



HHS Public Access

Author manuscript

Nat Biotechnol. Author manuscript; available in PMC 2022 January 17.

Published in final edited form as:

Nat Biotechnol. 2021 November ; 39(11): 1375–1384. doi:10.1038/s41587-021-00935-2.

Spatial transcriptomics at subspot resolution with BayesSpace

Edward Zhao^{1,2}, Matthew R. Stone³, Xing Ren¹, Jamie Guenthoer⁴, Kimberly S. Smythe⁵, Thomas Pulliam⁶, Stephen R. Williams⁷, Cedric R. Uyttingco⁷, Sarah E. B. Taylor⁷, Paul Nghiem^{5,6,8}, Jason H. Bielas^{3,9,10}, Raphael Gottardo^{1,2,✉}

¹Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA

²Department of Biostatistics, University of Washington, Seattle, WA, USA

³Fred Hutch Innovation Laboratory, Immunotherapy Integrated Research Center, Fred Hutchinson Cancer Research Center, Seattle, WA, USA

⁴Human Biology Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA

⁵Clinical Research Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA

⁶Department of Medicine, Division of Dermatology, University of Washington, Seattle, WA, USA

⁷10x Genomics, Pleasanton, CA, USA

⁸Seattle Cancer Care Alliance, Seattle, WA, USA

⁹Translational Research Program, Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA

¹⁰Department of Pathology, University of Washington, Seattle, WA, USA

Abstract

Recent spatial gene expression technologies enable comprehensive measurement of transcriptomic profiles while retaining spatial context. However, existing analysis methods do not address the limited resolution of the technology or use the spatial information efficiently. Here, we introduce BayesSpace, a fully Bayesian statistical method that uses the information from spatial neighborhoods for resolution enhancement of spatial transcriptomic data and for clustering

✉ **Correspondence and requests for materials** should be addressed to R.G. rgottard@fredhutch.org.

Author contributions

E.Z. and R.G. formulated the method and wrote the paper. M.R.S. and E.Z. developed software. E.Z., M.R.S. and X.R. analyzed data. J.G., K.S.S. and T.P. contributed to annotation and interpretation of cancer samples. C.R.U., S.R.W. and S.E.B.T. prepared and contributed to analysis of the IDC sample. P.N., J.H.B. and R.G. supervised the project.

Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41587-021-00935-2>.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41587-021-00935-2>.

Peer review information *Nature Biotechnology* thanks the anonymous reviewers for their contribution to the peer review of this work.

Reprints and permissions information is available at www.nature.com/reprints.

analysis. We benchmark BayesSpace against current methods for spatial and non-spatial clustering and show that it improves identification of distinct intra-tissue transcriptional profiles from samples of the brain, melanoma, invasive ductal carcinoma and ovarian adenocarcinoma. Using immunohistochemistry and an in silico dataset constructed from scRNA-seq data, we show that BayesSpace resolves tissue structure that is not detectable at the original resolution and identifies transcriptional heterogeneity inaccessible to histological analysis. Our results illustrate BayesSpace's utility in facilitating the discovery of biological insights from spatial transcriptomic datasets.

Knowledge of the spatial location of transcript expression can provide vital insights into biological function and pathology. Single-cell RNA sequencing (scRNA-seq) achieves high-throughput and high-resolution profiling of gene expression, but because tissue is dissociated for sample preparation, spatial information is not retained. Recent methods for high-throughput profiling of gene expression while retaining spatial information allow analyses to be made within the context of the biological tissue¹. Studies performed with the Spatial Transcriptomics (ST) platform and the improved Visium platform have already generated insights into diverse areas such as tumor heterogeneity^{2,3}, brain function⁴ and the pathophysiology of sepsis⁵. The primary technological limitation of these spatial gene expression platforms is resolution, with the unit of observation being spots that are 100µm in diameter on the ST platform and 55µm in diameter on the Visium platform. As such, the number of cells within a spot may range from one to 30 on the Visium platform and up to 200 on the older ST platform, depending on the biological tissue⁶. Alternative approaches include fluorescence in situ hybridization (FISH) technologies, such as seqFISH and multiplexed error-robust FISH, and other recently developed spatial sequencing methods, such as Slide-seq and ZipSeq⁷⁻¹⁰. While these methods provide increased resolution, most are lower throughput, less sensitive, rely on custom protocols or are not widely available.

Here, we propose BayesSpace, a computational method that uses the neighborhood structure in spatial transcriptomic data to increase the resolution to the subspot level (Fig. 1a). Our method draws from the existing literature for use of Bayesian statistics to achieve super-resolution images¹¹⁻¹³. In contrast to existing deconvolution methods using scRNA-seq data¹⁴⁻¹⁶, the enhanced-resolution modeling of BayesSpace, which approaches single-cell resolution with the Visium platform, does not require independent single-cell data and allows us to infer the spatial arrangement of subspots. While integration with scRNA-seq is appealing, it may be costly if matched samples are used or introduce bias if publicly available references are used. Furthermore, deconvolved mixtures are still only spatially resolved at the original scale of the ST or Visium technology, and the neighborhood structure of cell types cannot be recovered.

In addition, there is a need for new statistical methods for the analysis of spatial gene expression data that efficiently use the available spatial information. Clustering is an important step in the analysis of such data that allows downstream analyses, such as cell type or tissue annotation and differential expression, to provide unbiased biological insights. Existing analyses of spatial gene expression data often rely on clustering methods for non-spatial scRNA-seq data^{2,4}. The additional spatial information available from ST

and Visium can help address the analytical challenges of sparsity and noise by smoothing over adjacent spots, which are more likely to have similar transcriptomic profiles. Zhu et al.¹⁷ proposed a hidden Markov random field model (HMRF) for clustering of low-resolution in situ hybridization data into distinct spatial domains by jointly modeling gene expression and spatial neighborhood structure¹⁷. This approach was later adapted for use with high-throughput spatial transcriptomic data through selection of spatially differentially expressed genes before clustering¹⁸. Another recently developed spatial clustering algorithm is stLearn, which uses deep learning features extracted from histopathological images as well as expression of neighboring spots to spatially smooth data¹⁹.

BayesSpace enables spatial clustering by modeling a low-dimensional representation of the gene expression matrix and encouraging neighboring spots to belong to the same cluster via a spatial prior (Fig. 1b). Our method draws from previously developed spatial statistics methods for image analysis and microarray data^{20,21}. Compared with previous approaches, BayesSpace allows for a more flexible specification of the clustering structure and error term than alternative approaches. From a user perspective, BayesSpace is accessible in that it takes the widely used Bioconductor SingleCellExperiment object as input²², does not require the additional task of marker gene preselection and involves minimal parameter tuning.

Using several datasets, we show that BayesSpace improves the identification of spatially distributed tissue domains through spatial clustering and enhances the resolution of gene expression maps. We use immunohistochemistry as a ground truth in two cancer samples to validate that our enhanced-resolution clustering identifies a tissue structure consistent with cell surface markers, and we report examples of transcriptional heterogeneity in the tumor microenvironment not achievable by immunohistochemical analyses alone. Furthermore, using in silico spatial transcriptomic datasets generated from aggregating scRNA-seq data, we show that BayesSpace can recover the true spatial structure at near single-cell resolution.

Results

Spatial clustering improves identification of known layers in brain tissues.

Recently, Maynard et al.⁴ presented Visium spatial expression profiles of 12 dorsolateral prefrontal cortex (DLPFC) samples, as well as manual annotations of the six cortical layers and white matter for each sample as part of the spatialLIBD package⁴ (Fig. 2a). Maynard et al.⁴ annotated DLPFC layers by considering cytoarchitecture and selected gene markers. Here, we evaluate BayesSpace's ability to identify distinct layer-specific expression profiles and compare its performance to other spatial and non-spatial clustering methods. Specifically, we compare the performance of four non-spatial algorithms commonly applied to scRNA-seq data, *k*-means, mclust²³, Louvain²⁴ and SC3 (ref. ²⁵); two recently published spatial clustering algorithms, HMRF (as implemented in the Giotto package)¹⁸ and stLearn¹⁹; and the clustering partitions originally reported by Maynard et al.⁴ in the spatialLIBD package, which involve Walktrap clustering of spatial coordinates and principal components (PCs) calculated from highly variable genes (HVGs) or known layer-specific marker genes. Following the methodology of Maynard et al.⁴, we use the adjusted Rand

index (ARI) to quantify similarity between cluster labels and manual annotations, which are considered the ground truth.

BayesSpace substantially outperforms the original spatialLIBD clustering partitions, as well as all non-spatial clustering algorithms and spatial clustering methods developed for spatial transcriptomic data (Fig. 2b). BayesSpace and the non-spatial methods were applied on 15 PCs calculated from the top 2,000 HVGs. Spatial clustering methods Giotto and stLearn were implemented based on the original authors' recommended parameters (Supplementary Note). We also show Giotto and stLearn results using precomputed PCs from BayesSpace to provide a more controlled comparison, although we found that this did not improve either method's performance (Supplementary Fig. 1). As an example, in sample 151673, we found that only SC3 (ARI = 0.42), mclust (ARI = 0.42), stLearn (ARI = 0.37) and BayesSpace (ARI = 0.55) generated clusters that qualitatively followed the expected layer pattern (Fig. 2c). Most clustering partitions aside from BayesSpace exhibited substantial noise and lack of clear spatial separation between clusters. By contrast, BayesSpace leveraged spatial information to smooth the data and provided distinct layers of clusters. The t-distributed error model that BayesSpace uses is particularly robust against outliers in clusters, which may be driven by technical artifacts generated during sample preparation or downstream analyses (Supplementary Fig. 2). Additionally, BayesSpace's runtime and memory footprints are comparable to those of other spatial clustering methods, requiring 27 min of wall time and 9.6 GB of memory in this sample (Supplementary Fig. 3).

Increased resolution clustering leads to identification of known tissue structures missed by other methods.

We used BayesSpace to analyze a melanoma ST sample first annotated and described by Thrane et al.². As the manual annotation identified regions of melanoma, stroma and lymphoid tissue and left an additional area unannotated (Fig. 3a), we ran spatial clustering with $k = 4$ clusters (Fig. 3b). The resulting clusters corresponded well with the manually annotated tissue types. Furthermore, the melanoma tissue was split into the central region of the tumor and an outer ring of mixed tumor and lymphoid tissue. BayesSpace enhanced spatial clustering provided a higher-resolution map of the tissue types (Fig. 3c). Notably, the enhancement identified lymphoid regions along the tumor border and possible immune infiltration into the tumor that could not be discerned at the original resolution. These regions were also largely not identified by other clustering methods (Supplementary Fig. 4). While most clustering methods identified heterogeneity between the periphery and the center of the tumor, only SC3, Giotto and subspot-level BayesSpace identified lymphoid regions proximal to the tumor, with BayesSpace providing higher resolution and more robust signal (Supplementary Fig. 4). Finally, we also ran BayesSpace at the spot level using five and six clusters, identifying potential heterogeneity within the stroma region (Supplementary Fig. 4).

Using the enhanced PCs, we can generate high-resolution maps of individual genes or expression profiles for major cell types as described in the Methods. Differential expression analysis performed on enhanced-resolution gene expression indicated that the lymphoid regions had a distinct expression profile. We observed elevated expression of lymphocyte

markers such as *CD52* and *MS4A1* and lower expression of melanoma markers such as MCAM and *SPP1* relative to that of the surrounding tumor border (Supplementary Fig. 4). Enhanced-resolution differential expression analysis between the four clusters highlighted additional spatial variation in gene expression (Fig. 3d). In the stroma (cluster 2), expression levels were higher for extracellular matrix proteins such as those encoded by *DCN* and *COL3A1*. Furthermore, we revealed intratumor heterogeneity between the border and center of the tumor (clusters 1 and 3, respectively), with higher chemokine (*CXCL9*, *CXCL10*) activity at the border and elevated expression of genes related to cell proliferation (*HSPB1*) and metastasis (*ATPIA1*) at the center^{26,27}.

We defined tumor cell (*PMEL*), fibroblast (*COL1A1*), B cell (*CD19*, *MS4A1*), T cell (*CD2*, *CD3D*, *CD3E*, *CD3G*, *CD7*) and macrophage (*CD14*, *FCGR1A*, *FCGR1B*) expression profiles based on one or more marker genes from existing literature²⁸. The enhanced expression profiles provided noticeably higher spatial resolution (Fig. 3e). In particular, we could more clearly observe immune expression on the periphery of the tumor. The contrast between *PMEL* expression in the tumor, stroma and lymphoid tissue was also more apparent with enhanced resolution.

Immunohistochemistry validates enhanced-resolution clusters.

To validate our enhanced-resolution clustering and gene expression, we analyzed an unreported breast cancer sample: an estrogen receptor-positive (ER⁺), progesterone receptor-negative (PR⁻), human epidermal growth factor receptor (HER)2-amplified (HER⁺) invasive ductal carcinoma (IDC) prepared on the Visium platform with immunofluorescence staining for 4,6-diamidino-2-phenylindole (DAPI) (staining nuclei) and CD3 (staining T cells) (Supplementary Note and Supplementary Fig. 5). We additionally analyzed a dataset published by 10x Genomics: an endometrial adenocarcinoma of the ovary (ovarian cancer; OC) sequenced on the Visium platform and stained with immunofluorescence for DAPI, pan-cytokeratin (staining epithelial tissue) and CD45 (staining leukocytes) (Supplementary Fig. 6). After examination by a pathologist, out-of-focus and overexposed regions were excluded from the analysis (Methods and Supplementary Figs. 7 and 8). Cell segmentation of in-focus areas (IDC, $n = 2,929$ of 4,727 spots; OC, $n = 2,041$ of 3,493 spots) identified a median of 21 cells per spot in the IDC tissue and 19 cells per spot in the OC tissue, along with a median of three cells per subspot in both tissues (Supplementary Figs. 7 and 8).

We applied BayesSpace to cluster the IDC sample into ten clusters and the OC sample into eight clusters at spot and subspot resolution, selecting the number of clusters based on the negative log-likelihood curve (Supplementary Figs. 9 and 10). We analyzed anti-CD3 and anti-CD45 intensity in the in-focus area of each tissue section (Fig. 4a,f, respectively), finding that the immunofluorescence signal correlated well with the corresponding enhanced gene expression (Pearson's $r = 0.53$ in the IDC; Fig. 4b,g). In both samples, we identified clusters enriched for the respective immune immunofluorescence signal and dichotomized the clusters into CD3- or CD45-rich and CD3- or CD45-poor areas (Fig. 4c,h and Supplementary Figs. 9 and 10). From this, we identified regions of interest (ROI) between the spot-level and enhanced clustering: areas where the enhancement increased the observed heterogeneity and many subspots flipped from immune rich to

immune poor or vice versa. We highlight six of these ROI in Fig. 4d,i to demonstrate that enhanced clustering qualitatively improves concordance of clustering with the underlying immunohistochemical stain. Specifically, we present regions where, compared to the coarser spot-level clustering, the enhanced-resolution clustering detects subspots with high underlying immunofluorescence stain intensity and refines the boundary between immune-rich and immune-poor areas.

To quantify the improvement at enhanced resolution, we compared the distribution of immunofluorescence intensity between subspots that changed classification after enhancement (for example, immune rich at the spot level and immune poor after enhancement) and subspots that maintained their classification (for example, immune rich at both the spot and subspot level). We found a significant difference in the intensity of subspots that changed classification compared to those that maintained their spot-level status (Fig. 4e,j), indicating that BayesSpace's resolution enhancement improves the accuracy of expression-based clustering with respect to an orthogonal immunohistochemistry signal.

BayesSpace distinguishes intratumoral heterogeneity in IDC.

We further analyzed the IDC tissue section to identify clusters of biological relevance. Pathologist annotation identified regions of predominantly invasive carcinoma (IC), carcinoma in situ and benign hyperplasia, from which we derived ground-truth labels for each spot (Fig. 5a and Supplementary Fig. 11). The clusters were largely consistent with histopathological annotations (cluster purity = 0.839; Fig. 5b and Supplementary Figs. 9 and 11), and we identified five clusters that corresponded to annotated regions of predominantly IC (3–6 and 9), one cluster that encompassed all annotated regions of carcinoma in situ (8), one cluster that coincided with the annotated benign hyperplasia and an invasive-appearing area (2) and three clusters corresponding to predominantly non-tumor areas (1, 7 and 10; Supplementary Fig. 11). We note that, without hematoxylin and eosin (H&E) stains or an immunofluorescent stain for a tumor marker, the tumor–stroma interface could not be fully delineated histologically and BayesSpace's enhanced clustering identified heterogeneity within the tissue that was not reflected in the annotated boundaries but was clearly supported by key tumor marker genes (Fig. 5c–e). This further supports our previous validation with immunofluorescence (Fig. 4).

Spatial expression patterns of known marker genes and differential expression analysis between these clusters were largely in accord with clinical and histopathological annotations. Consistent with the clinical report of ER⁺PR⁻HER2⁺ IDC, we observed high expression levels of genes coding for HER2 (*ERBB2*) and ER (*ESR1*) through out the tumor clusters and minimal expression of the gene coding for PR (*PGR*) in the sample (Fig. 5c and Supplementary Fig. 12). The non-tumor clusters 1, 7 and 10 were characterized by the expression of immune genes, with *PTPRC* (leukocyte-common antigen CD45) highly expressed in these clusters. We found that these clusters corresponded to distinct spatial transcriptional patterns. Cluster 1 was enriched for signatures of cell-mediated immunity, including marker genes expressed by T cells (*CD4*, *CD8A*, *CD8B*) and macrophages (*CD14*, *CD68*), while clusters 7 and 10 were enriched for genes involved in humoral immunity, particularly those encoding immunoglobulin chains (for example, *IGHG3*; Fig. 5d and

Supplementary Figs. 12–15). Compared to other non-tumor clusters, cluster 7 was also enriched for expression of *ERBB2* and tumor-associated genes, such as *ZNF703*, suggesting that this cluster represents a mixture of tumor and immune cells. Analysis of non-tumor subspots (clusters 1, 7 and 10) with CIBERSORT was consistent with differential expression results, predicting subspots in cluster 1 to have a greater abundance of T cells, while clusters 7 and 10 had higher proportions of B and plasma cells (Supplementary Fig. 16).

We found similar heterogeneity within the invasive tumor clusters. Clusters 3, 5 and 6 displayed elevated expression of known markers of cell proliferation, including genes encoding Ki-67 (*MKI67*) and cyclins, as well as genes associated with tumor progression, invasion and proliferation, including *COL1A2* (refs. 29–31), *MUC1* (refs. 32–35) and *MMP11* (refs. 30,31,36) (Fig. 5e and Supplementary Figs. 13, 15 and 17). Clusters 4 and 9 showed increased expression of *ZNF703*, an oncogene in the more aggressive, ER⁺ luminal B breast cancer subtype^{37,38} as well as that of *GRB2*, a gene implicated in breast cancer tumorigenesis^{39,40} and *BAMBI*, encoding a pseudoreceptor for TGF- β ⁴¹, the signaling pathway of which is implicated in progression to invasion³² (Fig. 5e). These spatial expression patterns suggest a transcriptional heterogeneity among compartments of invasive tumor inaccessible to histopathological analysis, demonstrating the superiority of spatial transcriptomic data over immunofluorescence alone.

BayesSpace enhances gene expression patterns to near single-cell resolution on in silico spatial data.

We conducted several simulations to demonstrate that BayesSpace clustering and resolution enhancement outperform existing methods. In the first simulation, for which we simulated data modeled on two of our experimental datasets (see Methods for details), results showed that BayesSpace spot-level clustering consistently outperformed all other methods in both the simulated melanoma and ovarian datasets (Fig. 6a). Giotto, another spatial clustering method, also outperformed all non-spatial methods but provided slightly worse performance than BayesSpace. Among the non-spatial methods, mclust and Louvain clustering performed decently.

In the second simulation, we showed that BayesSpace enhanced-resolution clustering outperformed the optimal clustering that can be achieved at the spot level in melanoma and ovarian samples that were simulated at the subspot level (Fig. 6b). In each dataset, the enhanced clustering ARI exceeded the optimal spot-level clustering in all 20 simulated replicates. This indicates that BayesSpace is able to increase the resolution of data to better recapture finer details of the ground truth.

In the third simulation, we demonstrated that BayesSpace enhanced-resolution clustering can increase the resolution of data that were simulated from real, aggregated single cells (see Methods for details). BayesSpace captures the spatial distribution of clusters better than optimal spot-level clustering, as illustrated in the spatial representation of enhanced clustering results from one replicate (Fig. 6c). In regions with high mixing of cell types, there is little to no information available to resolve cluster labels at the subspot level, but BayesSpace is still able to closely approximate the overall tissue structure at the spot level. In these cases, although it is easy to miss isolated cells due to the signal being diluted out

from the aggregation of multiple cells at the spot level, we found that BayesSpace was still able to recover some of these populations. The simulation results further supported our melanoma analyses in which our enhanced analysis recovered lymphoid structure near the tumor that was not apparent at the spot level. In all, BayesSpace enhanced clusters recapture the ground truth better than all other methods, again highlighting the superior performance of our method (Fig. 6d) and showing that BayesSpace is able to successfully enhance the resolution of spot-level data.

Enhanced-resolution clustering resolves keratinocyte structure in squamous cell carcinoma.

Finally, we also used BayesSpace to analyze a squamous cell carcinoma Visium sample first described by Ji et al.⁴². H&E-stained tissue annotated by a pathologist revealed tumor borders and other major tissue structures (Supplementary Fig. 18). We defined expression profiles for the major cell types present in the sample based on known marker genes from the literature: keratinocytes (*KRT1*, *KRT5*, *KRT10*, *KRT14*), melanocytes (*MLANA*, *DCT*, *PMEL*), myeloid cells (*LYZ*) and T cells (*CD2*, *CD3D*, *CD3E*, *CD3G*, *CD7*)^{28,42}. Keratinocytes were further separated into basal keratinocytes (*KRT5*, *KRT14*) and suprabasal keratinocytes (*KRT1*, *KRT10*), as products of *KRT5* and *KRT14* form heterodimers that localize to the basal layer of the epidermis, while products of *KRT1* and *KRT10* form heterodimers that localize to the suprabasal layer⁴³. We show that our enhanced spatial gene expression plots delineate the border between the basal and suprabasal layers more precisely than spot-level plots (Supplementary Figs. 18 and 19) and similarly find that the enhanced expression of marker genes for melanocytes, myeloid cells and T cells better match the expected patterns based on annotated tissue structures (Supplementary Fig. 18).

Discussion

BayesSpace seamlessly integrates into the spatial transcriptomic analysis workflow by taking as input preprocessed data via the widely used Bioconductor SingleCellExperiment data structure. The output is likewise stored in a SingleCellExperiment object that can be used for downstream analyses. The methods are all implemented as an R package that is openly accessible on Bioconductor.

We have demonstrated the utility of BayesSpace in identifying spatial clusters with similar expression profiles and enhancing the resolution of spatial transcriptomics. BayesSpace overcomes both the challenge in efficiently using spatial information to inform the clustering of expression data and the limited resolution of current spatial transcriptomic technology. While there are similarities in the spatial prior specification between BayesSpace and Giotto (HMRF), we highlight several differences between the methods. BayesSpace is a spatial transcriptomic model-based clustering method that uses a t-distributed error model to identify spatial clusters that are more robust to the presence of outliers caused by technical noise. BayesSpace also uses Markov chain Monte Carlo (MCMC) to estimate model parameters, while HMRF uses expectation–maximization, which might not explore the space as efficiently⁴⁴. BayesSpace also differs from Giotto (HMRF) in that it uses a

fixed precision matrix rather than a variable precision matrix across clusters, which we found to improve the stability of estimates without compromising clustering performance (Supplementary Fig. 20) and in that it uses a more reliable method for detecting the spatial neighborhood network.

Studies have not achieved subspot resolution of spatial transcriptomic data without requiring the use of additional information aside from spatial coordinates. Immunohistochemical analyses in the IDC and OC tissue sections provide validation that our subspot model accurately refines and reflects the spatial structure of the underlying tissue. Enhancement of gene expression analysis at subspot resolution allows downstream differential expression analyses to compare finer and more biologically meaningful clusters. Our analyses of differential expression in the IDC tissue section identify transcriptional heterogeneity within regions of invasive tumor that appear histologically indistinct. While histological analysis of this tissue was limited by available immunofluorescent stains, notably lacking a tumor marker or H&E stains, our results suggest the potential for spatial transcriptomics and BayesSpace to capture previously uncharacterized spatial patterns of gene expression.

The resolution enhancement approaches single-cell resolution, with approximately three cells per subspot for data acquired with the Visium platform, without the need for external single-cell data. However, there is potential for the enhanced data to be integrated with external single-cell data through deconvolution or label-transfer methods. For example, it may be possible to enhance the resolution of spot-level cell-type proportion estimates by using a Dirichlet regression model with enhanced PCs as predictors. Integration with single-cell data has the potential to improve our ability to resolve cell types in dense and complex tissues, and it is a future direction of our research.

While our work focused on the ST and Visium platforms from 10x Genomics, BayesSpace should be applicable to other platforms in which spots are arranged on a lattice. Slight modifications may be needed so that our spatial model can be used with a different neighborhood structure. Because BayesSpace models a lower-dimensional representation of data (that is, principal component analysis (PCA)), it should also be applicable to other dimensional-reduction techniques such as uniform manifold approximation and projection and possibly be applied to other data types such as protein markers and multiomics. Finally, it may also be possible to extend BayesSpace to jointly cluster spots from multiple samples given appropriate data normalizations.

Methods

Data description.

We applied BayesSpace to samples from five spatial gene expression datasets, of which four were generated on the newer Visium platform. All Visium samples that were obtained directly from 10x Genomics were procured from BioIVT:ASTERAND. Details on dataset processing and availability are provided in the Supplementary Information. The first dataset included twelve human DLPFC samples from three individuals run on the Visium platform⁴. Briefly, each sample contained approximately 4,000 spots that were manually annotated to belong to one of six DLPFC layers or white matter. The second dataset involved melanoma

samples run on the ST platform². From this dataset, we analyzed the second replicate from biopsy 1 because it contained regions annotated as lymphoid tissue and was also described extensively in the original paper. Biopsy 1 contains 293 spots covered by tissue. The third dataset is publicly available from the 10x Genomics website and includes matching Visium spatial gene expression (3,493 spots) and immunofluorescence staining of an endometrial adenocarcinoma of the ovary. The sample was stained with an anti-cytokeratin antibody, an anti-human CD45 antibody and DAPI. The fourth dataset is from an IDC sample prepared on the Visium platform (4,727 spots) and stained with an anti-human CD3 antibody and DAPI. The final dataset included data from ten human skin squamous cell carcinomas profiled on either the ST or the Visium platform⁴². Among the two samples run on the Visium platform, we chose to analyze that from patient 4 (P4) as the data quality was higher as shown in the original paper. Sample P4 contains 722 spots covered by tissue.

Preprocessing and dimension reduction.

In all datasets, raw gene expression counts were log transformed and normalized using library size^{45,46}. PCA was then performed on the top 2,000 most HVGs. Two thousand HVGs provided the best clustering performance in our benchmarks (Supplementary Fig. 21). In downstream analyses, we modeled the top 15 PCs from the Visium libraries, and we modeled the top seven PCs from the sample prepared on the ST platform (melanoma). The choice to model PCs rather than the full gene expression profile allows for a more tractable probabilistic model, avoiding the need for cumbersome multivariate discrete distributions. PCs are commonly used in clustering analysis of gene expression data. Here, we recommend modeling the top 15 PCs to capture as much of the variability in the data as possible while limiting the rapid increase in space that occurs with higher dimensions, although users may choose to model a different number of PCs or HVGs using the BayesSpace R package. Modeling more than 15 PCs did not provide substantial improvements in clustering performance but increased runtime and memory usage in our benchmarks (Supplementary Fig. 21). In the melanoma sample, many of the higher PCs exhibited higher numbers of extreme outliers (Supplementary Fig. 22) and significantly less variance, suggesting that they most likely represent technical variability. Because the older ST technology has lower coverage, sequencing depth and throughput, fewer PCs are necessary for modeling.

Spatial clustering model.

BayesSpace implements a fully Bayesian model with a Markov random field before encouraging spots of the same cluster to be close to one another. Such models have been widely used in image analysis, including analyses of microarray images^{20,21}. ST and Visium spots are arranged on square and hexagonal lattices, which provide a natural way to define a neighborhood structure (Fig. 1b). For each spot i , a low d -dimensional representation y_i (for example, PCs) of the gene expression vector can be obtained. We model the data as follows:

$$(y_i | z_i = k, w_i) \sim N(y_i; \mu_k, w_i^{-1} \Lambda^{-1})$$

$z_i \in \{1, \dots, q\}$ denotes the latent cluster that i belongs to, μ_k denotes the mean vector for cluster k , Λ denotes the precision matrix, and w_i denotes an unknown (observation-

specific) scaling factor. We assume a common (fixed) precision matrix across clusters because the number of unknown parameters in the precision matrix quickly rises with higher numbers of clusters and numbers of PCs modeled. In practice, we found that the variable precision model often required strong priors for parameter estimation. We also assume that the common precision matrix is unconstrained as there is correlation between PCs after conditioning on cluster, even though PCs are marginally uncorrelated (Supplementary Fig. 20). On real data, variable and independent precision models both performed poorly relative to the unconstrained, fixed precision model.

The number of clusters q is determined by prior biological knowledge when available or otherwise by the elbow of the pseudo-log-likelihood plot (Supplementary Figs. 9 and 10). We place the following priors on μ_k , Λ and w_i :

$$\begin{aligned}\mu_k &\stackrel{\text{i.i.d.}}{\sim} N(\mu_0, \Lambda_0^{-1}), \\ \Lambda &\stackrel{\text{i.i.d.}}{\sim} \text{Wishart}_d(\alpha, \text{diag}(\beta)_d^{-1}), \\ w_i &\stackrel{\text{i.i.d.}}{\sim} \Gamma\left(\frac{\nu}{2}, \frac{\nu}{2}\right),\end{aligned}$$

where μ_0 , Λ_0 , α and β are fixed hyperparameters. By default, we set μ_0 to be the empirical mean vector of the data, which is generally the zero vector for PCA input. Λ_0 is set to 0.01 times the identity matrix to provide a weak prior that will be dominated by the data when there are spots assigned to the cluster. Similarly, we set $\alpha = 1$ and $\beta = 0.01$ to provide a weak prior for the precision matrix. ν denotes a fixed degrees-of-freedom parameter to control the heaviness of tails and was set to $\nu = 4$, which was previously shown to overcome the influence of outlier spots during clustering²¹. We also assume that y_i and w_i are independent. As such, when marginalizing over w_i , our normal likelihood becomes a multivariable t distribution with a mean of 0 and covariance matrix $\frac{\nu}{\nu-2}\Lambda^{-1}$. This formulation allows us to use a simple Gibbs sampling for updating most of the parameters because the observations are normally distributed when conditioning on w_i . w_i values can also be interpreted as weights; the model will simply estimate a small weight value for any potential outlying data value. This provides robustness against outliers that can be commonly encountered in these types of data (Supplementary Fig. 2). Estimation of parameters is carried out using an MCMC method. We initialize z using a non-spatial clustering method such as `mclust` by default²³. Alternative initializations can also be supplied as a label vector. Next, iteratively and sequentially, each μ_k , Λ and w_i is updated via Gibbs sampling, and each z_i is updated via the Metropolis–Hastings algorithm. Specifically, each z_i is updated by taking into account both the likelihood and spatial prior information. The Markov random field prior is given by the Potts model:

$$\pi(z_i) = \exp\left(\frac{\gamma}{|\langle ij \rangle|} \times 2 \sum_{\langle ij \rangle} I(z_i = z_j)\right),$$

where $\langle ij \rangle$ denotes all spots j that are neighbors of i , I represents the indicator function, and γ is a fixed parameter controlling the strength of the smoothing. In this way, neighboring spots are encouraged to belong to the same cluster. Further details on the MCMC algorithm

are provided in the Supplementary Information. Model fitting diagnostics are provided in Supplementary Figs. 2 and 20.

Spatial clustering model at enhanced resolution.

To enhance the resolution of the clustering map, we segmented each spot into subspots and again leveraged spatial information using the Potts model spatial prior. Specifically, we segmented each ST spot into nine subspots and each Visium spot into six subspots (Fig. 1b). For ST, we used nine subspots to help increase the resolution of data from lower-resolution technology, because ST spots are 100 μm in diameter, while Visium spots are 55 μm in diameter. This translates into more than a threefold difference in area. In the IDC and OC samples, Visium spots are estimated to contain a median of around 20 cells; therefore, subspots will generally represent the expression of a few cells, rather than that of potentially dozens of cells at the spot level (Supplementary Figs. 7 and 8).

Relative to the spot-level clustering method, model specification and parameter estimation is largely similar for enhanced-resolution clustering, although the unit of analysis is now the subspot rather than the spot. As gene expression is not observed at the subspot level, it is modeled as another latent variable that is also estimated through MCMC. The latent expression of each subspot j that is part of spot i is denoted as γ_{ij}^* , initialized to be y_i and then updated via the Metropolis–Hastings algorithm. In each iteration and for each spot, the new proposal is given by $\gamma_{ij}^{*'} = \gamma_{ij}^* + \epsilon_{ij}$ for each subspot, such that the error $\epsilon_{ij} \sim \mathcal{N}(0, \sigma^2 I_d)$, where σ^2 is a small fixed parameter and $\sum_j \epsilon_{ij} = 0$. In effect, this jitters the latent expression value of each subspot within a spot while keeping the total expression of the spot fixed. The proposal is accepted or rejected based on the conditional likelihood of the proposal given the other parameters. We set σ^2 such that the acceptance rate ranges from 25% to 40% of iterations on average to maximize the efficiency of the Metropolis–Hastings algorithm⁴⁷. A weak Gaussian prior is placed on the latent expression to ensure that the jittered values do not drift too far away from y_i . Aside from replacing y_i with γ_{ij}^* , all other steps of the MCMC algorithm remain the same as in the spot-level clustering method. Model fitting diagnostics are provided in Supplementary Fig. 20. Intuitively, the enhancement procedure reassigns the total expression within a spot to its constituent subspots by leveraging spatial information, ultimately generating a higher-resolution spatial clustering map.

Mapping high-resolution PCs to high-resolution gene expression space.

While BayesSpace can provide higher-resolution maps of spatial transcriptomic patterns, the modeling is carried out on the PC space, and an additional step is necessary to map the PC values back to the original log-normalized gene expression space. BayesSpace implements two options for predicting high-resolution gene expression: linear regression and nonlinear regression using XGBoost (default)⁴⁸. In either case, a model is trained for each gene for which the outcome is the measured gene expression at the spot level and the predictors are the PCs generated from the original data. The fitted model can then be used to predict gene expression from the high-resolution PCs estimated using enhanced-resolution clustering. The enhanced gene expression values can be visualized spatially and analyzed via differential expression methods (Fig. 1a). In our analyses, we used the two-sided

Wilcoxon rank-sum test as implemented in Seurat to identify the top differentially expressed genes, and also we used Seurat for heatmap visualization of the centered and scaled gene expression values⁴⁹.

Simulations.

Using several simulations, we evaluated the performance of BayesSpace. The first simulation compared BayesSpace spot-level clustering to other non-spatial and spatial clustering methods: *k*-means, Louvain, mclust, SC3 and Giotto. We could not evaluate stLearn in simulation due to the need for an image as input. The simulated data were based on the melanoma and OC samples introduced in the earlier results. Eight replicates of simulated melanoma and OC PCs were generated from *t*-distributions with means, precision and spot labels determined by spot-level clustering results of the real melanoma and OC samples, respectively (Fig. 3b and Supplementary Fig. 23). Other clustering methods were implemented as described in the Supplementary Information with the true cluster number provided as input. BayesSpace was also implemented with the true cluster number provided as input. Performance was assessed using the ARI between ground-truth spot labels and clustering results.

In the second simulation, we evaluated the performance of BayesSpace subspot-level enhanced clustering. We simulated 20 replicates from *t*-distributions with means, precision and labels based on real melanoma and OC samples, but, unlike for the previous simulation, we generated subspots using the enhanced clustering results as the ground truth (Figs. 3c and 4d). The simulated subspot-level PCs were averaged to provide spot-level PCs that were given as input to BayesSpace. We can use the modal ground-truth label of the subspots within each spot to generate an optimal spot-level clustering for each dataset (Supplementary Fig. 23). The ARI between this optimal spot-level clustering and the subspot-level ground truth represents the highest ARI that can be achieved when all subspots within a spot must belong to the same cluster, as is the case with spot-level clustering.

In the third simulation, we sampled data from real single cells rather than simulating PCs. Here, we sampled single cells from scRNA-seq profiling of patients with high-grade serous OC (HGSOc)⁵⁰. The single cells can be sampled into subspots on the OC Visium sample, providing another way to evaluate the performance of BayesSpace clustering and enhancement relative to other methods without relying on model-based data generation. Given the limited number of single cells, we used only the positions from a portion of the OC Visium sample. Ground-truth cluster labels were derived from expert single-cell level annotation of tumor and stroma compartments within the immunofluorescence stain image associated with the OC sample. In each subspot, the ground truth was assigned using the modal annotation of the single cells located within the subspot. Consequently, the ground-truth assignment takes into account gaps between spots in spatial transcriptomic technologies, and the clusters represent realistic biological spatial domains.

To add complexity to the simulation, we separated the tumor compartment into two ground-truth clusters and introduced two additional intratumoral clusters that represent heterogeneity within tumors. Thus, the simulation included a total of five spatial ground-truth clusters, including the stroma compartment cluster. The single-cell sampling strategy is

shown in Supplementary Table 1, with single cells randomly drawn from single-cell clusters into corresponding spatial clusters in each of the eight simulation replicates. As raw counts were not available in the HGSOc dataset, pseudocounts were obtained by back transforming log-normalized counts, and simulated data were generated by aggregating across all subspots within a spot. The data were then processed to generate PCs as described for real data in the Methods. Because HGSOc single-cell clusters are very well separated, we also added random noise to each simulated PC equal to 25% of its variance, thus adding additional complexity to our simulation. This process also made our simulation more realistic when comparing the generated PCs to PCs derived from experimental data (Supplementary Fig. 22).

Reporting Summary.

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

Data availability

Datasets analyzed in this paper are available in raw form from their original authors (see details in the Supplementary Note), and the SingleCellExperiment objects that we prepared for analysis with BayesSpace are available through the BayesSpace package. Raw count matrices, images and spatial data from the IDC sample are accessible on the 10x Genomics website at <https://support.10xgenomics.com/spatial-gene-expression/datasets>.

Code availability

BayesSpace is available as a Bioconductor package at <http://www.bioconductor.org/packages/release/bioc/html/BayesSpace.html>, and the source code is publicly available at <https://github.com/edward130603/BayesSpace>.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgements

This research was supported by funding from the National Institutes of Health (P01-CA225517, P30-CA015704 to R.G. and P.N.; T32-CA080416, F30-CA254168 to T.P.), the Immunotherapy and Data Science Integrated Research Centers at Fred Hutchinson to E.Z., M.R.S., X.R. and J.H.B. and the Scientific Computing Infrastructure at Fred Hutchinson funded by ORIP grant S10OD028685. We thank M. Lin and P.L. Porter for their pathological review of J.G.'s histological annotations, K.J. Cheung from the Fred Hutchinson Public Health Sciences and Human Biology Divisions for his suggestions in our analysis of the IDC sample, A. Moshiri from the UW Division of Dermatology for his review of T.P.'s histopathological annotations and Q. Nguyen and X. Tan at the University of Queensland for their assistance in applying stLearn.

Competing interests

R.G. has received consulting income from Juno Therapeutics, Takeda, Infotech Soft, Celgene and Merck, has received research support from Janssen Pharmaceuticals and Juno Therapeutics and declares ownership in Ozette Technologies and stock ownership in 10x Genomics. S.R.W., C.R.U. and S.E.B.T. are employees of and hold shares in 10x Genomics. All other authors declare no conflicts of interest.

References

1. Ståhl PL et al. Visualization and analysis of gene expression in tissue sections by spatial transcriptomics. *Science* 353, 78–82 (2016). [PubMed: 27365449]
2. Thrane K, Eriksson H, Maaskola J, Hansson J & Lundeberg J Spatially resolved transcriptomics enables dissection of genetic heterogeneity in stage III cutaneous malignant melanoma. *Cancer Res.* 78, 5970–5979 (2018). [PubMed: 30154148]
3. Berglund E et al. Spatial maps of prostate cancer transcriptomes reveal an unexplored landscape of heterogeneity. *Nat. Commun.* 9, 2419 (2018). [PubMed: 29925878]
4. Maynard KR et al. Transcriptome-scale spatial gene expression in the human dorsolateral prefrontal cortex. *Nat. Neurosci.* 24, 425–436 (2021). [PubMed: 33558695]
5. Janosevic D et al. The orchestrated cellular and molecular responses of the kidney to endotoxin define a precise sepsis timeline. *eLife* 10, e62270 (2021). [PubMed: 33448928]
6. Saiselet M et al. Transcriptional output, cell types densities and normalization in spatial transcriptomics. *J. Mol. Cell Biol.* 12, 906–908 (2020). [PubMed: 32573704]
7. Lubeck E, Coskun AF, Zhiyentayev T, Ahmad M & Cai L Single-cell in situ RNA profiling by sequential hybridization. *Nat. Methods* 11, 360–361 (2014). [PubMed: 24681720]
8. Chen KH, Boettiger AN, Moffitt JR, Wang S & Zhuang X Spatially resolved, highly multiplexed RNA profiling in single cells. *Science* 348, aaa6090 (2015). [PubMed: 25858977]
9. Rodriques SG et al. Slide-seq: a scalable technology for measuring genome-wide expression at high spatial resolution. *Science* 363, 1463–1467 (2019). [PubMed: 30923225]
10. Hu KH et al. ZipSeq: barcoding for real-time mapping of single cell transcriptomes. *Nat. Methods* 17, 833–843 (2020). [PubMed: 32632238]
11. Gavin J & Jennison C A subpixel image restoration algorithm. *J. Comput. Graph. Stat.* 6, 182–201 (1997).
12. Ripley BD The use of spatial models as image priors. In *Spatial Statistics and Imaging: Papers from the Research Conference on Image Analysis and Spatial Statistics held at Bowdoin College, Brunswick, Maine, Summer 1988* 20, 309–340 (Institute of Mathematical Statistics, 1991).
13. Tipping ME & Bishop CM Bayesian image super-resolution. In *Proc. 15th Int. Conf. Neural Information Processing Systems* (eds. Becker S, Thrun S & Obermayer K) 1303–1310 (2002).
14. Andersson A et al. Single-cell and spatial transcriptomics enables probabilistic inference of cell type topography. *Commun. Biol.* 3, 565 (2020). [PubMed: 33037292]
15. Cable DM et al. Robust decomposition of cell type mixtures in spatial transcriptomics. *Nat. Biotechnol.* 10.1038/s41587-021-00830-w (2021).
16. Elosua-Bayes M, Nieto P, Mereu E, Gut I & Heyn H SPOTlight: seeded NMF regression to deconvolute spatial transcriptomics spots with single-cell transcriptomes. *Nucleic Acids Res.* 10.1093/nar/gkab043 (2021).
17. Zhu Q, Shah S, Dries R, Cai L & Yuan GC Identification of spatially associated subpopulations by combining scRNAseq and sequential fluorescence in situ hybridization data. *Nat. Biotechnol.* 36, 1183–1190 (2018).
18. Dries R et al. Giotto: a toolbox for integrative analysis and visualization of spatial expression data. *Genome Biol.* 22, 78 (2021). [PubMed: 33685491]
19. Pham DT et al. stLearn: integrating spatial location, tissue morphology and gene expression to find cell types, cell–cell interactions and spatial trajectories within undissociated tissues. Preprint at *bioRxiv* 10.1101/2020.05.31.125658 (2020).
20. Besag J On the statistical analysis of dirty pictures. *J. R. Stat. Soc. Ser. B* 48, 259–279 (1986).
21. Gottardo R, Besag J, Stephens M & Murua A Probabilistic segmentation and intensity estimation for microarray images. *Biostatistics* 7, 85–99 (2006). [PubMed: 16049139]
22. Amezquita RA et al. Orchestrating single-cell analysis with Bioconductor. *Nat. Methods* 17, 137–145 (2020). [PubMed: 31792435]
23. Fraley C, Raftery AE & Murphy TB mclust version 4 for R: normal mixture modeling for model-based clustering, classification, and density estimation. *R. J* 8, 289–317 (2012).

24. Blondel VD, Guillaume JL, Lambiotte R & Lefebvre E Fast unfolding of communities in large networks. *J. Stat. Mech. Theory Exp.* 2008, P10008 (2008).
25. Kiselev VY et al. SC3: consensus clustering of single-cell RNA-seq data. *Nat. Methods* 14, 483–486 (2017). [PubMed: 28346451]
26. Wang HX et al. HSPB1 deficiency sensitizes melanoma cells to hyperthermia induced cell death. *Oncotarget* 7, 67449–67462 (2016). [PubMed: 27626679]
27. Mathieu V et al. The sodium pump α 1 sub-unit: a disease progression-related target for metastatic melanoma treatment. *J. Cell. Mol. Med.* 13, 3960–3972 (2009). [PubMed: 19243476]
28. Tirosh I et al. Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq. *Science* 352, 189–196 (2016). [PubMed: 27124452]
29. Mori K et al. CpG hypermethylation of collagen type I α 2 contributes to proliferation and migration activity of human bladder cancer. *Int. J. Oncol.* 34, 1593–1602 (2009). [PubMed: 19424577]
30. Knudsen ES et al. Progression of ductal carcinoma in situ to invasive breast cancer is associated with gene expression programs of EMT and myoepithelia. *Breast Cancer Res. Treat.* 133, 1009–1024 (2012). [PubMed: 22134623]
31. Lee S et al. Differentially expressed genes regulating the progression of ductal carcinoma in situ to invasive breast cancer. *Cancer Res.* 72, 4574–4586 (2012). [PubMed: 22751464]
32. Hu M et al. Regulation of in situ to invasive breast carcinoma transition. *Cancer Cell* 13, 394–406 (2008). [PubMed: 18455123]
33. Hattstrup CL & Gendler SJ MUC1 alters oncogenic events and transcription in human breast cancer cells. *Breast Cancer Res.* 8, R37 (2006). [PubMed: 16846534]
34. Besmer DM et al. Pancreatic ductal adenocarcinoma mice lacking mucin 1 have a profound defect in tumor growth and metastasis. *Cancer Res.* 71, 4432–4442 (2011). [PubMed: 21558393]
35. Behrens ME et al. The reactive tumor microenvironment: MUC1 signaling directly reprograms transcription of CTGF. *Oncogene* 29, 5667–5677 (2010). [PubMed: 20697347]
36. Zhang X et al. Insights into the distinct roles of MMP-11 in tumor biology and future therapeutics (Review). *Int. J. Oncol.* 48, 1783–1793 (2016). [PubMed: 26892540]
37. Holland DG et al. *ZNF703* is a common luminal B breast cancer oncogene that differentially regulates luminal and basal progenitors in human mammary epithelium. *EMBO Mol. Med.* 3, 167–180 (2011). [PubMed: 21337521]
38. Sircoulomb F et al. *ZNF703* gene amplification at 8p12 specifies luminal B breast cancer. *EMBO Mol. Med.* 3, 153–166 (2011). [PubMed: 21328542]
39. Daly R, Binder M & Sutherland R Overexpression of the *Grb2* gene in human breast cancer cell lines. *Oncogene* 9, 2723–2727 (1994). [PubMed: 8058337]
40. Tari AM, Hung MC, Li K & Lopez-Berestein G Growth inhibition of breast cancer cells by Grb2 downregulation is correlated with inactivation of mitogen-activated protein kinase in EGFR, but not in ErbB2, cells. *Oncogene* 18, 1325–1332 (1999). [PubMed: 10022814]
41. Onichtchouk D et al. Silencing of TGF-signalling by the pseudoreceptor BAMBI. *Nature* 401, 480–485 (1999). [PubMed: 10519551]
42. Ji AL et al. Multimodal analysis of composition and spatial architecture in human squamous cell carcinoma. *Cell* 182, 497–514 (2020). [PubMed: 32579974]
43. Sümer C, Boz Er AB & Dinçer T Keratin 14 is a novel interaction partner of keratinocyte differentiation regulator: receptor-interacting protein kinase 4. *Turk. J. Biol.* 43, 225–234 (2019). [PubMed: 31582880]
44. Liu W Unsupervised learning approaches for the finite mixture models: EM versus MCMC. In 2010 International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery 2010, 498–501 (IEEE, 2010).
45. Lun ATL, Bach K & Marioni JC Pooling across cells to normalize single-cell RNA sequencing data with many zero counts. *Genome Biol.* 17, 75 (2016). [PubMed: 27122128]
46. McCarthy DJ, Campbell KR, Lun ATL & Wills QF Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics* 33, 1179–1186 (2017). [PubMed: 28088763]

47. Gelman A, Roberts GO & Gilks WR Efficient Metropolis jumping rules. *Bayesian Stat.* 5, 599–607 (1996).
48. Chen T & Guestrin C XGBoost: a scalable tree boosting system. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* 785–794 (Association for Computing Machinery, 2016).
49. Stuart T et al. Comprehensive integration of single-cell data. *Cell* 177, 1888–1902 (2019). [PubMed: 31178118]
50. Izar B et al. A single-cell landscape of high-grade serous ovarian cancer. *Nat. Med.* 26, 1271–1279 (2020). [PubMed: 32572264]

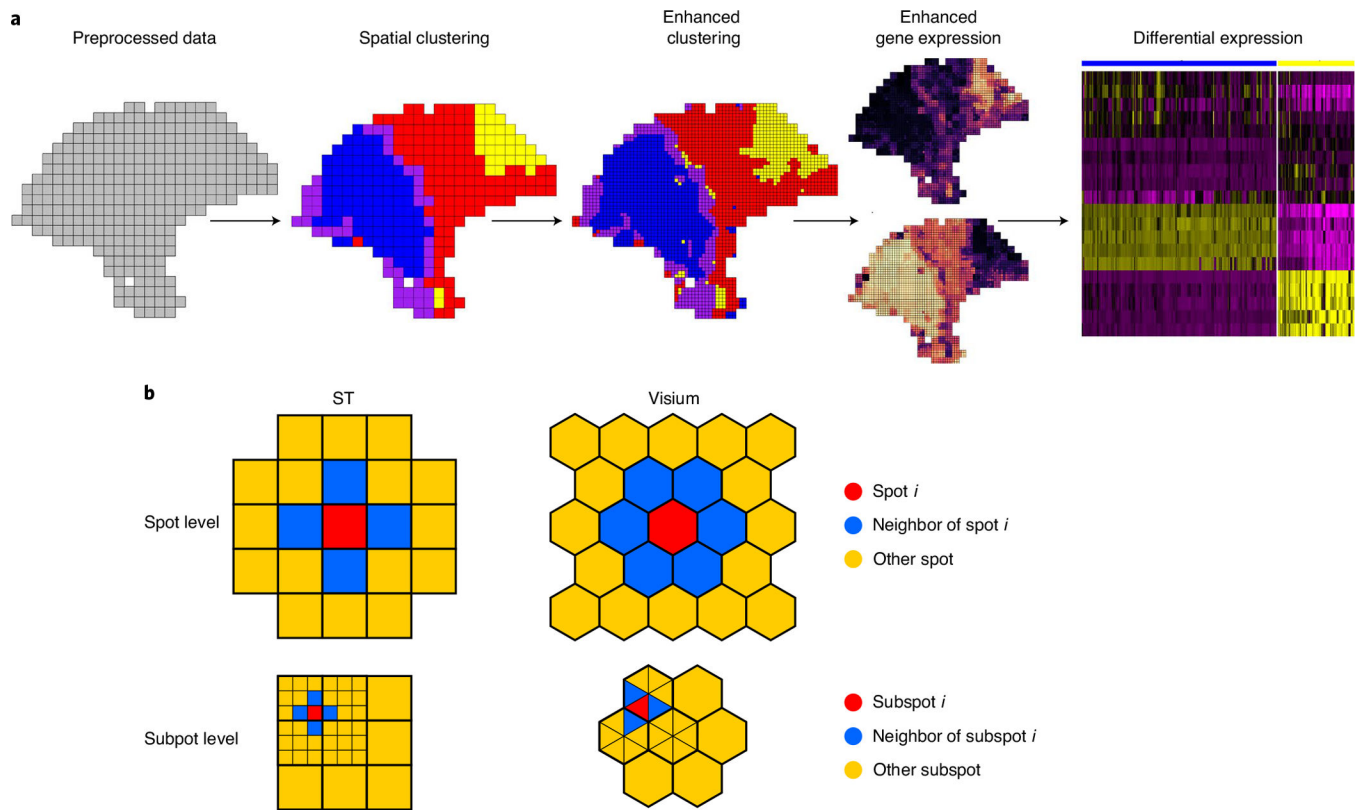


Fig. 1. The BayesSpace workflow.

a. The BayesSpace workflow begins with preprocessed ST or Visium data. Data are spatially clustered to infer regions with similar expression profiles. These clusters can be refined via enhanced clustering to provide a higher-resolution spatial map. Enhanced clustering also provides the basis for predicting gene expression at the higher resolution, which can be used in further differential expression analyses. **b.** From geometric representations of spatial distribution of spots in the ST and Visium technologies, neighbors can be identified for each spot based on shared edges (top). Each spot can be subdivided into subspots, which again have natural edge-based neighbors (bottom).

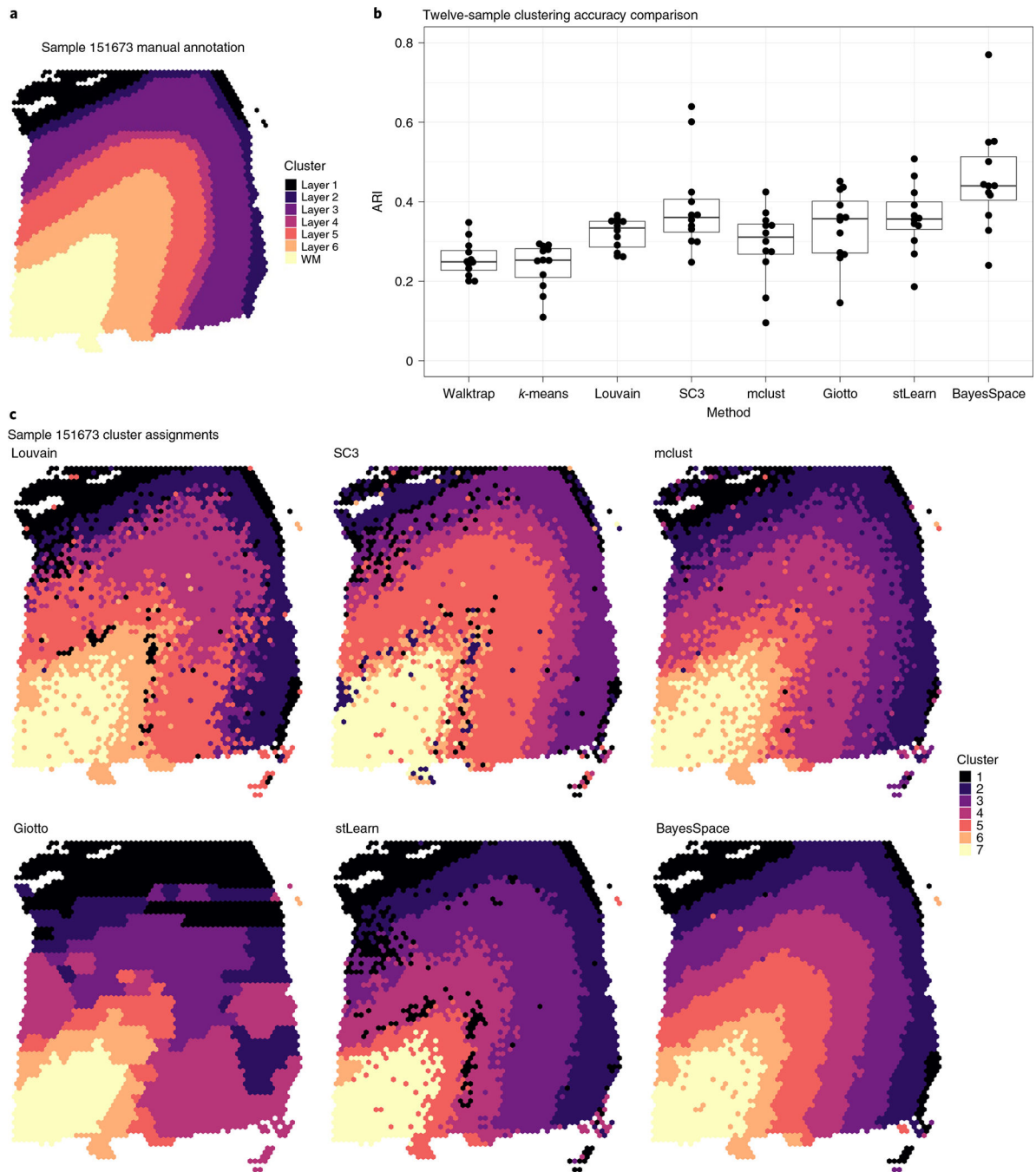


Fig. 2. BayesSpace improves computational resolution of layers in the DLPCF.

a. Ground truth. We highlight the manually annotated six DLPCF layers and white matter (WM) in sample 151673 from the spatialLIBD dataset. Annotated layers for the remaining samples can be found in the original publication⁴. **b.** Summary of clustering accuracy in all twelve samples. The ARI is used to compare similarity between cluster labels from each method against the manually annotated layers for all twelve samples. In the boxplot, the center line, box limits and whiskers denote the median, upper and lower quartiles and 1.5×

interquartile range, respectively. **c**, Cluster assignments generated by non-spatial (top) and spatial (bottom) methods for sample 151673.

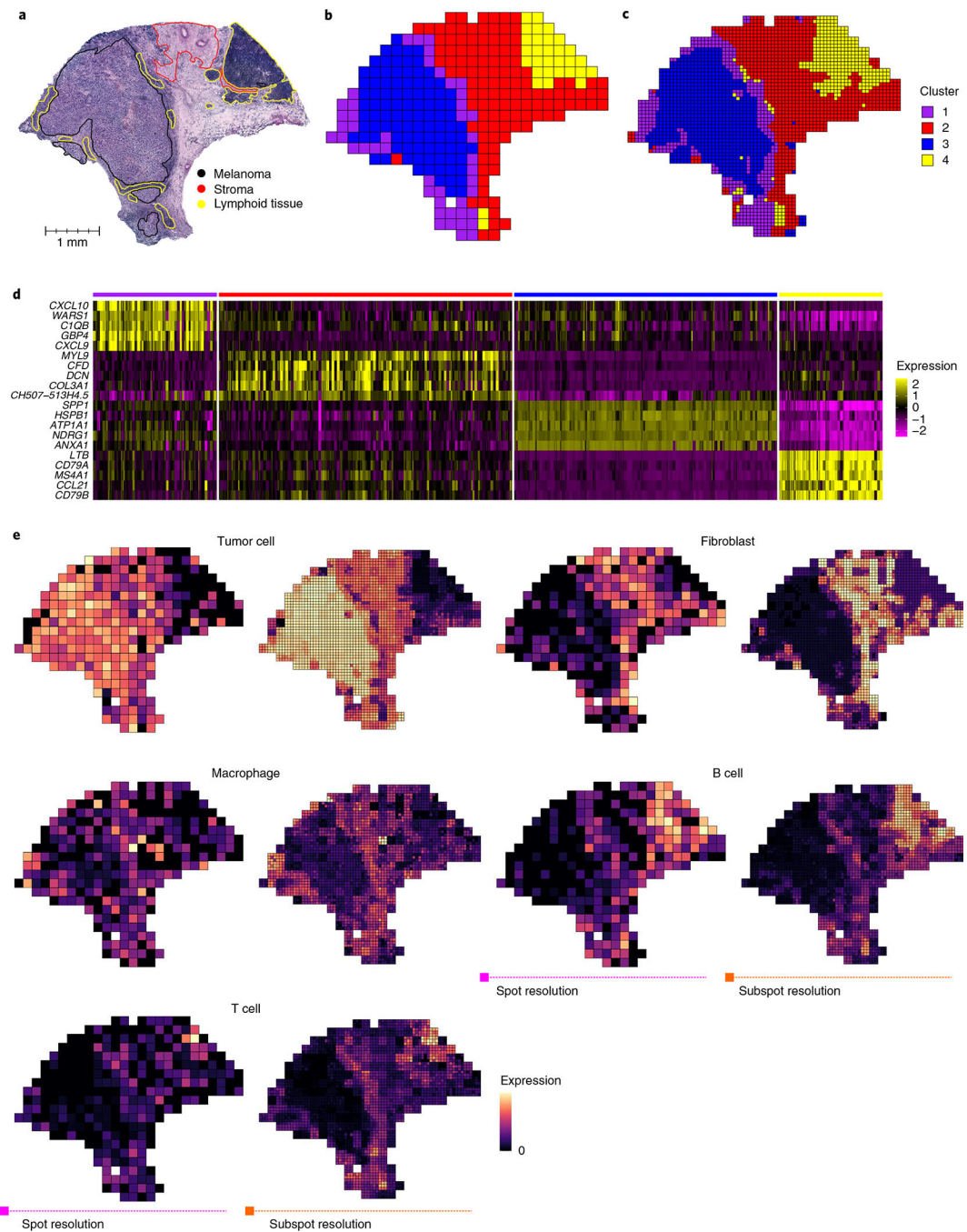


Fig. 3. Enhanced-resolution clustering identifies tumor-proximal lymphoid tissue in a melanoma sample.

a, The original histopathological annotations of H&E-stained tissue ($N = 1$ tissue section, $n = 293$ spots) revealed a section of melanoma (black) adjacent to tumor-proximal lymphoid tissue (yellow) and a region of stroma (red), separating these from a larger section of tumor-distal lymphoid tissue (yellow)². Adapted from ref. ² with permission from the American Association for Cancer Research. Spatial clustering (**b**) and enhancement (**c**) generate biologically meaningful spatial domains corresponding to the original annotations.

Enhanced-resolution clustering identified tumor-proximal lymphoid tissue (cluster 4, yellow), which was not resolved at spot-level clustering. **d**, Differential expression analysis between the four clusters highlighted spatial differences in the expression of immune genes, cancer markers and genes encoding extracellular matrix proteins. **e**, For each of the five major cell types, the observed total spot-level expression (as measured by the summed log-normalized counts) of the defined marker genes (left) is shown alongside the corresponding enhanced-resolution expression (right). We show spatial expression plots for tumor cells (*PMEL*), fibroblasts (*COL1A1*), macrophages (*CD14*, *FCGR1A*, *FCGR1B*), B cells (*CD19*, *MS4A1*) and T cells (*CD2*, *CD3D*, *CD3E*, *CD3G*, *CD7*).

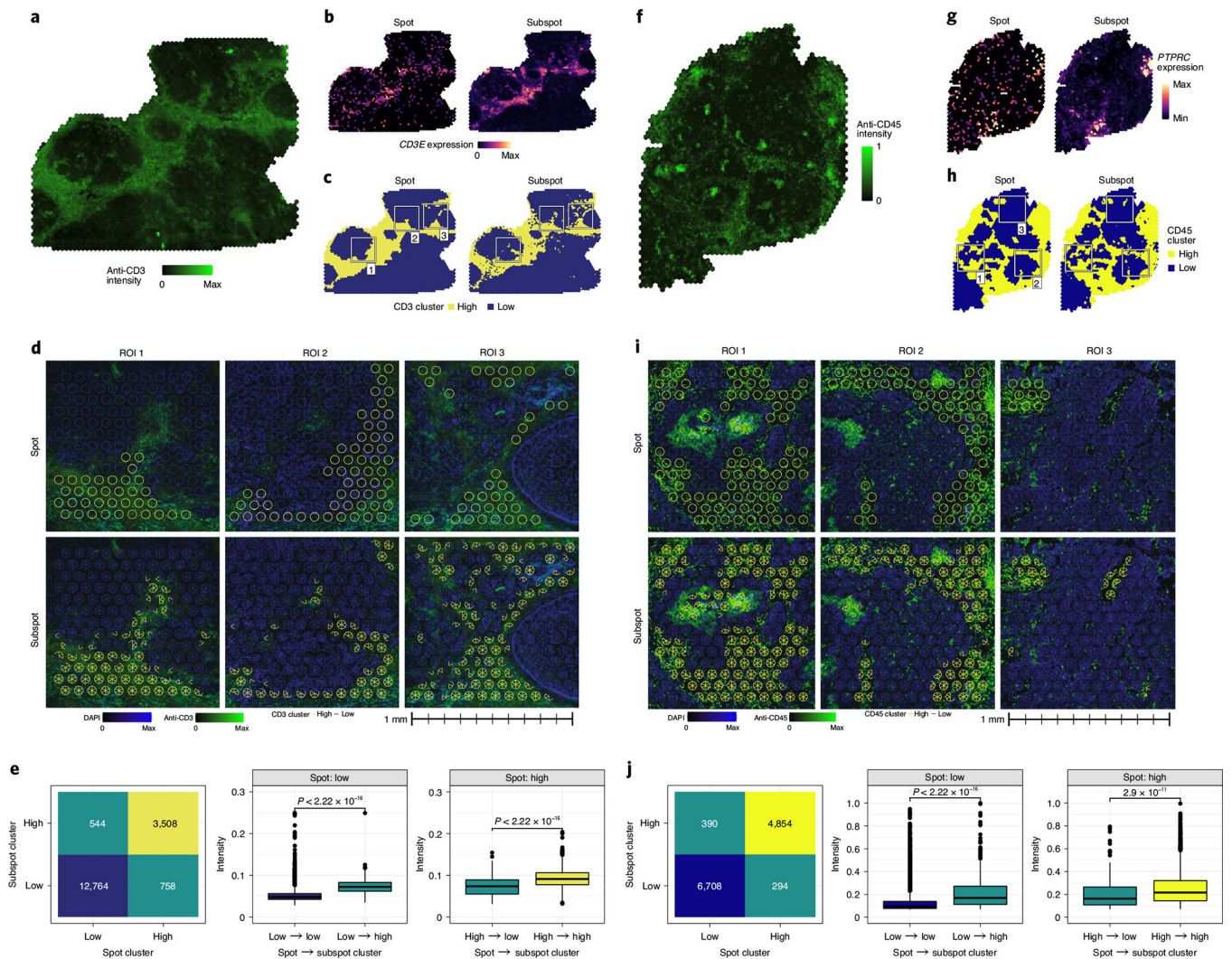


Fig. 4. Immunohistochemistry validates BayesSpace enhancement in an IDC sample and an OC sample.

a, Average intensity of the anti-CD3 immunofluorescent stain in the IDC. Intensity was scaled to the range (0, 1) for visualization. **b**, Log-normalized gene expression of *CD3E* measured on the Visium platform (left, 'spot') and enhanced with BayesSpace (right, 'subspot'). **c**, Dichotomized clustering of Visium gene expression values. After clustering the tissue section into ten clusters, the clusters were binned by their median anti-CD3 stain intensity into CD3 'high' and CD3 'low' clusters, shown here. White squares outline three ROI where the enhanced clustering revealed areas of increased heterogeneity. **d**, Zoomed-in views of the $n = 3$ ROI. Each panel shows a 1-mm² area of the immunofluorescence image. DAPI intensity is shown in blue, and anti-CD3 intensity is shown in green. Overlaid on each panel in the top row is the spot-level clustering. Each circle corresponds to the position and size (55- μ m diameter) of a spot on the Visium array and is colored based on whether it belongs to a CD3 'high' (yellow) or CD3 'low' (blue) cluster. The bottom row contains a similar overlay of the enhanced-resolution subspot clustering, where the circles are now subdivided into six wedges corresponding to the positions of subspots in the BayesSpace

model. As in the spot overlay, the subspots are colored based on their cluster membership. **e**, Summary of subspot reassignment after enhancement. On the left, we show a contingency table describing the number of subspots ($n = 17,574$) that belong to a CD3 'high' or 'low' cluster at the spot level and at the subspot level. Using two-sided Wilcoxon rank-sum tests, we also show that anti-CD3 intensity in subspots that are reassigned to a 'high' cluster is significantly higher ($P < 2.22 \times 10^{-16}$) than that in those that remain in a 'low' cluster (center) and that subspots that are reassigned to a 'low' cluster have a significantly lower ($P < 2.22 \times 10^{-16}$) anti-CD3 intensity than that in those that remain in a 'high' cluster (right). **f–j**, Panels for the OC mirror those for the IDC, with anti-CD45 intensity replacing anti-CD3 intensity and *PTPRC* (CD45) gene expression replacing that of *CD3E*. In **e**, we show $n = 12,246$ subspots. In **i**, we show $n = 3$ ROI. In **j**, using two-sided Wilcoxon rank-sum tests, we show that anti-CD45 intensity in subspots that are reassigned to a 'high' cluster is significantly higher ($P < 2.22 \times 10^{-16}$) than that in those that remain in a 'low' cluster (center) and that subspots that are reassigned to a 'low' cluster have a significantly lower ($P = 2.9 \times 10^{-11}$) anti-CD45 intensity than that in those that remain in a 'high' cluster (right). All reported P values are unadjusted values.

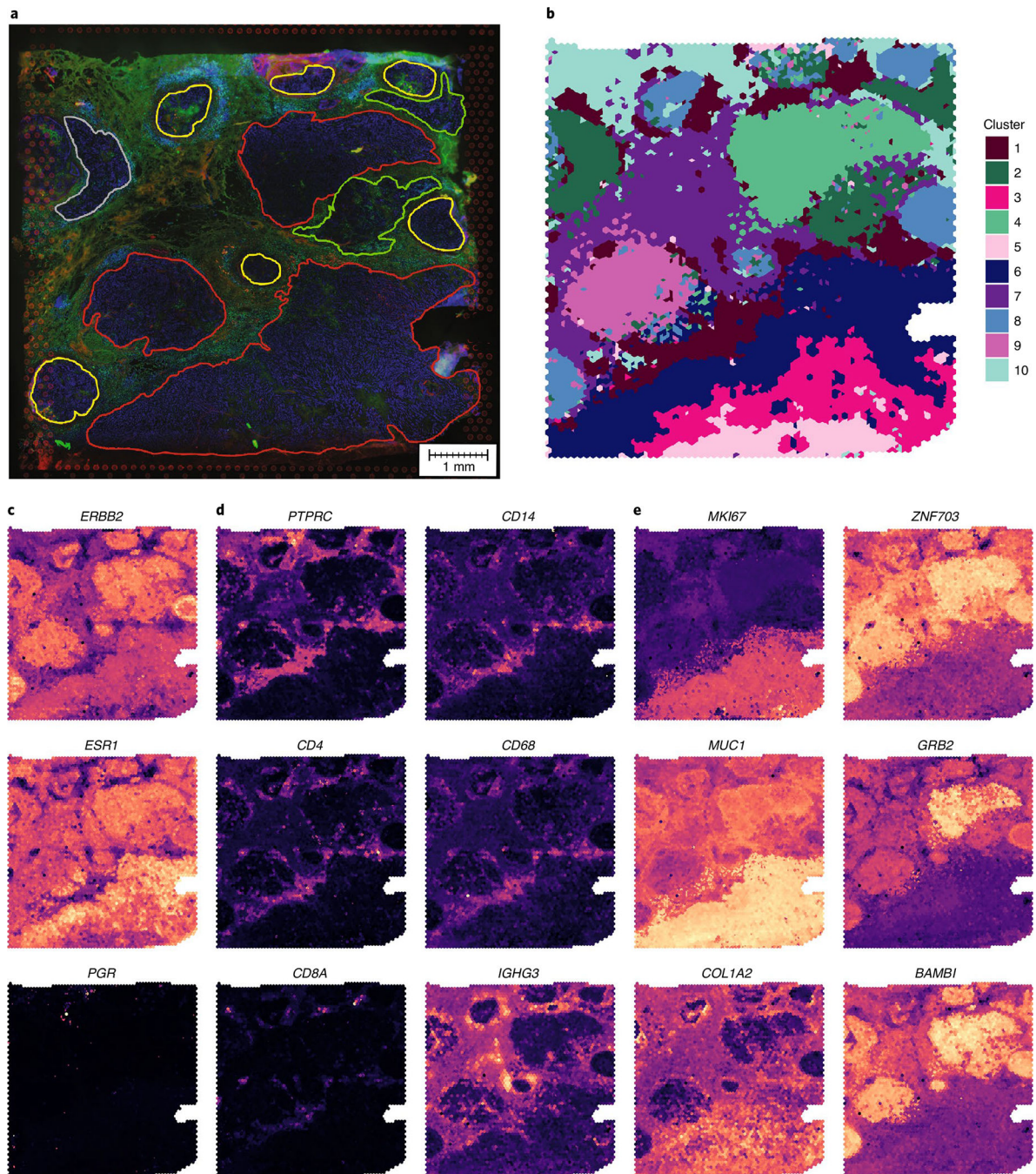


Fig. 5. BayesSpace identifies transcriptional heterogeneity within an IDC.

a, Immunofluorescent imaging of the tissue section ($N = 1$ tissue section, $n = 4,727$ spots) and histopathological annotations. DAPI intensity is shown in blue, anti-CD3 intensity is shown in green, and the Visium fiducial frame is shown in red. Annotated regions of IC are outlined in red, those of carcinoma in situ are outlined in yellow, those of benign hyperplasia are outlined in green, and those of unclassified tumor are outlined in gray. **b**, Enhanced BayesSpace clustering. **c**, Spatial expression of genes coding for HER2 (*ERBB2*) and ER (*ESR1*) and PR (*PGR*). **d**, Spatial expression of immune genes *PTPRC* (CD45),

CD4, *CD8A*, *CD14*, *CD68* and *IGHG3*. e. Spatial expression of proliferation marker *MKI67* (Ki-67), markers of tumor progression *MUC1* and *COL1A2*, the oncogene *ZNF703*, *GRB2* (coding for the growth factor receptor protein) and *BAMBI* (coding for transforming growth factor (TGF)- β pseudoreceptor).

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

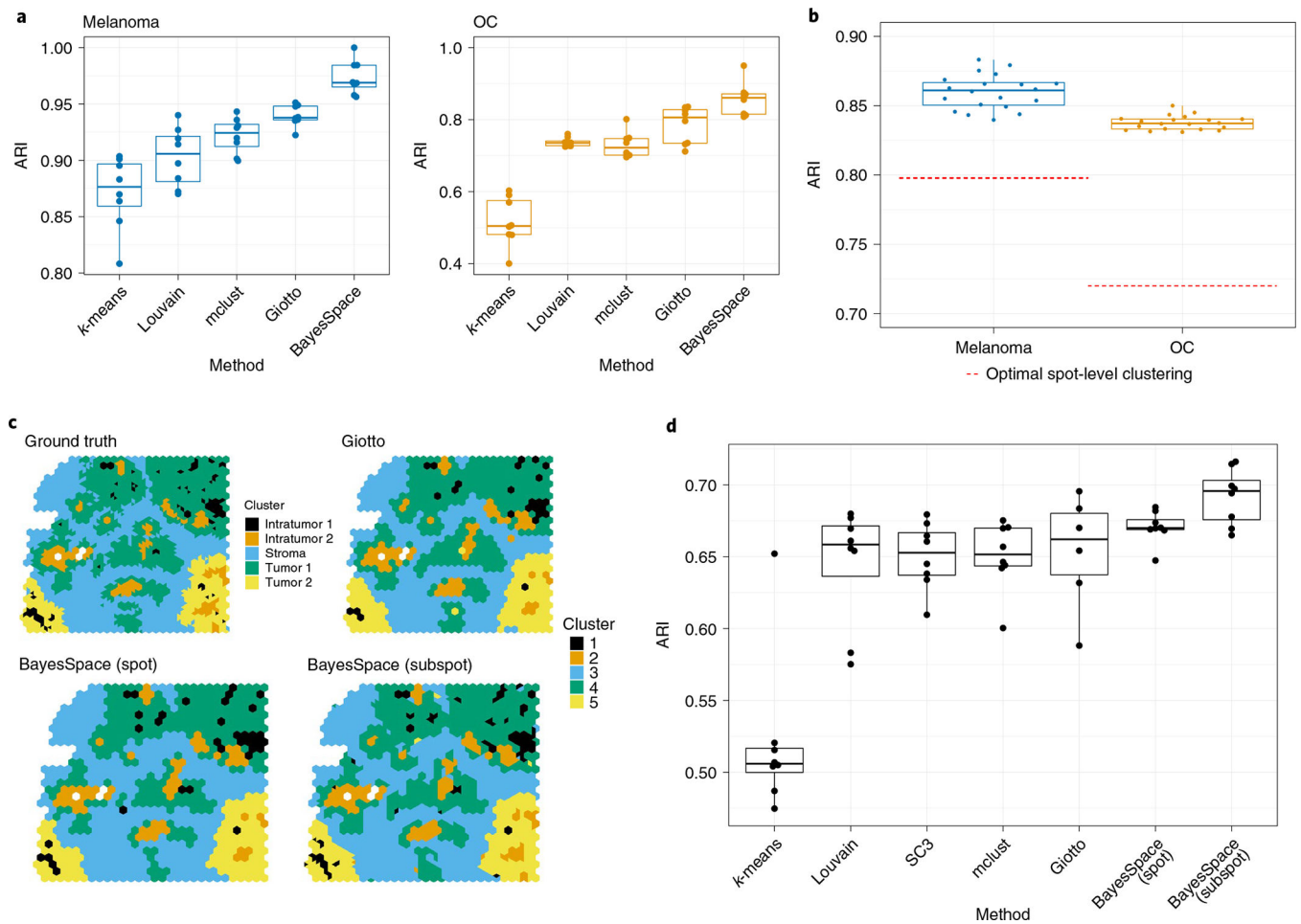


Fig. 6. BayesSpace outperforms spatial and non-spatial clustering methods with simulated data. **a**, In $N = 8$ replicates simulated from the melanoma sample and $N = 8$ replicates simulated from the OC sample, BayesSpace spot-level clustering outperforms other clustering methods. **b**, In $N = 20$ replicates for the simulation performed at the subspot level, BayesSpace enhanced clustering outperforms the optimal spot-level clustering (red dotted line). **c**, In the third simulation using single-cell data, the ground truth is derived from expert annotation of an immunofluorescence staining image corresponding to the OC sample (top left). Examples of clustering partitions generated by BayesSpace at the spot and subspot levels as well as by the next best method (Giotto) are also shown. **d**, BayesSpace clustering at the spot level slightly outperforms competing methods, while BayesSpace enhancement to the subspot level generally provides substantially higher performance than that of other methods in recapturing ground-truth clusters among the $N = 8$ simulation replicates. In all boxplots, the center line, box limits and whiskers denote the median, upper and lower quartiles and $1.5\times$ interquartile range, respectively.