

# Decoding Spatial Tissue Architecture: A Scalable Bayesian Topic Model for Multiplexed Imaging Analysis

Xiyu Peng<sup>1,2\*</sup>, James W. Smithy<sup>3</sup>, Mohammad Yosofvand<sup>1</sup>,  
Caroline E. Kostrzewa<sup>1</sup>, MaryLena Bleile<sup>1</sup>, Fiona D. Ehrich<sup>1</sup>, Jasme Lee<sup>1</sup>,  
Michael A. Postow<sup>3</sup>, Margaret K. Callahan<sup>4</sup>, Katherine S. Panageas<sup>1\*</sup>,  
Ronglai Shen<sup>1\*</sup>

<sup>1</sup>Department of Epidemiology and Biostatistics, Memorial Sloan Kettering Cancer Center, New York, 10065, NY, USA.

<sup>2</sup>Department of Statistics, Texas A&M University, College Station, 77843, TX, USA.

<sup>3</sup>Department of Medicine, Memorial Sloan Kettering Cancer Center, New York, 10065, NY, USA.

<sup>4</sup>Neag Comprehensive Cancer Center, UConn Health, Farmington, 06030, CT, USA.

\*Corresponding author(s). E-mail(s): [pengx@stat.tamu.edu](mailto:pengx@stat.tamu.edu); [panageak@mskcc.org](mailto:panageak@mskcc.org); [shenr@mskcc.org](mailto:shenr@mskcc.org);

Contributing authors: [smithyj@mskcc.org](mailto:smithyj@mskcc.org); [yosofvm@mskcc.org](mailto:yosofvm@mskcc.org); [kostrzc@mskcc.org](mailto:kostrzc@mskcc.org); [bleilem@mskcc.org](mailto:bleilem@mskcc.org); [ehrichf@mskcc.org](mailto:ehrichf@mskcc.org); [leej22@mskcc.org](mailto:leej22@mskcc.org); [PostowM@mskcc.org](mailto:PostowM@mskcc.org); [mcallahan@uchc.edu](mailto:mcallahan@uchc.edu);

## Abstract

Recent progress in multiplexed tissue imaging is advancing the study of tumor microenvironments to enhance our understanding of treatment response and disease progression. Cellular neighborhood analysis is a popular computational approach for these complex image data. Despite its popularity, there are significant challenges, including high computational demands that limit feasibility for large-scale applications and the lack of a principled strategy for integrative analysis across images. This absence hampers the precise and consistent identification of spatial features and tracking of their dynamics over disease progression. To overcome these challenges, we introduce *SpaTopic*, a spatial topic model designed to decode high-level spatial architecture across multiplexed tissue images. This algorithm integrates both cell type and spatial information within a topic modelling framework, originally developed for natural language processing and adapted for computer vision. Spatial information is incorporated into the flexible design of documents, representing densely overlapping regions in images. The model employs an efficient collapsed Gibbs sampling algorithm for both statistical and computational inference. We benchmarked the performance against five state-of-the-art algorithms through various case studies using different single-cell spatial transcriptomic and proteomic imaging platforms across different tissue types. Our findings demonstrate that *SpaTopic* consistently identifies biologically and clinically significant spatial “topics” such as tertiary lymphoid structures (TLSs) and tracks dynamic changes in spatial features over disease progression. Its computational efficiency and broad applicability across various molecular imaging platforms will enhance the analysis of large-scale tissue imaging datasets.

**Keywords:** Multiplexed tissue imaging, Spatial molecular profiling, Tumor microenvironment, Cellular neighborhoods, Topic models

## 40 Introduction

41 Recent advancements in multiplexed tissue imaging allow the profiling of RNA and protein expression  
42 in situ across thousands to millions of single cells within a whole-slide tissue context [1–5]. These  
43 technologies generate high-dimensional molecular imaging data, offering significant opportunities for a  
44 spatially resolved understanding of cellular heterogeneity and organization within tissues. Compared  
45 to other single-cell technologies (such as single-cell RNA-seq, flow cytometry), multiplexed imaging  
46 provides unique opportunities to examine spatial patterns of diverse cell types and characterize the  
47 tissue microenvironment of interest, which may play an essential role in understanding disease progres-  
48 sion, tissue development, and mechanisms of treatment response [1, 2, 4–7]. One recent discovery in  
49 cancer, partly enabled by multiplexed spatially resolved omics data, is the presence of tertiary lymphoid  
50 structures (TLSs) in tumor tissues and its role in the adaptive antitumor immune response [8–11]. TLSs  
51 have been identified in a wide range of human cancers [9] and have demonstrated a promising positive  
52 association with improved outcomes in cancer patients who underwent immunotherapy [8].

53  
54 While promising, the complex cellular architecture revealed by whole-slide multiplexed tissue imaging  
55 presents significant analytical challenges. Pathology images of tissue samples affected by certain diseases,  
56 such as cancer, are particularly complex, displaying abnormal cellular structures and significant varia-  
57 tion between tumor samples. Currently, most analyses focus on individual images, examining elements  
58 such as cell densities and inter-cellular distances [1, 2, 6], or conducting basic spatial domain analyses  
59 that primarily focus on binarized tissue compartments, such as tumor versus stroma [12]. Associating  
60 these features with outcomes requires manual and heuristic aggregation across images. While promising,  
61 a significant hurdle in spatial pattern analysis is deciphering biologically and clinically relevant patterns  
62 from the complex architecture within tissue across various slides.

63  
64 In recent literature, cell neighborhood (niche) analysis is emerging as a popular approach. This analysis  
65 pipeline typically consists of two primary steps by first identifying neighborhood features for each  
66 single cell using either a K-nearest-neighbor (KNN) graph or a defined radius, and then applying a  
67 clustering algorithm, such as k-means, Louvain, or Latent Dirichlet Allocation (LDA) [2, 6, 7, 13–15].  
68 Seurat v5 [16] for instance, clusters cells using k-means based on similar cell type compositions, offering  
69 a straightforward niche analysis method. There are different variants of the approach depends on how  
70 to incorporate spatial information into the clustering process. UTAG [13] averages marker expression  
71 within the neighborhood for clustering, while BankSY [17] further refines this by combining local mean  
72 expression with individual cell expression. Spatial-LDA [14] incorporates spatial priors into clustering to  
73 allow proximity-closed cells to share similar cell neighborhoods. More recently, graph neural networks  
74 have been employed to discern cell neighborhood patterns, such as CytoCommunity [18]. However, deep  
75 learning methods like CytoCommunity require significant computational resources, posing challenges for  
76 individual labs, particularly for large-scale image analysis. Other studies adapt computational methods  
77 designed for spatial transcriptomics to analyze tissue imaging data [19–21], such as those intended for  
78 10x Visium, face limitations due to high computational costs [17] and are generally restricted to single  
79 tissue sections with fewer spots [13, 19]. These methods struggle with large-scale images, like whole-slide  
80 multiplexed data containing millions of cells, and are challenging to adapt for modern imaging platforms  
81 like Nanostring CosMx and 10x Xenium.

82  
83 Highly interpretable and scalable machine learning methods are in great need for analyzing molecular  
84 tissue imaging data. In this work, we propose *SpaTopic*, a Bayesian topic model designed to identify  
85 and interpret spatial tissue architecture across various multiplexed images by considering both the cell  
86 types and their spatial arrangement of cells (Figure 1A). We adapt an approach originally developed for  
87 image segmentation in computer vision [22], incorporating spatial information into the flexible design of  
88 regions (image partitions, analogous to documents in language modeling). Unlike standard image pixels,  
89 the basic units of analysis in multiplexed tissue images are cells, which are not uniformly distributed  
90 due to the complexity of human tissue samples, posing a unique challenge. To address these challenges,  
91 we refined the original model used for image segmentation by using a nearest-neighbor kernel function  
92 to boost computational efficiency, as well as a unique initialization strategy for robustness. In addition,  
93 we also provide an efficient implementation of the spatial topic model in our R package *SpaTopic*.

95 *SpaTopic* offers a scalable solution for cell neighborhood and domain analysis on large-scale, multi-image  
96 datasets, efficiently handling data without requiring the extraction of cell neighborhood information for  
97 each individual cell – a process that becomes computationally demanding and inefficient with millions of  
98 cells. Unlike the rigid clustering strategies of other methods, *SpaTopic* identifies ‘topics’—tissue microen-  
99 vironment features—through a probabilistic distribution over cell types and *across* diverse tissue images  
100 using a generative model. We demonstrate our method can accurately identify and quantify interpretable  
101 and biologically meaningful topics from imaging data without human intervention. We also present mul-  
102 tiple case studies encompassing tissue images from mouse spleen, non-small cell lung cancer, healthy  
103 lung, and melanoma tissue samples. Finally, we highlight an example of a TLS-like topic and its correla-  
104 tion with outcomes from *SpaTopic* analysis across different platforms, as well as a multi-stage example  
105 showing dynamic changes in spatial tissue architecture across varying disease stages.

## 106 Results

### 107 Overview of *SpaTopic*, a Bayesian probabilistic model for highly scalable 108 and interpretable spatial topic analysis *across* multiplexed tissue images

109 *SpaTopic* is designed as a flexible spatial analysis module within the current imaging analysis workflow  
110 (Figure 1B). Its main objective is to identify biologically meaningful topics *across* multiplexed images  
111 using unsupervised learning. Here, “topics” refer to latent spatial features defined by distinct cell type  
112 compositions within tissue microenvironment neighborhoods. *SpaTopic* incorporates spatial data into a  
113 Latent Dirichlet Allocation model, assuming that each cell in an image arises from a mixture of spa-  
114 tially resolved topics, with each topic being a distribution over distinct cell types. Combining cell type  
115 information with spatial orientation, this method enables the automated and simultaneous detection of  
116 immunological patterns *across* multiple images. Subsequent analyses can further link these topics with  
117 patient data, such as treatment response and survival.

118 We adopt a Bayesian approach for inference to model the uncertainties inherent in tissue spatial patterns.  
119 *SpaTopic* requires cell types and their locations as input, with the cell types determined by the users’  
120 preferred phenotyping algorithm tailored to the specific marker panel of the dataset. The algorithm  
121 generates two key statistics for further analysis: 1) topic content, a spatially-resolved topic distribution  
122 over cell types, and 2) topic assignment for each cell within the images. After Gibbs sampling, the topic  
123 assignment of each cell is determined by the topic with the highest posterior probability. Cell types  
124 enriched in the same topic tend to be spatially correlated across images, leading to the identification of  
125 recurrent patterns of cell-cell interactions.

126 We developed an R package to efficiently implement the *SpaTopic* pipeline as outlined in Figure 1A,  
127 which details the primary steps of the pipeline (See the Methods section). Figure 1C displays a graphical  
128 representation of the spatial topic model. The key inputs for *SpaTopic* are the cell type annotations  $\mathcal{C}$  and  
129 their locations  $\mathcal{X}$  across all images. Here,  $Z_{gi}$  denotes the topic assignment, and  $D_{gi}$  indicates the region  
130 assignment of cell  $i$  in image  $g$ . Analogous to how computer vision algorithms segment images by spatially  
131 co-occurring pixel patterns with similar color, intensity or texture for object detection, *SpaTopic* identifies  
132 topics as clusters of spatially co-occurring cell types (shown in Figure 1D), potentially corresponding to  
133 biologically meaningful cellular structures (e.g., tertiary lymphoid structure). The process involves the  
134 following steps:

- 137 • Initialization: Anchor cells are chosen as regional centers via spatially stratified sampling. For each  
138 image, a KNN graph is constructed between anchor cells and all other cells: For each cell, we retrieve  
139 its top  $m$  closest anchor cells. The initial region assignments of cells are made based on proximity to  
140 region centers.
- 141 • Collapsed Gibbs sampling: for every individual cell, there are two main steps per iteration:
  - 142 – Sample topic assignment  $Z_{gi}$  conditional on its region assignment  $D_{gi}$  and cell type  $c_{gi}$ , as well as  
143 the topic distribution of the region  $D_{gi}$  and the cell type distribution of the topic  $Z_{gi}$ .

- Sample region assignment  $D_{gi}$  conditional on current topic assignment  $Z_{gi}$ , distance of the cell  $\mathbf{x}_{gi}^c$  to the region center  $\mathbf{x}_{D_{gi}}^d$ , and the topic distribution of the region  $D_{gi}$ . The spatial information is weakly incorporated with a kernel function.
- After Gibbs sampling, the output includes the posterior probabilities  $Z_{gi}$  of each cell and the per-topic cell type distribution  $\{\hat{\beta}_k\}$ . Each cell in the image is assigned to a topic with the highest posterior probability  $P(Z_{gi}|\mathcal{C}, \mathcal{X})$ .

We applied *SpaTopic* to multiple datasets from diverse imaging platforms, including spatial proteomics data from Co-detection by Indexing (CODEX), Multiplexed ImmunoFluorescence (mIF), and Imaging Mass Cytometry (IMC) platforms, as well as spatial transcriptomics data from Nanostring CosMx (Table S1). In the next few sections, we apply *SpaTopic* to analyze tissue imaging data from a variety of spatial molecular profiling platforms and benchmark analysis of *SpaTopic* with other popular algorithms for spatial domain/niche analysis, including Seurat v5 [16], Spatial-LDA [14], CytoCommunity [18], UTAG [13], and BankSY [17] (Table S2). The benchmark datasets contain between 0.1 to 1 million cells per image; making it challenging to apply methods with high computational costs. In contrast, *SpaTopic* processes these large-scale images in just a few minutes.

## ***SpaTopic* identifies global and local spatial features of human lung cancer tissue with higher precision and interpretability**

We applied our method to a single non-small cell lung cancer (NSCLC) tissue image generated using a 960-plex CosMx RNA panel on the Nanostring CosMx Spatial Molecular Imager platform, which is publicly available on the Nanostring website. We selected a Lung5-1 sample containing approximately 100,000 cells, with 38 cell types annotated using Azimuth [23] based on the human lung reference v1.0 (Figure 2A).

To illustrate the general tissue architecture, Figure 2A displays the distribution of the top 10 main cell types and the expression patterns of key genes including *KRT17*, *C1QA*, *IL7R*, *TAGLN*, *MS4A1*. These genes serve as markers for tumor cells (*KRT17*), macrophages, CD4 T cells, stroma cells, and B cells, respectively (Figure 2B). Our results demonstrate that *SpaTopic* identified seven distinct topics from the complex image (Figure 2A), with each topic representing a unique spatial niche characterized by a specific cell-type composition, as detailed in Figure 2C. For example, Topic 2 is predominantly composed of tumor cells, indicating the tumor region in the image, while other topics correspond to distinct immune-enriched stromal regions. Topic 4 represents a stromal region enriched with macrophages. Notably, Topic 3 captures tertiary lymphoid-like structures in the lung tissue, consisting of B cells, CD4 T cells, and smaller proportions of dendritic cells and CD8 T cells. This composition aligns with the current understanding of cell types in tertiary lymphoid structures, which are strong predictive biomarkers associated with a good prognosis and response to immunotherapy in non-small cell lung cancer [24].

We compared results from *SpaTopic* with Seurat v5, Spatial-LDA, CytoCommunity, BankSY, and UTAG. BankSY and UTAG directly use cell-level gene expression as input, whereas the other four methods, including *SpaTopic*, rely on cell-type annotations. All methods can detect the global structure of the image and classify tumor and stromal regions. However, BankSY and UTAG appear to miss the lymphoid structure, likely because they do not use the detailed information provided by cell-type annotations. Reference-based cell type annotation typically offers more detailed information and can be more robust for noisy data when matched with a single-cell reference [25, 26]. *SpaTopic* distinctly identified the lymphoid structure as Topic 3, comprising a mix of CD4 T cells and B cells (Figure 2F). Additionally, when we focused on two local tumor tissue regions (Figures 2D and 2E), *SpaTopic* identified the tumor region with higher precision (Topic 2), more consistently matching the expression pattern of *KRT17*, a lung cancer marker gene. *SpaTopic* and UTAG are the only two methods showing the consistency (between tumor domain and *KRT17* expression) higher than 0.8 across the entire image (Figure 2G), which aligns with the visual measure in Figure 2D and 2E.



## 194 *SpaTopic* identifies tertiary lymphoid structures from whole-slide melanoma 195 tissue imaging

196 We applied *SpaTopic* to a whole-slide melanoma tissue image obtained from our internal multiplexed  
197 immunofluorescent (mIF) imaging platform, which uses a 12-plex marker panel [27]. This analysis covered  
198 a whole-slide soft tissue image containing 0.4 million cells, annotated into seven major cell types (CD4  
199 T cells, Tumor/Epithelial, B cells, CD8 T cells, Macrophages, Regulatory T (Treg) cells, and Others).  
200 The categorization was based on the expression of six lineage markers: CK/SOX10, CD3, CD8, CD20,  
201 CD68, and Foxp3. Cells were annotated as ‘Other’ if they showed negative expression for all six markers.

202  
203 Despite using fewer markers compared to the Nanostring CosMx platform, *SpaTopic* identified five  
204 distinct topics (Figure 3A): Topic 1 (tumor), Topic 2 (CD4 immune zone), Topic 3 (stroma), Topic  
205 4 (immune-enriched tumor-stroma boundary), and Topic 5 (tertiary lymphoid structures). The tissue  
206 structures revealed by these topics visually correspond to the histological pattern seen in the co-  
207 registered H&E image (Figure 3B) and the merged raw mIF images (Figure 3C) with three key markers:  
208 CD3 (T cells), CD20 (B cells), and PANCK/SOX10 (tumor cells). Figure 3D demonstrates that topic 5  
209 (tertiary lymphoid structures) mainly consists of B cells, CD4 T cells, a few CD8 T cells, and Treg cells,  
210 consistent with the TLS-like pattern identified in the Nanostring dataset discussed earlier. Due to the  
211 lack of a dendritic cell marker in the mIF dataset, dendritic cells could not be identified and included in  
212 topic 5. This analysis demonstrates that *SpaTopic* can consistently detect the same biologically relevant  
213 patterns across various tumor tissues and imaging platforms, which may be clinically significant, as  
214 tertiary lymphoid structure have been recognized as a promising biomarker for cancer immunotherapy.

## 216 *SpaTopic* recovers spatial domain from cell type spatial organization in 217 healthy lung tissue

218 We further demonstrate that *SpaTopic* can effectively distill signals from noisy cell type annotations  
219 and identify clear tissue architecture based solely on the spatial arrangement of cells. To illustrate this,  
220 we applied *SpaTopic* to the IMC dataset from the UTAG paper [13], which includes 26 small regions of  
221 interest (ROIs) images from healthy lung tissue. For comparison, we used the UTAG result provided in  
222 the paper [13] without rerunning UTAG.

223  
224 Our analysis shows that *SpaTopic* can recover tissue architectures directly from the spatial distribution  
225 of cell type annotations, yielding results consistent with manual annotations (Figure 4A). *SpaTopic* per-  
226 forms comparably to UTAG using only cell type annotations (Figure 4B), as indicated by the adjusted  
227 Rand index, which shows similar performance levels. Additionally, Figure 4C illustrates the topic content  
228 and cell type composition for each topic identified by *SpaTopic*. This demonstrates *SpaTopic*’s capability  
229 to perform domain analysis without discarding existing cell type annotations, offering valuable flexibility  
230 for datasets with cell-type annotations or for incorporating any existing cell-type annotation method.  
231 Unlike UTAG, which learns spatial tissue architecture directly from cell features due to noisy cell type  
232 annotations, we demonstrate that *SpaTopic* can effectively identify tissue architecture from these annota-  
233 tions. Thus *SpaTopic* is a robust alternative that leverages existing data without the need for additional  
234 cell-level features.

## 235 *SpaTopic* identifies disease-specific topics and tracks topic evolution in 236 mouse spleen over disease progression

237 We also applied *SpaTopic* to a CODEX mouse spleen dataset [2] to demonstrate its proficiency in iden-  
238 tifying spatial topics across multiple images. This dataset includes nine images: three control normal  
239 BALBc spleens (BALBc 1-3) and six MRL spleens (samples 4-9) at varying disease stages—early (MRL  
240 4-6), intermediate (MRL 7-8), and late (MRL 9) (Figure 5A). Using a 30-plex protein marker panel, the  
241 study identified 27 major splenic-resident cell types across the nine tissue images. We use the cell type  
242 annotation in the original paper [2].

244 *SpaTopic* identified six topics from approximately 0.7 million cells across the nine images, highlighting  
245 the dramatic changes in spatial tissue structures associated with disease progression from normal spleen  
246 to spleen tissue at different disease stages (Figure 5A). Figure 5B and 5C highlight per-topic cell type  
247 compositions, aiding in labeling each topic. The normal spleen tissue samples predominantly comprised  
248 three topics: Topic 1 (red pulp), Topic 2 (periarteriolar lymphoid sheath, PALS), and Topic 3 (B-follicle).  
249 Figures S2 and S3 show the cell type distribution and domain annotations from the original paper,  
250 demonstrating *SpaTopic*'s ability to capture the main structures consistent with these annotations, as  
251 compared to other methods (Figures S2 and S4). With an increasing number of topics, *SpaTopic* also  
252 successfully delineated the marginal zone from the B-follicle (Figure S2).

253  
254 Topics identified by *SpaTopic* were comparable across normal and diseased spleens, allowing us to  
255 identify condition-specific topics and quantify changes in topic proportions as the disease progressed.  
256 In contrast to normal spleens, MRL spleens showed a decrease in B cells and F4/80(+) macrophages  
257 but an increase in granulocytes and erythroblasts within the red pulp region, indicating inflammation  
258 or systemic infection in the spleen tissue. This shift was marked by the predominance of Topic 6 in  
259 MRL spleens, superseding Topic 1 (red pulp). Topic 4 emerged in the mouse spleen tissue affected by  
260 autoimmune disease, characterized by a high abundance of CD106+ stroma cells, indicative of leukocyte  
261 recruitment to inflamed areas. This topic also shows a high concentration of immune cells, including  
262 CD4 and CD8 T cells. Unique to MRL spleens, Topic 5 is characterized by an enrichment of B220+  
263 double negative (DN) T cells and conventional CD4 T cells, predominating in tissues during advanced  
264 stages of the disease, indicating a shift in the immune cell landscape. These dynamics indicate immune  
265 surveillance or dysregulation in the spleen tissue with MRL/lpr progression [2].

266  
267 Furthermore, *SpaTopic*'s capability to identify topics based on the spatial proximity of cell types  
268 suggests that cell types grouped within the same topic are likely close to each other and prone to  
269 interaction. Figure 5D illustrates the changes in topic proportions throughout the course of the disease.  
270 The distinct contributions of cell types to each topic are highlighted in Figure 5E, selected based on  
271 their specific composition and evaluated based on the lift and FREX metrics [28, 29] (Figures S5). Cell  
272 types are clustered into topics that exhibit similar dynamics across different slides.

## 274 ***SpaTopic* is highly scalable on large-scale modern images**

275 To benchmark the scalability of *SpaTopic* as the number of cells in images increases, we conducted tests  
276 using simulated datasets of varying scales. Figure 6A demonstrates that our method is scalable with an  
277 increasing number of cells within a single image, compared to Seurat v5. Figure 6B further confirms  
278 the high scalability of *SpaTopic* when evaluating the user time of all methods on real datasets. For the  
279 Nanostring CosMx NSCLC image with around 0.1 million cells, *SpaTopic* runs within 1 minute on a  
280 standard MacBook Air. *SpaTopic* is in the same tier as Seurat v5. BankSY and UTAG are in the second  
281 tier since they use similar strategies. CytoCommunity, limited by GPU support, was run with reduced  
282 epochs and only on CPU for the NSCLC dataset, which compromised its performance and underscored  
283 its impracticality for labs without extensive computing resources.

## 285 **Discussion**

286 In summary, we introduced *SpaTopic*, a spatial topic model designed to identify and quantify biologically  
287 relevant topics across multiple multiplexed tissue images. This represents a novel approach to apply-  
288 ing language modeling techniques to decipher the tissue microenvironment from tissue imaging data.  
289 *SpaTopic* stands out as one of the few unsupervised learning methods capable of discerning clinically  
290 relevant spatial patterns [13, 17, 19]. Unlike other methods that rely on hard clustering strategies for  
291 analyzing samples, *SpaTopic* is a probabilistic model-based approach using Bayesian inference methods  
292 to identify complex tissue architectures. The model generates two key outputs: The first of these, the  
293 *topic content* maps the cell type composition in spatial niches, allowing direct interpretation of the  
294 corresponding topic (e.g., TLS); The second output, *topic assignment* for each single cell allows the  
295 quantification of each topic in individual tissue samples for subsequent association analysis with patient

296 outcome. Application to multiple datasets along with benchmark analysis show that *SpaTopic* achieves  
297 higher precision in defining global and local spatial niches and higher sensitivity at capturing complex  
298 structures such as TLS. Notably, our method is highly scalable to large-scale imaging data with efficient  
299 runtime, handling millions of cells on a standard laptop.

300  
301 *SpaTopic* is designed as a flexible spatial analysis module within the current imaging analysis workflow.  
302 A standard image analysis pipeline includes cell segmentation, data normalization/batch correction,  
303 cell phenotyping/clustering, and the analysis of cell type content and spatial relationships. Downstream  
304 statistical analysis typically starts with cell-level metadata derived from image analysis. Due to varied  
305 marker panels and molecular imaging platforms, a one-size-fits-all solution for cell phenotyping across  
306 diverse platforms seems unlikely. In practice, we find that reference-based cell annotation works best on  
307 single-cell imaging data, rather than unsupervised clustering. *SpaTopic* does not specify any upstream  
308 method, and thus can be seamlessly integrated with other cell phenotyping modules tailored for datasets  
309 from different platforms. This design offers users adaptability, accommodating datasets from different  
310 panel designs.

311  
312 In our proposed analysis pipeline for imaging data, we separate cell phenotyping from cell neighbor-  
313 hood/domain analysis for image-based spatial data, with *SpaTopic* directly taking cell types as input.  
314 This key difference sets *SpaTopic* apart from UTAG and BankSY, which use protein/gene expression  
315 as input for niche/domain analysis. UTAG performs dimension reduction before message passing,  
316 while BankSY engineers new spatial features for each cell before dimension reduction. We propose  
317 that treating cell phenotyping and neighborhood/domain analysis as distinct steps is a better analysis  
318 strategy for datasets generated by image-based technology with selected marker panels. Using cell type  
319 annotations as input for cell neighborhood analysis enhances the interpretability of different tissue  
320 microenvironments and undoubtedly increases the computational efficiency when analyzing large-scale  
321 images. The performance of *SpaTopic* may rely on the accuracy of cell phenotyping. A better strategy  
322 for cell phenotyping is to annotate cells directly from cell images instead of using summary statistics,  
323 such as mean marker expression or gene count data. As part of the analysis pipeline, we are developing  
324 an image-based deep learning method for cell phenotyping, incorporating subcellular information, as  
325 well as domain knowledge [30].

326  
327 For multi-sample analysis, addressing the batch effect is a key challenge. Our proposed analysis pipeline  
328 seeks to mitigate the batch effect during cell phenotyping using a reference-based cell phenotyping  
329 method. For spatial transcriptomics data, a supervised classification method with a reliable single-cell  
330 reference can mitigate batch effects and inherent noises in the imaging data. Batch effect is more critical  
331 for algorithms that directly consider the gene expression data as input. When analyzing the mouse  
332 spleen dataset, we used Combat [31] for batch correction across multiple images before applying UTAG  
333 and BankSY. However, Combat appears to over-correct for batch effects (Figure S4), thus failing to dis-  
334 tinguish between normal and diseased red pulp tissue. This might stem from the substantial differences  
335 between normal and diseased tissues.

336  
337 Modern datasets from platforms like 10x Xenium and Nanostring CosMx require scalable computational  
338 methods to handle their size and complexity. Existing spatial domain analysis methods, originally  
339 designed for 10x Visium spatial transcriptomics data and optimized for datasets with thousands of cells  
340 or spots per slide, find it challenging to handle these more advanced, datasets with millions of cells per  
341 image. *SpaTopic* meets this need by efficiently managing neighborhood calculations and constructing  
342 the KNN graph only among  $m$  anchor cells instead of all  $n$  cells in the image. This reduces the time  
343 complexity from  $O(n \log n)$  to  $O(m \log m)$ , where  $m \ll n$ . Additionally, *SpaTopic* maintains linear time  
344 complexity relative to the number of cells and iterations with collapsed Gibbs sampling and uses an  
345 approximate fast approach for constructing the KNN graph. These optimizations ensure *SpaTopic*'s  
346 computational efficiency, making it accessible on standard laptops and practical for analyzing large-scale  
347 imaging data from platforms like the 10x Xenium and Nanostring CosMx.

348  
349 Moreover, advances in technology now enable the quantification of immune cell spatial diversity and the  
350 characterizing of tumor microenvironments in 3D tissues [32]. While *SpaTopic* can be adapted to infer

351 immunological topics from 3D tissue, a refined strategy is needed to select anchor cells in the 3D spaces, as  
352 the spatial information obtained by *SpaTopic* primarily stems from the relationships between the anchor  
353 cells and other cells. Incorporating a hierarchical Dirichlet prior on topic distributions across regions  
354 would allow regions within the same image to share priors while differing across images. Furthermore,  
355 optimizing the initialization strategy is needed when applying *SpaTopic* to extremely large datasets with  
356 hundreds of images. These improvements would broaden the applicability of *SpaTopic*.

## 357 Methods

### 358 *SpaTopic*

#### 359 Notations

360 We assume there are total  $V$  cell types that contribute to  $K$  different tissue microenvironments (topics)  
361 across  $G$  multiplexed images. Let  $c_{gi}$  be the  $i$ th cell at the location  $\mathbf{x}_{gi}^c = (x_{gi1}^c, x_{gi2}^c)$ ,  $g = 1, 2, \dots, G$ ,  $i =$   
362  $1, 2, \dots, n_g$ , on the  $g$ th image with total  $n_g$  cells. Let  $c_{gi} = v$  if the cell has been classified to the  $v$ th cell  
363 type. Let  $\mathcal{C} = \{c_{gi}\}_{i=1,2,\dots,n_g}^{g=1,2,\dots,G}$  and  $\mathcal{X} = \{\mathbf{x}_{gi}\}_{i=1,2,\dots,n_g}^{g=1,2,\dots,G}$  denote all observed cell types and cell locations  
364 across all  $G$  images.

#### 366 Model

367 In a conventional LDA model, each image is treated as an individual document, employing a bag-of-  
368 words approach without accounting for spatial information. This approach is similar to our prior work  
369 on longitudinal flow cytometry data analysis [28]. Here, in order to incorporate spatial information  
370 within images, we introduce a spatial topic model, *SpaTopic*, integrating spatial data into the founda-  
371 tional LDA framework. This spatial topic framework was first proposed for image segmentation [22],  
372 instead of viewing each image as a singular document, we treat each image consisting of densely placed  
373 overlapping regions (documents). Unlike the conventional LDA model where relationships between  
374 documents and words are known and fixed, the word-document relationship here is unknown: each cell  
375 (word) is flexible to be assigned to all possible regions (documents). This flexible region (document)  
376 design allows us to identify spatial structure with irregular shape.

377 For *SpaTopic*, we introduce a new hidden variable  $D_{gi}$  to denote cell region (document) assignment.  
378 Thus, each cell is associated with two hidden variables: the latent topic assignment  $Z_{gi} \in \{1, 2, \dots, K\}$   
379 and the latent region assignment  $D_{gi} \in \{1, 2, \dots, M\}$ ,  $M = \sum_g M_g$ , where  $M_g$  denote the number of  
380 regions on the image  $g$ . During the initialization, we pre-selected anchor cells as region centers. Let  
381  $\mathcal{X}^d = \{\mathbf{x}_d^d\}_{d=1,2,\dots,M}$  be the set of all  $M$  region centers across all images. Let  $\theta_d$  be the proportion of  
382 region  $d$  over  $K$  topics and  $\beta_k$  be the proportion of topic  $k$  over  $V$  cell types. Hyperparameters  $\psi$  and  
383  $\alpha$  specify the nature of the Dirichlet priors of  $\{\beta_k\}$  and  $\{\theta_d\}$ , respectively.

384  
385 Then we are ready to describe our generative model:

- 387 • For each topic  $k$ , sample  $\beta_k$  (topic weights over  $V$  cell types) from a Dirichlet prior  $\beta_k \sim \text{Dir}(\psi)$ .
- 388 • For each image region  $d$  (centered at  $\mathbf{x}_d^d$ ), sample topic proportion  $\theta_d \sim \text{Dir}(\alpha)$
- 389 • For each cell, the  $i$ th cell in the image  $g$ :
  - 390 – Sample its region assignment  $D_{gi}$  from a uniform prior over possible documents (regions) in the
  - 391 image  $g$ .
  - Sample the location  $\mathbf{x}_{gi}^c$  conditional on its region assignment  $D_{gi}$  with a kernel function based on
  - the distance between the cell location  $\mathbf{x}_{gi}^c$  and the region center  $\mathbf{x}_d^d$ .

$$\mathbf{x}_{gi}^c | D_{gi} = d \propto K(\mathbf{x}_{gi}^c, \mathbf{x}_d^d).$$

- 392 – Sample topic assignment  $Z_{gi} | D_{gi} = d \sim \text{Multi}(\theta_d, 1)$ .
- 393 – Sample cell type  $c_{gi} | Z_{gi} = k \sim \text{Multi}(\beta_k, 1)$ .

394 Hyperparameters  $\alpha$  and  $\psi$  should be chosen based on the belief on  $\{\theta_d\}$  and  $\{\beta_k\}$  in a Bayesian perspec-  
395 tive. In our application, both  $\alpha$  and  $\psi$  are set very small by default (default:  $\alpha_k = .01, \forall k$ ;  $\psi_v = .05, \forall v$ )  
396 to encourage the sparsity in region-topic distributions  $\{\theta_d\}$  and topic-celltype distributions  $\{\beta_k\}$ .

#### 397 Nearest-neighbor Exponential Kernel

398 The flexible relationships between regions and cells in *SpaTopic* allow each cell to be assigned to any  
399 one of its proximate regions. We employ a nearest-neighbor Gaussian kernel to capture the spatial  
400 correlation between cells and their respective regions, as previously used in the nearest-neighbor Gaussian



401 process [33]. For computational efficiency, especially with large-scale images, we restrict our consideration  
 402 to the top nearest-neighbor regions for each cell. Let  $\mathcal{N}(\mathbf{x}_{gi}) \subset \mathcal{X}^d$  be the collection of  $m$  closed region  
 403 centers to the cell  $\mathbf{x}_{gi}$  (default:  $m = 5$ ). In practice, the commonly used squared exponential Gaussian  
 404 kernel function decays too rapidly. This rapid decay often results in cells predominantly being linked to  
 405 their closest region, irrespective of their cell types. Let  $\sigma$  be the lengthscale that controls the strength of  
 406 decay of correlation with distance in the kernel function. Thus, drawing inspiration from [34], instead of  
 407 the squared exponential kernel, we used the following exponential kernel,

$$K(\mathbf{x}_{gi}^c, \mathbf{x}_d^d) \propto \mathbb{1}\{\mathbf{x}_d^d \in \mathcal{N}(\mathbf{x}_{gi}^c)\} \exp\{-\|\mathbf{x}_{gi}^c - \mathbf{x}_d^d\|_2/\sigma\}, \quad (1)$$

408 where  $\|\mathbf{x}_{gi}^c - \mathbf{x}_d^d\|_2$  represents the Euclidean distance between the cell location  $\mathbf{x}_{gi}^c$  and the region  
 409 center  $\mathbf{x}_d^d$ . We fix  $\sigma$  for computational efficiency, but it can also be sampled during the Gibbs sampling.  
 410 Increasing  $\sigma$  would reduce the strength of the spatial correlation, resulting in a diminished spatial effect  
 411 when assigning cells to regions.

### 413 Collapsed Gibbs Sampling

414 We use collapsed Gibbs Sampling for model inference. The collapsed Gibbs Sampling algorithm was first  
 415 introduced as the Bayesian approach of Latent Dirichlet Allocation [35]. This method's comprehensive  
 416 derivation and implementation can be found in the paper [36]. Similar to [22], we further adapted and  
 417 extended the algorithm for our proposed spatial topic model. It's noteworthy that during the collapsed  
 418 Gibbs sampling process, the parameters  $\beta_k$  and  $\theta_d$  are integrated out and are not explicitly sampled.  
 419 Instead, our focus is on the two hidden variables associated with each cell: the topic assignment  $Z_{gi}$  and  
 420 the region (or document) assignment  $D_{gi}$ . These variables undergo iterative sampling using the collapsed  
 421 Gibbs Sampler:

- 422 1. Sample topic assignment  $Z_{gi}$  conditional on region assignment  $D_{gi}$  with [35]

$$P(Z_{gi} = k \mid D_{gi} = d, c_{gi} = v, \mathcal{D}_{-gi}, \mathcal{Z}_{-gi}, \mathcal{C}_{-gi}, \boldsymbol{\psi}, \boldsymbol{\alpha}) \propto \frac{n_{k,-gi}^{(v)} + \psi_v}{\sum_{t=1}^V n_{k,-gi}^{(t)} + \psi_t} \frac{n_{d,-gi}^{(k)} + \alpha_k}{\sum_{k'=1}^K n_{d,-gi}^{(k')} + \alpha_{k'}} \quad (2)$$

423 where  $n_{k,-gi}^{(v)}$  refers the number of times that cell type  $v$  has been observed with topic  $k$  and  $n_{d,-gi}^{(k)}$   
 424 refers the number of times that topic  $k$  has been observed in region  $d$ , both excluding the current cell  
 425  $gi$ , the  $i$ th cell on the  $g$ th image. The first ratio expresses the probability of cell type  $v$  under topic  
 426  $k$ , and the second ratio expresses the probability of topic  $k$  in region  $d$ .  $\mathcal{D}_{-gi}$ ,  $\mathcal{Z}_{-gi}$ , and  $\mathcal{C}_{-gi}$  denote  
 427 collections of  $\mathcal{D}$ ,  $\mathcal{Z}$ , and  $\mathcal{C}$  excluding cell  $c_{gi}$ .

- 428 2. Sample  $D_{gi}$  conditional on  $Z_{gi}$  with

$$\begin{aligned} & P(D_{gi} = d \mid Z_{gi} = k, \mathcal{D}_{-gi}, \mathcal{Z}_{-gi}, \mathbf{x}_{gi}^c, \mathbf{x}_d^d, \boldsymbol{\alpha}, \sigma) \\ & \propto P(Z_{gi} = k \mid \mathcal{Z}_{-gi}, D_{gi} = d, \mathcal{D}_{-gi}, \boldsymbol{\alpha}) P(\mathbf{x}_{gi}^c \mid D_{gi} = d, \mathbf{x}_d^d, \sigma) P(D_{gi} = d) \end{aligned}$$

429 According to [36],  $P(Z_{gi} = k \mid \mathcal{Z}_{-gi}, D_{gi} = d, \mathcal{D}_{-gi}, \boldsymbol{\alpha})$  can be obtained by integrating out  $\theta_d$ , that

$$P(Z_{gi} = k \mid \mathcal{Z}_{-gi}, D_{gi} = d, \mathcal{D}_{-gi}, \boldsymbol{\alpha}) = \frac{n_{d,-gi}^{(k)} + \alpha_k}{\sum_{k'=1}^K n_{d,-gi}^{(k')} + \alpha_{k'}}.$$

430 We can further omit  $P(D_{gi} = d)$  due to uniform prior. Thus  $D_{gi}$  can be sampled based on the following  
 431 conditional distribution:

$$P(D_{gi} = d \mid Z_{gi} = k, \mathcal{D}_{-gi}, \mathcal{Z}_{-gi}, \mathbf{x}_{gi}^c, \mathbf{x}_d^d, \boldsymbol{\alpha}, \sigma) \propto K(\mathbf{x}_{gi}^c, \mathbf{x}_d^d) \frac{n_{d,-gi}^{(k)} + \alpha_k}{\sum_{k'=1}^K n_{d,-gi}^{(k')} + \alpha_{k'}} \quad (3)$$

---

**Algorithm 1** Collapsed Gibbs Sampling for *SpaTopic*

---

1. Identify  $M$  anchor cells (located at  $\{\mathbf{x}_d^d\}_{d=1,2,\dots,M}$ ) as the region centers across images.
  2. For each image, pre-compute a  $k$ -nearest-neighbor graph between all cells and the selected region centers.
  3. Initialize topic assignment  $Z_{gi}$  and region assignment  $D_{gi}$  for each cell. Compute region-topic counts  $n_d^{(k)}$  and topic-celltype counts  $n_k^{(v)}$ .
  4. Gibbs sampling over burn-in and sampling period. For each cell, do
    - (a) Update counts  $n_d^{(k)}$  and  $n_k^{(v)}$  excluding the current  $Z_{gi}$  and  $D_{gi}$ .
    - (b) Sample topic assignment  $Z_{gi}$  conditional on region assignment  $D_{gi}$  based on equation (2).
    - (c) Sample region assignment  $D_{gi}$  conditional on topic assignment  $Z_{gi}$  based on equation (3).
    - (d) Update counts  $n_d^{(k)}$  and  $n_k^{(v)}$  with the updated  $Z_{gi}$  and  $D_{gi}$ .
  5. Check convergence. If converged during burn-in and  $L$  posterior samples drawn, output posterior samples and parameters estimated based on equation (4) and (5). If not, increase the number of iterations for burn-in.
- 

## Initialization

During the initialization, we employ a spatially stratified sampling approach to randomly select anchor cells from each image, which will serve as region centers. The number of anchor cells selected from each image is determined by a predetermined region radius  $r$  (default:  $r = 400$ ), as well as the image size. The radius should be set with the consideration of the image resolution and complexity of the images, and an adequate number of cells are expected within each region since it is crucial for estimating topic distribution  $\theta_d$  precisely. In practice, for whole-slide imaging, we expect at least 100 cells per region on average. For each individual image, an  $m$ -nearest-neighbor graph will be constructed between all cells and the chosen anchor cells. For computational efficiency, distances between each cell and its top  $m$ -nearest anchor cells will be pre-computed before Gibbs sampling.

The performance of *SpaTopic* depends on anchor cells selected in the initialization, especially on images with highly complex spatial structures. Thus, we take a warm start approach rather than starting Gibbs sampling from a random initialization. This involves running multiple Gibbs sampling initializations (default: `ninit = 10`), each having a unique set of anchor cells. After a few iterations (default: `niter_init = 100`), only the one with the highest log-likelihood is retained and continued.

## Implementation

We implemented *SpaTopic* in Rcpp and made it an R package *SpaTopic* (officially available on CRAN after Jan 17, 2024). The complete algorithm is shown in Algorithm 1. For the Gibbs sampling, we have set the default parameters as follows: `iter = 200`, `burnin = 1000`, `thin = 20` (200 Gibbs sampling draws are made with the first 1000 iterations discarded and then every 20th iteration kept). We can infer topic distributions across all images using the posterior samples drawn from the Gibbs sampling. For each of these posterior samples, the predictive distributions of parameters  $\{\beta_k\}$  and  $\{\theta_d\}$  are obtained as follows:

$$\hat{\beta}_{kv} = \frac{n_k^{(v)} + \psi_v}{\sum_{t=1}^V n_k^{(t)} + \psi_t}, \quad (4)$$

$$\hat{\theta}_{dk} = \frac{n_d^{(k)} + \alpha_k}{\sum_{k'=1}^K n_d^{(k')} + \alpha_{k'}}. \quad (5)$$

Moreover, we also keep the posterior distribution of  $Z_{gi}$  from all posterior samples for each individual cell. Notably,  $D_{gi}$  has been marginalized during this process and each cell in the end is assigned to the topic with the highest posterior probability. Thus we are also able to visualize the spatial distribution of cell topics in the images.

## 461 Model Selection

462 The likelihood of the topic model is intractable to compute in general, but we can approximate the model  
463 log-likelihood in terms of model parameters  $\{\beta_k\}$  and  $\{\theta_d\}$  [37]. With the law of total probabilities, we  
464 take into account uncertainties both in cells' region and topic assignment, then the log-likelihood of the  
465 spatial topic model can be presented as

$$466 \quad \ell(\mathcal{C}, \mathcal{X}) = \sum_g \sum_{i=1}^{N_g} \log \left[ \sum_{k=1}^K \sum_{d=1}^M \sum_{v=1}^V \mathbb{I}(c_{gi} = v) \theta_{dk} \beta_{kv} \eta_{gi}^d \right], \quad (6)$$

466 where  $\eta_{gi}^d = P(\mathbf{x}_{gi}^c \mid D_{gi} = d, \mathbf{x}_d^d) P(D_{gi} = d) \propto K(\mathbf{x}_{gi}^c, \mathbf{x}_d^d)$ .

467  
468 We use the Deviance Information Criterion (DIC) [38] to select the number of topics, a generalization of  
469 the Akaike Information Criterion (AIC) in Bayesian model selection:

$$470 \quad DIC = p_D + \overline{D(\mathcal{C}, \mathcal{X})}, \quad (7)$$

470 where the Deviance is defined as  $D(\mathcal{C}, \mathcal{X}) = -2\ell(\mathcal{C}, \mathcal{X})$  and  $p_D = \frac{1}{2} \overline{Var(D(\mathcal{C}, \mathcal{X}))}$ .

471  
472 DIC requires calculating the log-likelihood for every posterior sample, which is time-consuming. To  
473 determine the optimum number of topics, we run *SpaTopic* with a varied number of topics (2-9 in  
474 practice) and collect a few posterior samples (such as the first 20 posterior samples) after convergence  
475 (with `trace=1`). The number of topics was selected based on DIC (7). Otherwise, we only output the  
476 deviance and the log-likelihood of the final posterior sample (default: `trace=0`).

## 478 Comparing to other methods

479 We compared the performance of *SpaTopic* with five other niche analysis methods: spatial-LDA, Seurat-  
480 v5, UTAG, CytoCommunity, and BankSY. For BankSY and UTAG, we used protein or gene expression  
481 data and cell spatial coordinates as inputs, while the other methods used existing cell-type annotations  
482 and cell spatial coordinates. We followed the pre-processing procedures and parameters described in the  
483 original papers and tutorials for each method, with some hyperparameters slightly adjusted for computa-  
484 tional efficiency on large datasets or when clear guidelines for tuning parameters were available. Details  
485 of these adjustments and the rationale for not using the default settings are described in this section.

486  
487 All methods were initially run using R Studio (for R-based methods) or Jupyter Lab (for Python-based  
488 methods) on a standard MacBook Air (M2, 2022). If a method could not be run on a standard Mac  
489 due to memory constraints, we used our high-performance computing server with a single-core CPU  
490 and 200GB of assigned memory. For the Nanostring CosMx NSCLC dataset, both CytoCommunity and  
491 UTAG were run on the server due to high memory usage. Additionally, for the CODEX mouse spleen  
492 dataset, UTAG can be run on the Mac only without the default parallel mode due to memory constraints.

493  
494 ***SpaTopic* (v1.1.0).** We ran *SpaTopic* with `region_radius = 400, 150, 300` for the NSCLC, the mouse  
495 spleen, and the melanoma datasets, respectively, allowing around 100 cells per region on average dur-  
496 ing initialization, which is necessary for accurately estimating the topic-region distribution. We chose  
497 length-scale `sigma = 20` for the mouse spleen dataset and used the default parameters for the NSCLC  
498 dataset. Posterior samples were collected after the convergence of the Gibbs sampling chain, with a  
499 burn-in period of 2000 iterations for the NSCLC dataset and 1500 iterations for the mouse spleen  
500 dataset. For the Melanoma dataset, *SpaTopic* was run with a burn-in period of 2000 iterations. For the  
501 healthy lung dataset with 26 small ROIs, *SpaTopic* was run with `sigma = 5` and `radius = 60` to identify  
502 the complex local structures. In addition, we increase the number of initializations to 200 times to  
503 increase the robustness of identifying consensus patterns across ROIs while increasing the running time.

504  
505 **Seurat-v5 (v5.0.2).** We used the default niche analysis in Seurat v5, specifically the `BuildNicheAssay()`  
506 function in the Seurat R package. Seurat v5 employs k-means clustering to group cell neighborhood

507 features, which are derived from the shared-nearest-neighbor graph (default neighbors.k = 30), a variant  
508 of the k-nearest-neighbor graph, as part of its image-based spatial data analysis pipeline. We ran Build-  
509 NicheAssay() with all default parameters except for the NSCLC datasets, for which we set neighbors.k  
510 = 100. We found that increasing neighbors.k from 10 to 100 (testing neighbors.k = 10, 30, 50, 100)  
511 significantly improved the algorithm’s performance on this dataset, with results presented in Figure S1.

512

513 **Spatial-LDA (v0.1.3).** When working on mouse spleen datasets, we used the same parameters as the  
514 authors used in the original methodology paper, though we now use neighborhoods of all cells as the  
515 input, not only B cells. For the NSCLC datasets, we also use neighborhoods of all cells as the input  
516 but set radius = 400 to extract neighborhood cell type compositions. To reduce the computational  
517 complexity for both datasets, we set the threshold = 0.01 for ADMM Primal-Dual optimizer. Finally,  
518 we output the topic weights for every cell and assign every cell to a topic with the maximal weight.

519

520 **CytoCommunity (Github version obtained on 2024 February).** CytoCommunity (unsupervised  
521 version) was run on a CPU with 200GB of assigned memory and evaluated only on the NSCLC dataset  
522 due to its demand for large-memory GPU resources and the unsupervised version’s inability to learn  
523 Tissue Cell Neighborhoods (TCNs) across multiple images (TCNs learned from individual images are  
524 not comparable). We set KNN-K = 300 for 0.1M cells, as suggested in the original paper. For large  
525 image data, the second step of CytoCommunity is time-consuming when trained on a CPU. Therefore,  
526 we greatly reduced num\_RUN to 10 and Num\_Epoch to 100 per run while ensuring the final loss was  
527 less than -0.2 for each run. Other parameters were set to their defaults.

528

529 **UTAG (v0.1.1).** UTAG was primarily developed for protein expression data with limited marker  
530 channels. For the Nanostring CosMx NSCLC datasets with 960 genes, we used typical pre-processing  
531 steps suggested by Scanpy (v1.9.8) for analyzing scRNA-seq datasets. These steps included filtering  
532 low-prevalence genes, log transformation, and retaining only highly variable genes. We then performed  
533 z-score normalization, truncated at 10 standard deviations, followed by PCA. Only the top 50 principal  
534 components were used as input for UTAG. UTAG was run under multiple clustering resolutions [0.05,  
535 0.1, 0.3, 0.5] and mix\_dist = 60, with an image resolution of 0.18 microns per pixel, since the authors  
536 suggested setting mix\_dist between 10 and 20 microns in the user manual. For the CODEX mouse spleen  
537 dataset (with intensity values already transformed), we performed z-score normalization truncated at  
538 10 standard deviations, followed by Combat batch correction [31] and a second z-score normalization  
539 truncated at 10 standard deviations, a similar procedure as introduced in the UTAG paper for prepro-  
540 cessing IMC data [13]. We also set mix\_dist = 60, with an imaging resolution of 0.188 microns per pixel.

541

542 **BankSY (v0.99.9).** In contrast to UTAG, BankSY is specifically designed to analyze spatial tran-  
543 scriptomics datasets. We ran BankSY with lambda = 0.8 to identify spatial domains, as recommended,  
544 with other parameters set to default, as described in the GitHub tutorial. For the NSCLC dataset, we  
545 followed the same pre-processing procedures outlined in the domain analysis tutorial, using k\_geom =  
546 30, npcs = 50, and clustering resolutions of 0.1, 0.2, 0.3, and 0.5. For the mouse spleen datasets, we  
547 used the same input as UTAG, after batch correction and normalization. We followed the tutorial for  
548 multi-sample analysis, running the results under npcs = 30 since the dataset has only 30 markers.

549

## 550 Data Preprocessing

551 **Nanostring CosMx Human NSCLC.** The Nanostring CoxMx NSCLC dataset is  
552 available on the Nanostring Website (<https://nanosttring.com/products/cosmx-spatial-molecularimager/ffpe-dataset/nsclc-ffpe-dataset/>). For our analysis, we selected Lung5-1  
553 sample and annotated about 0.1M cells into 38 cell types using Azimuth [23] with  
554 a human lung reference v1.0 (<https://azimuth.hubmapconsortium.org/references/>). We  
555 used the same cell annotations from the Seurat image analysis pipeline tutorial  
556 ([https://satijalab.org/seurat/articles/seurat5\\_spatial\\_vignette\\_2.html](https://satijalab.org/seurat/articles/seurat5_spatial_vignette_2.html)). Since healthy lung tissue was  
557 used as the reference, the ‘basal’ cells were re-labeled as tumor cells since they are the most closed  
558 cell type. We checked that the tumor locations indicated by the reference-based cell annotations are  
559 generally consistent with the tumor region labeled by the Nanostring company.

560

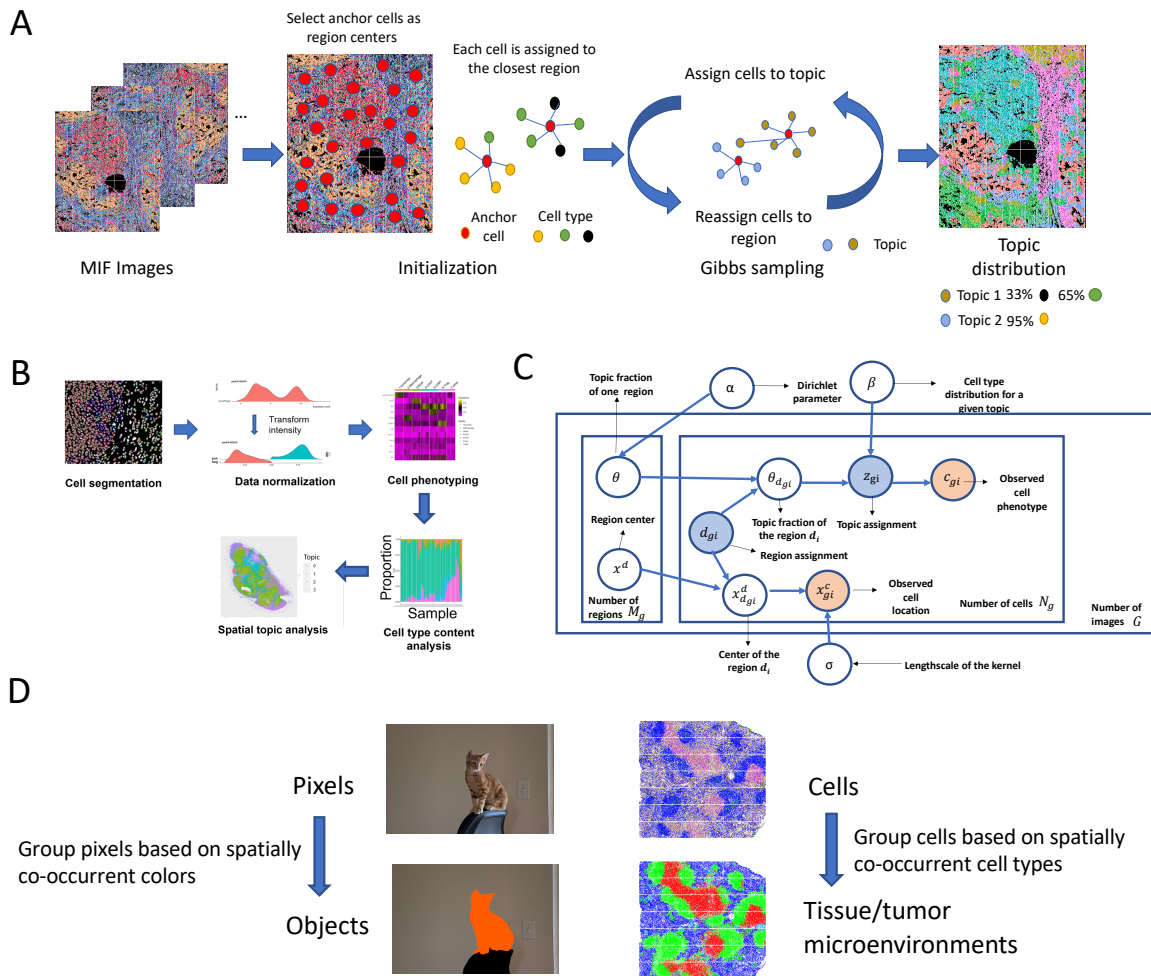
561 **CODEX Mouse Spleen.** We used the cell type annotation, marker expression level, and  
562 imaging coordinates from the original paper [2]. The image dataset can be downloaded from  
563 <https://data.mendeley.com/datasets/zjnpwh8m5b/1>. For cell coordinates, we only use the X and Y axes  
564 of the samples, ignoring Z axis. However, the result is similar when considering all three dimensions.  
565 **IMC Healthy Lung.** We used the cell type annotation, marker expression level, cell imaging coordi-  
566 nates, and cell UTAG domain labels in the original paper [13]. This image dataset can be downloaded  
567 from <https://zenodo.org/records/6376767>.  
568 **mIF Melanoma.** This is one of the whole-slide images from our internal mIF melanoma tissue  
569 samples [27]. Those whole tissue sections were stained using Ultivue UltiMapper I/O Immuno8 Kit  
570 (Cambridge, MA, USA) containing CD8, PD-1, PD-L1, CD68, CD3, CD20, FoxP3, and pancytokeratin  
571 + SOX10 (panCK-SOX10) followed by opal tyramide staining containing TCF1/7, TOX, Ki67, LAG-3.  
572 The whole imaging preprocessing pipeline has been previously described [27]. Here, we used only the cell  
573 phenotypes (classified based on marker expression of CD8, panCK-SOX10, CD68, CD3, CD20, FoxP3)  
574 and cell locations as the input of *SpaTopic*.

## 575 Simulation

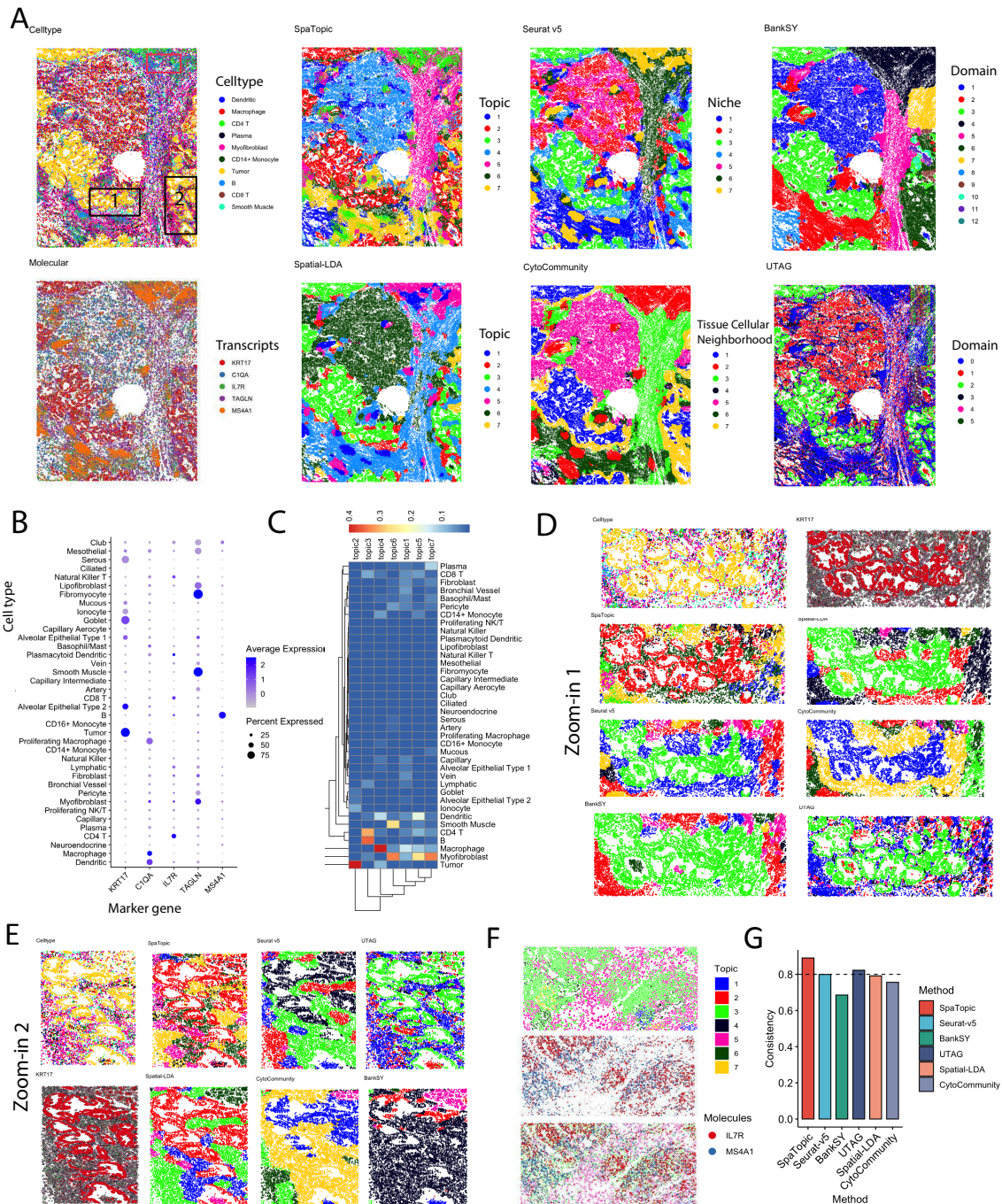
576 We tested methods on simulated datasets of different scales to benchmark the scalability of *SpaTopic*  
577 with an increasing number of cells in images. We randomly sampled 10k, 40k, 90k, 160k, and 250k pixels  
578 from an image, similar to the simulation method described in [21], to represent cell locations. We did not  
579 simulate gene expression levels for every individual cell. Instead, for each domain, we randomly sampled  
580 cells with domain-specific cell type distributions, with parameters simulated from  $Dirichlet(1, 1, 1, 1, 1)$ ,  
581 anticipating five distinct cell types per domain. Five unique datasets were generated for each simulation  
582 scenario. We also scaled the X and Y axes to maintain consistent cell densities across all simulation  
583 scenarios.



584 **Figures**

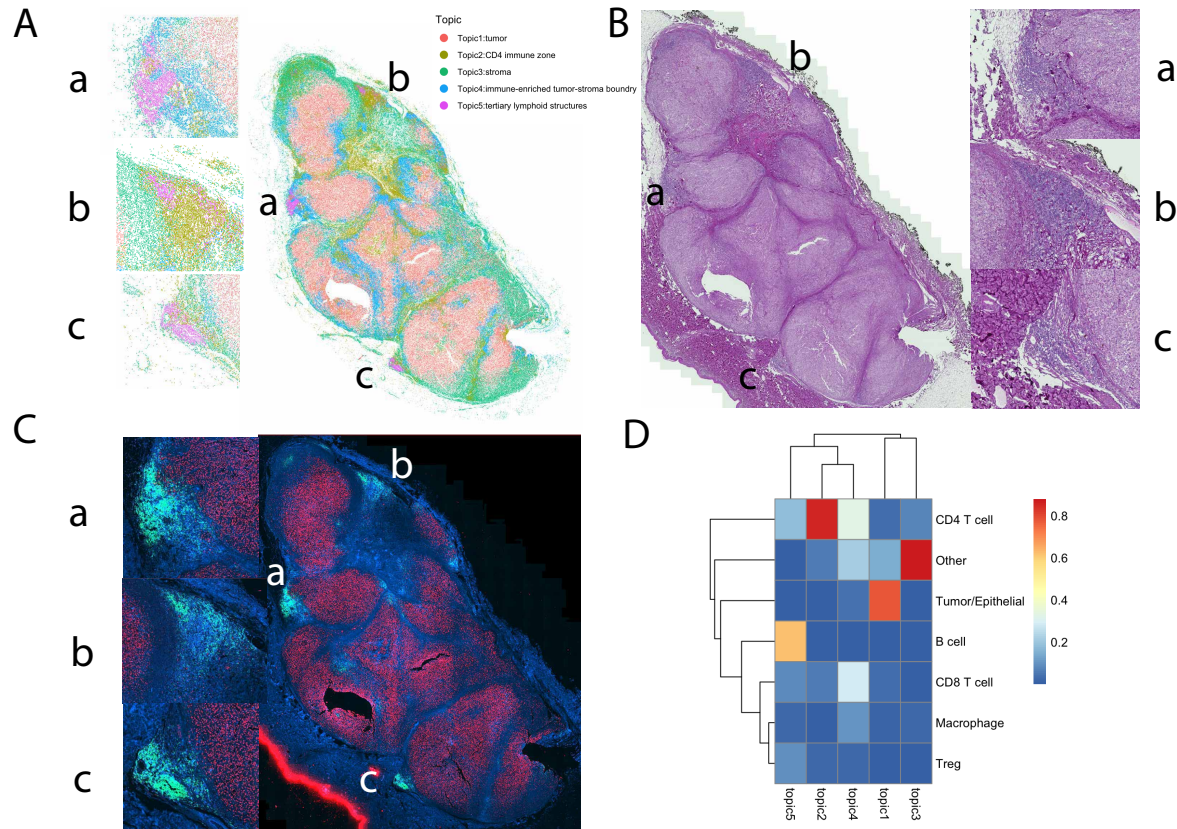


**Fig. 1** *SpaTopic* unsupervisedly identifies distinct tissue microenvironments across images, utilizing topic model concepts in computer vision. **A.** Overview of *SpaTopic*. *SpaTopic* identifies biologically relevant topics across multiple images, while each topic is a distribution of cell types, reflecting the spatial tissue architecture across images. **B.** Image analysis pipeline designed for multiplexed immunofluorescence images. *SpaTopic* is designed as a critical step for spatial analysis after cell phenotyping. **C.** Graphic representation for *SpaTopic*. The observed and hidden variables are colored orange and blue accordingly. **D.** *SpaTopic* groups cells in an unsupervised manner based on spatially co-occurring cell types, similar to image segmentation based on spatially co-occurring colors in the photo.

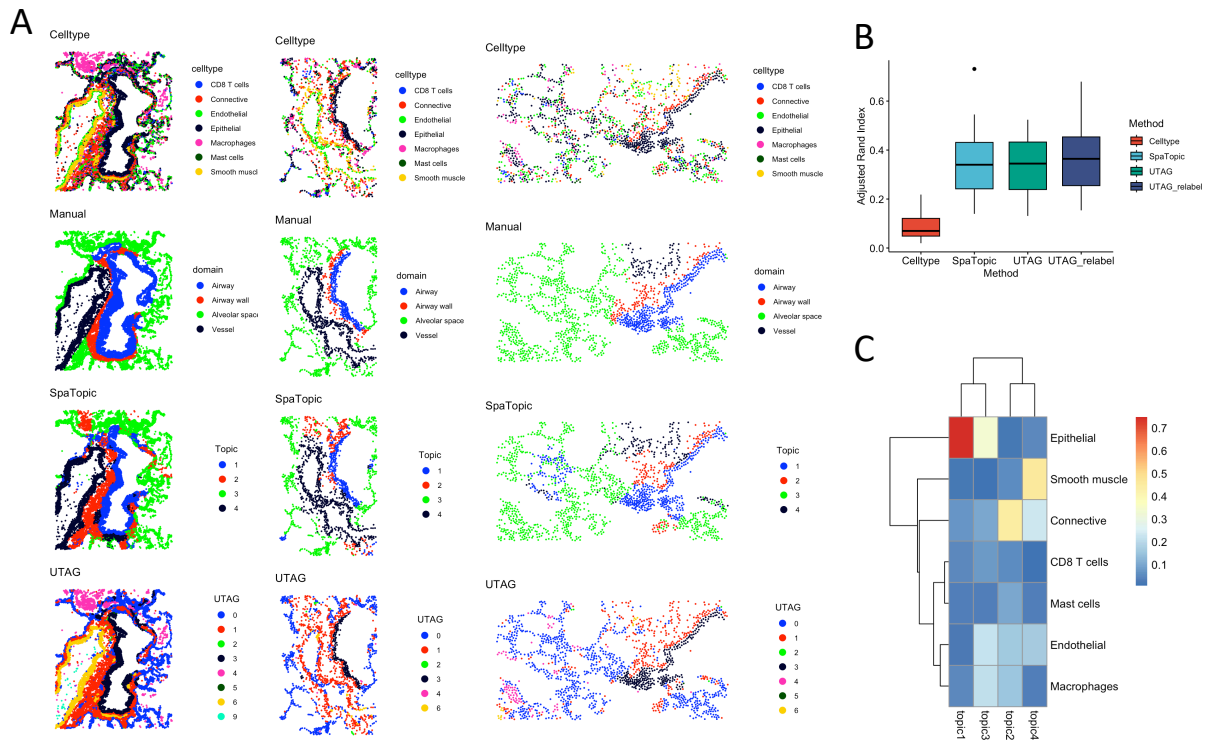


**Fig. 2** *SpaTopic* better detects tumor microenvironment in a nanostring human non-small cell lung cancer tissue image. A. We compare *SpaTopic*, Seurat v5, Spatial-LDA, CytoCommunity, Banksy, and UTAG results on the human lung tumor tissue samples. We also visualize the distribution of the top 10 most abundant cell types and five unique mRNA molecules (*KRT17*, *C1QA*, *IL7R*, *TAGLN*, *MS4A1*), showing the tissue architecture. We only show up to a total of 20k molecules due to limitations in visualization. B. Dot plots showing gene marker expression across all 38 annotated cell types. *KRT17*, *C1QA*, *IL7R*, *TAGLN*, and *MS4A1* are marker genes for tumor, macrophage, CD4 T, stroma, and B cells, respectively. C. Heatmap shows per-topic cell type composition. Topic 2 represents tumor regions. The other topics represent distinct immune-enriched stroma regions, including topic 3, which captures the lymphoid structure in the lung tissue consisting of B cells and CD4 T cells, and topic 4, which is a macrophage-enriched stroma region. D. and E. *SpaTopic* can better capture the local structure of the lung tumor tissue. F. Topic 3 (green) captures the lymphoid structures, consistent with the distribution of *IL7R* (CD4 T cells, red) and *MS4A1* (B cells, blue). G. We compare the consistency of different results, presenting the percentage of cells in the identified tumor domains expressing the *KRT17* gene. *SpaTopic* and UTAG generally show higher consistency than other methods.

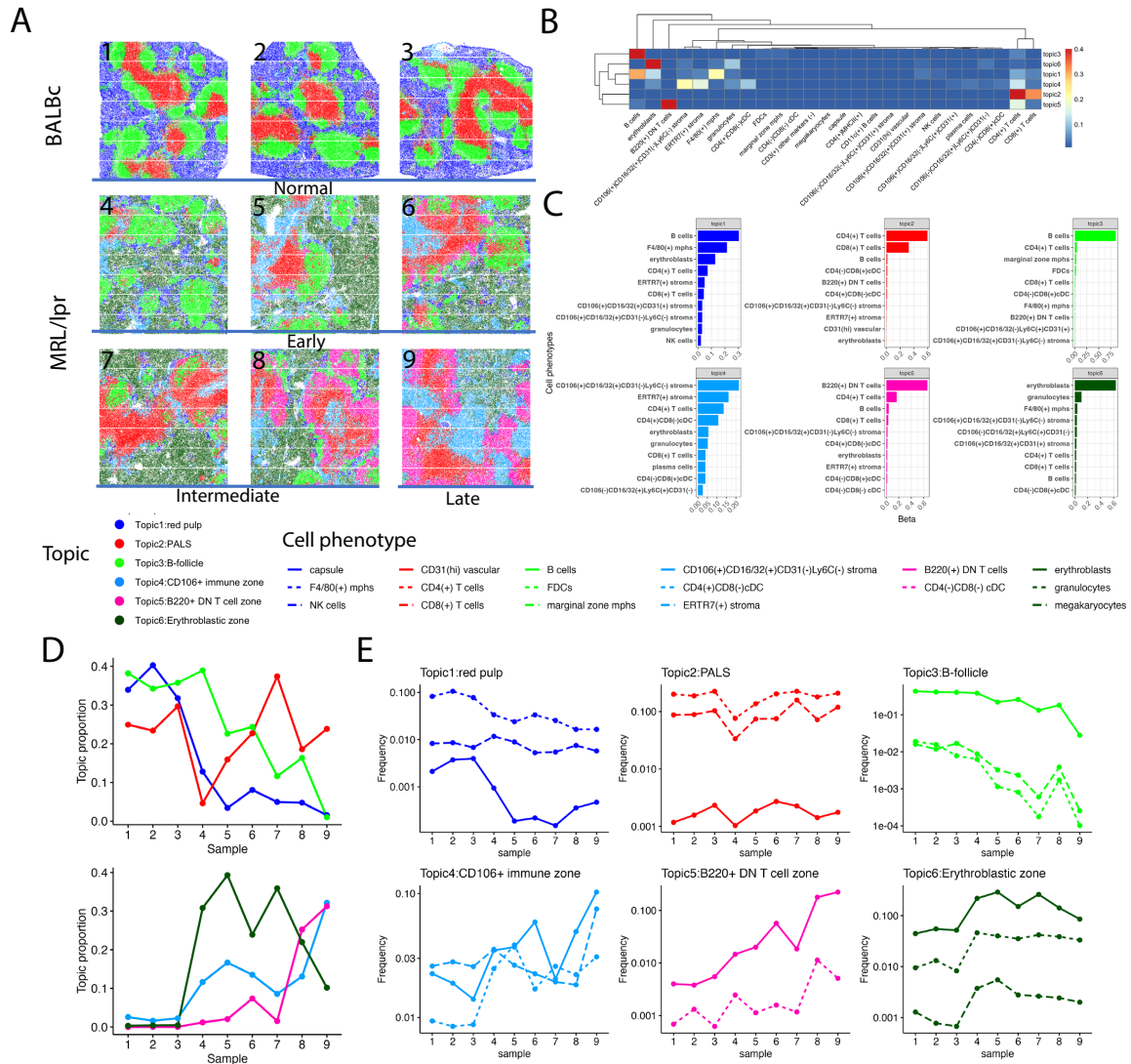




**Fig. 3** *SpaTopic* identifies tertiary lymphoid structures from a whole-slide melanoma tissue sample. A. *SpaTopic* identifies five topics from the whole-slide melanoma tissue sample: topic 1 for tumor region, topic 2 for CD4 T cell region, topic 3 for stroma region, topic 4 for immune-enriched stroma-tumor boundary, topic 5 for the potential tertiary lymphoid structures, with three Region of Interests (ROIs) highlighted in the subfigures. B. H&E staining images for the whole-slide melanoma tissue sample and three ROIs. C. Merged mIF image for the whole-slide melanoma tissue sample and three ROIs with three channels: PANCK/SOX10 (red), CD3 (royal blue), and CD20 (green). D. Heatmap shows per-topic cell type composition for the five topics identified by *SpaTopic*. Topic 5 (tertiary lymphoid structures) mainly consists of B cells and CD4 T cells, with a small proportion of CD8 T cells and regulatory T (Treg) cells.

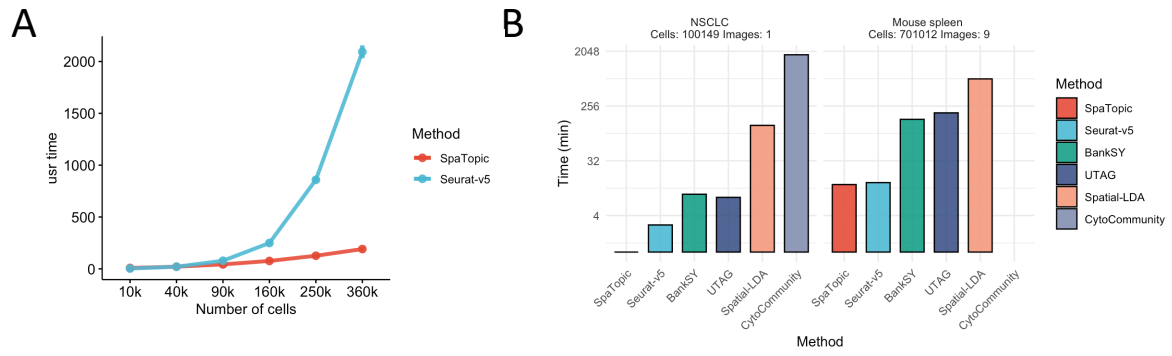


**Fig. 4** *SpaTopic* recovers spatial domain architecture from cell type spatial layout in healthy lung. A. Spatial distribution of cell type annotations, manual domain annotations, and spatial domains recovered by *SpaTopic* and UTAG from the healthy lung tissue samples. B. Consistency comparing manual spatial domain to the cell type annotation, *SpaTopic*, and UTAG with and without ad hoc relabel across 26 images. UTAG results are obtained from the original publication. C. Heatmap shows per-topic cell type composition for the four topics identified by *SpaTopic*.



**Fig. 5** *SpaTopic* captures main dynamics in tissue architecture of normal and diseased mouse spleen. A. Six topics were identified by *SpaTopic* across nine mouse spleen samples representing normal (BALB/c 1-3) and different disease stages: early (MRL 4-6), intermediate (MRL 7-8), and late (MRL 9). B. Heatmap shows per-topic cell type composition for the six main topics identified by *SpaTopic*. Based on cell type compositions, the first three topics are labeled as red pulp, PALS (periarteriolar lymphoid sheath), and B-follicle in normal mouse spleen tissue, while the other three topics are unique to the disease stages. C. Barplots show the top 10 per-topic cell types for the six main topics identified by *SpaTopic*. D. Dynamic change in the topic proportion of the six topics during disease progression. Normal spleen samples are primarily characterized by topics 1, 2, and 3, which reflect red pulp (mixed of B cells, erythroblasts, and F4/80(+) mphs), B-follicle (most B cells), and PALS (mixed of CD8 T cells and CD4 T cells), respectively. There is an increase in Topic 1 and depletion of Topic 6 in MRL samples, representing much fewer B cells and F4/80(+) mphs but more granulocytes and erythroblasts in the red pulp regions. Topic 1 (mainly B220(+) DN T cells and CD4(+) T cells) is enriched in tissue at late disease stage. E. Dynamic change of key immunological cell types within each topic, identified by FREX (omega = 0.9) and lift metrics (See Figure S5).





**Fig. 6** *SpaTopic* is scalable to large-scale images and can be run on a regular laptop within minutes. A. Runtime of *SpaTopic* (region radius  $r = 60$ ) and *Seurat-v5* on simulated datasets for increasing cell numbers. B. Runtime of *SpaTopic*, *Seurat-v5*, *BankSY*, *UTAG*, *Spatial-LDA*, and *CytoCommunity* on large-scale nanostring and mouse spleen datasets. All methods were benchmarked on a standard MacBook Air (M2, 2022) unless exceeding the memory limitation.

585 **Supplementary information.** The supplemental information was provided, including supplemental  
586 tables and figures.

587 **Acknowledgements.** We would like to thank computational support from MSK-MIND. This work  
588 is supported in part by the MSKCC Society, the V Foundation, the Parker Institute for Cancer  
589 Immunotherapy, NIH P30 CA008748, NIH R01 CA276286, and the MSK-MIND consortium.

590 **Declaration of Interests.** J.W.S. Research funding—IO Biotech (Inst), Regeneron (Inst), Daiichi  
591 Sankyo (Inst); Consulting or advisory role—IO Biotech; M.A.P. Consulting or Advisory Role - Bristol-  
592 Myers Squibb; Cancer Expert Now; Chugai Pharma; Eisai; Erasca, Inc; Intellisphere; Merck; MJH  
593 Associates; Nektar; Novartis; Pfizer; WebMD; Research Funding - Array BioPharma (Inst); Bristol-Myers  
594 Squibb (Inst); Infinity Pharmaceuticals (Inst); Merck (Inst); Novartis (Inst); Rgenix (Inst); K.S.P. Stock  
595 ownership in 23 and Me, Vincerx, Eyepoint, & Kyverna; C.E.K, Stock ownership in Johnson & Johnson;  
596 M.K.C. BMS— Research support (Inst), advisory role/consulting; Medimmune - advisory role/consulting;  
597 Immunocore—advisory role/consulting; Merus—family member employee; X.P., J.L., M.Y., M.B., R.S.,  
598 F.D.E. No disclosures;

599 **Data availability.** All public datasets we used in the study can be downloaded online, with analysis  
600 details described in the Method section. Our in-house Melanoma dataset will be made public available  
601 with the analysis paper [27].

602 **Code availability.** The R package is available on Github (<https://github.com/xiyupeng/SpaTopic/>)  
603 with a tutorial (<https://xiyupeng.github.io/SpaTopic/>). The R package is also available on CRAN  
604 (<https://cloud.r-project.org/package=SpaTopic>). The first version of the R package was officially released  
605 on CRAN on Jan 17, 2024.

606 **Author contribution.** X.P. contributed to the original draft, developed the statistical model, and  
607 wrote the software. X.P., J.W.S., R.S., and K.S.P. developed the initial study concept. X.P., R.S., K.S.P.,  
608 and J.L. developed the algorithm. X.P., C.E.K contributed to the R package. X.P. J.W.S., C.E.K, F.E.,  
609 M.Y., M.B. analyzed the data. R.S., K.S.P., M.A.P, and M.K.C. oversaw all data generation and analysis.  
610 X.P., J.W.S., F.E., J.L., M.B., R.S., and K.S.P. edited the manuscript. All authors reviewed and approved  
611 the final manuscript.

## 612 References

- 613 [1] Keren, L. *et al.* A Structured Tumor-Immune Microenvironment in Triple Negative Breast Cancer  
614 Revealed by Multiplexed Ion Beam Imaging. *Cell* **174**, 1373–1387.e19 (2018).
- 615 [2] Goltsev, Y. *et al.* Deep Profiling of Mouse Splenic Architecture with CODEX Multiplexed Imaging.  
616 *Cell* **174**, 968–981.e15 (2018).
- 617 [3] Ko, J. *et al.* Spatiotemporal multiplexed immunofluorescence imaging of living cells and tissues with  
618 bioorthogonal cycling of fluorescent probes. *Nature Biotechnology* **40**, 1654–1662 (2022).
- 619 [4] Hoch, T. *et al.* Multiplexed imaging mass cytometry of the chemokine milieu in melanoma  
620 characterizes features of the response to immunotherapy. *Science Immunology* **7**, eabk1692 (2022).
- 621 [5] Moldoveanu, D. *et al.* Spatially mapping the immune landscape of melanoma using imaging mass  
622 cytometry. *Science Immunology* **7**, eabi5072 (2022).
- 623 [6] Nirmal, A. J. *et al.* The Spatial Landscape of Progression and Immunoediting in Primary Melanoma  
624 at Single-Cell Resolution. *Cancer Discovery* **12**, 1518–1541 (2022).
- 625 [7] McCaffrey, E. F. *et al.* The immunoregulatory landscape of human tuberculosis granulomas. *Nature*  
626 *Immunology* **23**, 318–329 (2022).
- 627 [8] Helmink, B. A. *et al.* B cells and tertiary lymphoid structures promote immunotherapy response.  
628 *Nature* **577**, 549–555 (2020).

- 629 [9] Mature tertiary lymphoid structures predict immune checkpoint inhibitor efficacy in solid tumors  
630 independently of PD-L1 expression. *Nature Cancer* **2**, 794–802 (2021).
- 631 [10] Cabrita, R. *et al.* Tertiary lymphoid structures improve immunotherapy and survival in melanoma.  
632 *Nature* **577**, 561–565 (2020).
- 633 [11] Fridman, W. H. *et al.* B cells and tertiary lymphoid structures as determinants of tumour immune  
634 contexture and clinical outcome. *Nature Reviews Clinical Oncology* **19**, 441–457 (2022).
- 635 [12] Feng, Y. *et al.* Spatial analysis with SPIAT and spaSim to characterize and simulate tissue  
636 microenvironments. *Nature Communications* **14**, 2697 (2023).
- 637 [13] Kim, J. *et al.* Unsupervised discovery of tissue architecture in multiplexed imaging. *Nature Methods*  
638 **19**, 1653–1661 (2022).
- 639 [14] Chen, Z., Soifer, I., Hilton, H., Keren, L. & Jovic, V. Modeling multiplexed images with spatial-lda  
640 reveals novel tissue microenvironments. *Journal of Computational Biology* (2020).
- 641 [15] Patrick, E. *et al.* Spatial analysis for highly multiplexed imaging data to identify tissue  
642 microenvironments. *Cytometry Part A* **103**, 593–599 (2023).
- 643 [16] Hao, Y. *et al.* Dictionary learning for integrative, multimodal and scalable single-cell analysis.  
644 *Nature Biotechnology* **42**, 293–304 (2024).
- 645 [17] Singhal, V. *et al.* BANKSY unifies cell typing and tissue domain segmentation for scalable spatial  
646 omics data analysis. *Nature Genetics* **56**, 431–441 (2024).
- 647 [18] Hu, Y. *et al.* Unsupervised and supervised discovery of tissue cellular neighborhoods from cell  
648 phenotypes. *Nature Methods* **21**, 267–278 (2024).
- 649 [19] Li, Z. & Zhou, X. BASS: multi-scale and multi-sample analysis enables accurate cell type clustering  
650 and spatial domain detection in spatial transcriptomic studies. *Genome Biology* **23**, 1–35 (2022).
- 651 [20] Chidester, B., Zhou, T., Alam, S. & Ma, J. SpiceMix enables integrative single-cell spatial modeling  
652 of cell identity. *Nature Genetics* **55**, 78–88 (2023).
- 653 [21] Shang, L. & Zhou, X. Spatially aware dimension reduction for spatial transcriptomics. *Nature*  
654 *Communications* **13**, 1–22 (2022).
- 655 [22] Wang, X. & Grimson, E. Platt, J., Koller, D., Singer, Y. & Roweis, S. (eds) *Spatial latent dirichlet*  
656 *allocation*. (eds Platt, J., Koller, D., Singer, Y. & Roweis, S.) *Advances in Neural Information*  
657 *Processing Systems*, Vol. 20 (Curran Associates, Inc., 2007).
- 658 [23] Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573–3587 (2021).
- 659 [24] Rakaee, M. *et al.* Tertiary lymphoid structure score: a promising approach to refine the TNM staging  
660 in resected non-small cell lung cancer. *British Journal of Cancer* **124**, 1680–1689 (2021).
- 661 [25] Aran, D. *et al.* Reference-based analysis of lung single-cell sequencing reveals a transitional  
662 profibrotic macrophage. *Nature Immunology* **20**, 163–172 (2019).
- 663 [26] Meng, G. *et al.* imply: improving cell-type deconvolution accuracy using personalized reference  
664 profiles. *Genome Medicine* **16**, 65 (2024).
- 665 [27] Smithy, J. W. *et al.* Spatial assessment of stromal b cell aggregates predicts response to checkpoint  
666 inhibitors in unresectable melanoma. *medRxiv* (2024).

- 667 [28] Peng, X. *et al.* A topic modeling approach reveals the dynamic T cell composition of peripheral  
668 blood during cancer immunotherapy. *Cell Reports Methods* **3**, 100546 (2023).
- 669 [29] Roberts, M. E., Stewart, B. M. & Tingley, D. stm: An R Package for Structural Topic Models.  
670 *Journal of Statistical Software* **91**, 1–40 (2019).
- 671 [30] Yosofvand, M. *et al.* Spatial immunophenotyping from whole-slide multiplexed tissue imaging using  
672 convolutional neural networks. *bioRxiv* (2024).
- 673 [31] Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data using  
674 empirical Bayes methods. *Biostatistics* **8**, 118–127 (2006).
- 675 [32] Kuett, L. *et al.* Three-dimensional imaging mass cytometry for highly multiplexed molecular and  
676 cellular mapping of tissues and the tumor microenvironment. *Nature Cancer* **3**, 122–133 (2022).
- 677 [33] Datta, A., Banerjee, S., Finley, A. O. & Gelfand, A. E. Hierarchical nearest-neighbor gaussian  
678 process models for large geostatistical datasets. *Journal of the American Statistical Association*  
679 **111**, 800–812 (2016).
- 680 [34] Weber, L. M., Saha, A., Datta, A., Hansen, K. D. & Hicks, S. C. nnSVG for the scalable identification  
681 of spatially variable genes using nearest-neighbor Gaussian processes. *Nature Communications* **14**,  
682 4059 (2023).
- 683 [35] Griffiths, T. L. & Steyvers, M. Finding scientific topics. *Proceedings of the National Academy of*  
684 *Sciences* **101**, 5228–5235 (2004).
- 685 [36] Heinrich, G. *Parameter estimation for text analysis* (2009).
- 686 [37] Newman, D., Asuncion, A., Smyth, P. & Welling, M. Distributed algorithms for topic models. *The*  
687 *Journal of Machine Learning Research* **10**, 1801–1828 (2009).
- 688 [38] Gelman, A., Carlin, J. B., Stern, H. S. & Rubin, D. B. *Bayesian Data Analysis: Second Edition*  
689 *Texts in Statistical Science* (CRC Press, 2004).