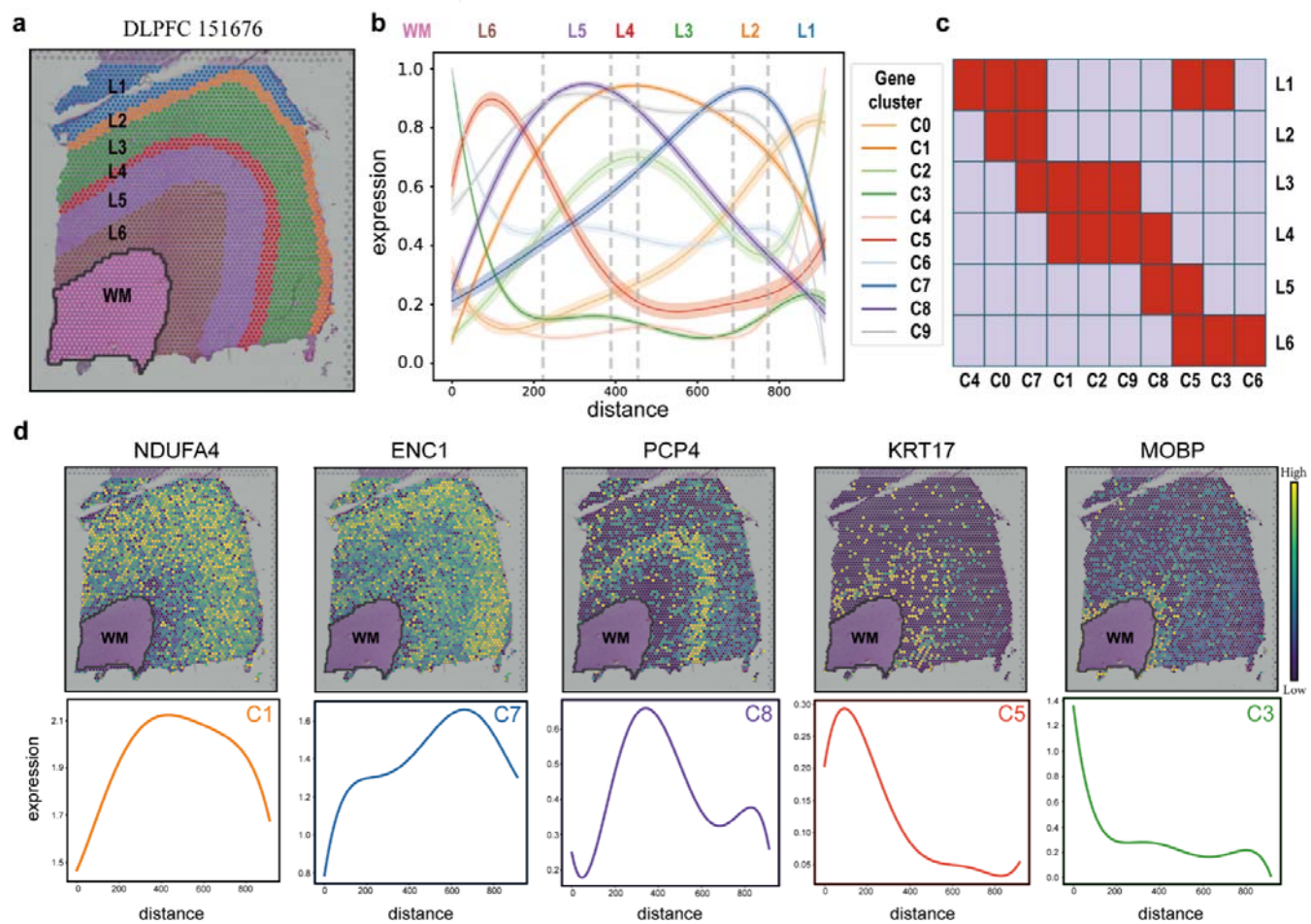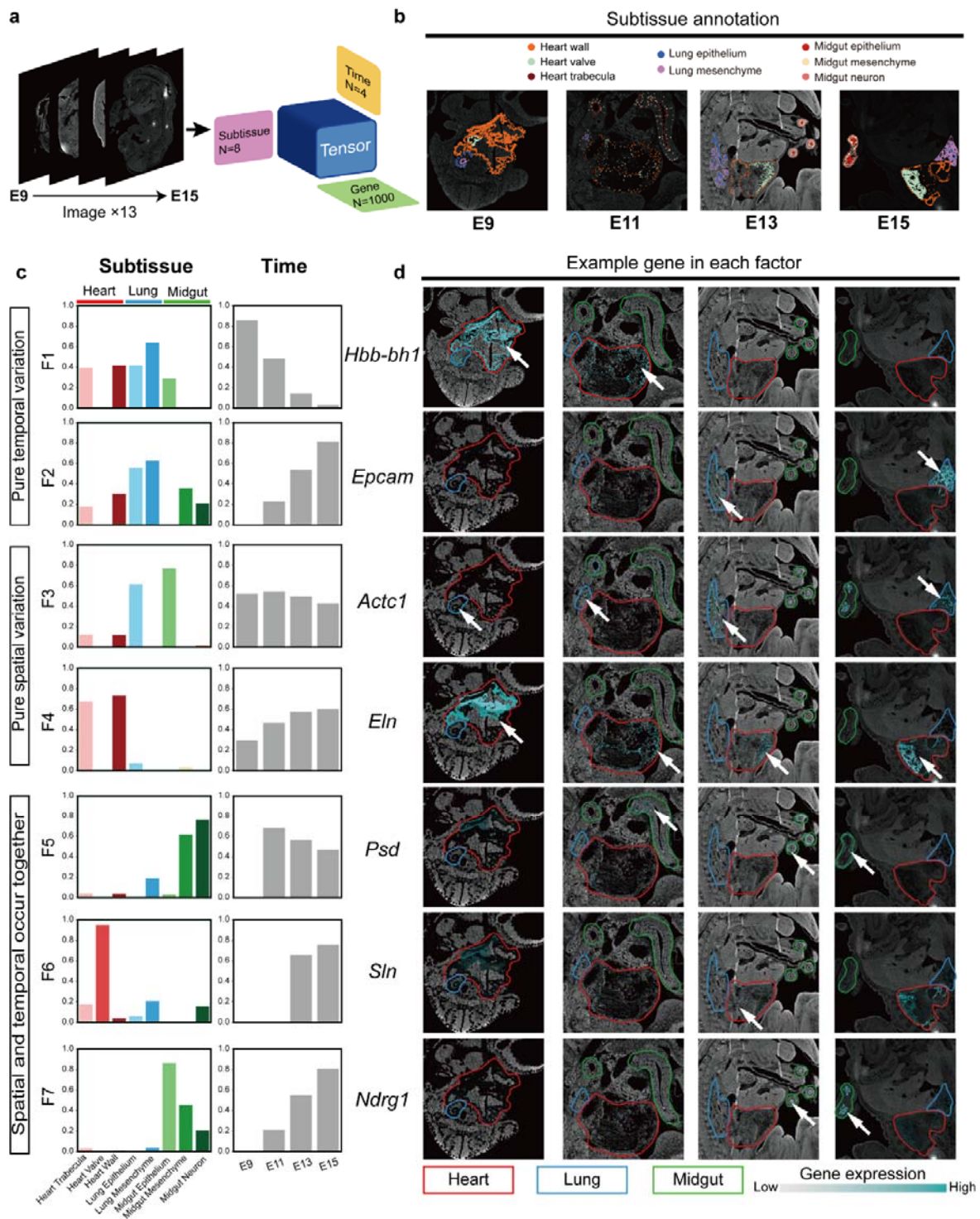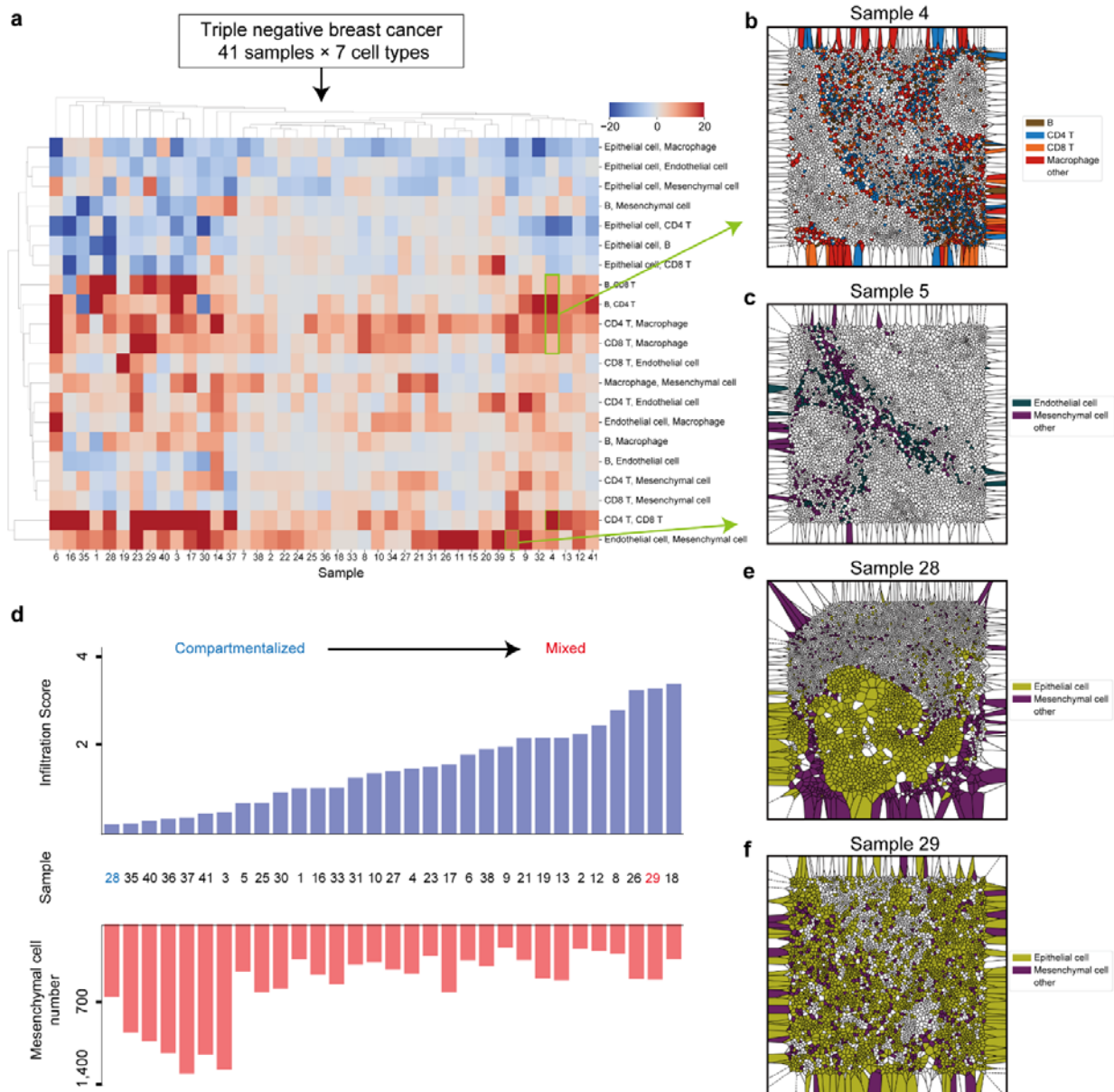
# Figure 1

# Figure 2

# Figure 3

# Figure 4

# Figure 5

# Figure 6