

RESEARCH ARTICLE

Open Access



CoSTA: unsupervised convolutional neural network learning for spatial transcriptomics analysis

Yang Xu¹ and Rachel Patton McCord^{2*}

*Correspondence:
rmccord@utk.edu

² Department
of Biochemistry & Cellular
and Molecular Biology,
University of Tennessee,
Knoxville, TN, USA
Full list of author information
is available at the end of the
article

Abstract

Background: The rise of spatial transcriptomics technologies is leading to new insights about how gene regulation happens in a spatial context. Determining which genes are expressed in similar spatial patterns can reveal gene regulatory relationships across cell types in a tissue. However, many current analysis methods **do not take full advantage of the spatial organization of the data**, instead treating pixels as independent features. Here, we present CoSTA: a novel approach to learn spatial similarities between gene expression matrices via convolutional neural network (ConvNet) clustering.

Results: By analyzing simulated and previously published spatial transcriptomics data, we demonstrate that CoSTA learns spatial relationships between genes in a way that emphasizes broader spatial patterns rather than pixel-level correlation. CoSTA provides a **quantitative** measure of expression pattern similarity between each pair of genes rather than only classifying genes into categories. We find that CoSTA identifies narrower, but biologically relevant, sets of significantly related genes as compared to other approaches.

Conclusions: The deep learning CoSTA approach provides a different angle to spatial transcriptomics analysis by **focusing on the shape of expression patterns, using more information about the positions of neighboring pixels** than would an overlap or pixel correlation approach. CoSTA can be applied to any spatial transcriptomics data represented in matrix form and may have future applications to datasets such as histology in which images of different genes are from similar but not identical biological sections.

Keywords: Spatial transcriptomics, Gene clustering, Convolutional neural network

Background

Spatial transcriptomics has recently gained extensive attention from the scientific community. Different technologies have enabled high resolution measurements of how gene regulation is spatially organized across a tissue or thousands of single cells [1]. Analyses of these data have the potential to reveal spatial regulatory relationships between genes. However, current analysis pipelines often treat each pixel in an expression matrix as an independent feature, thus losing spatial information. For example, the seqFISH+ technique



© The Author(s). 2021. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

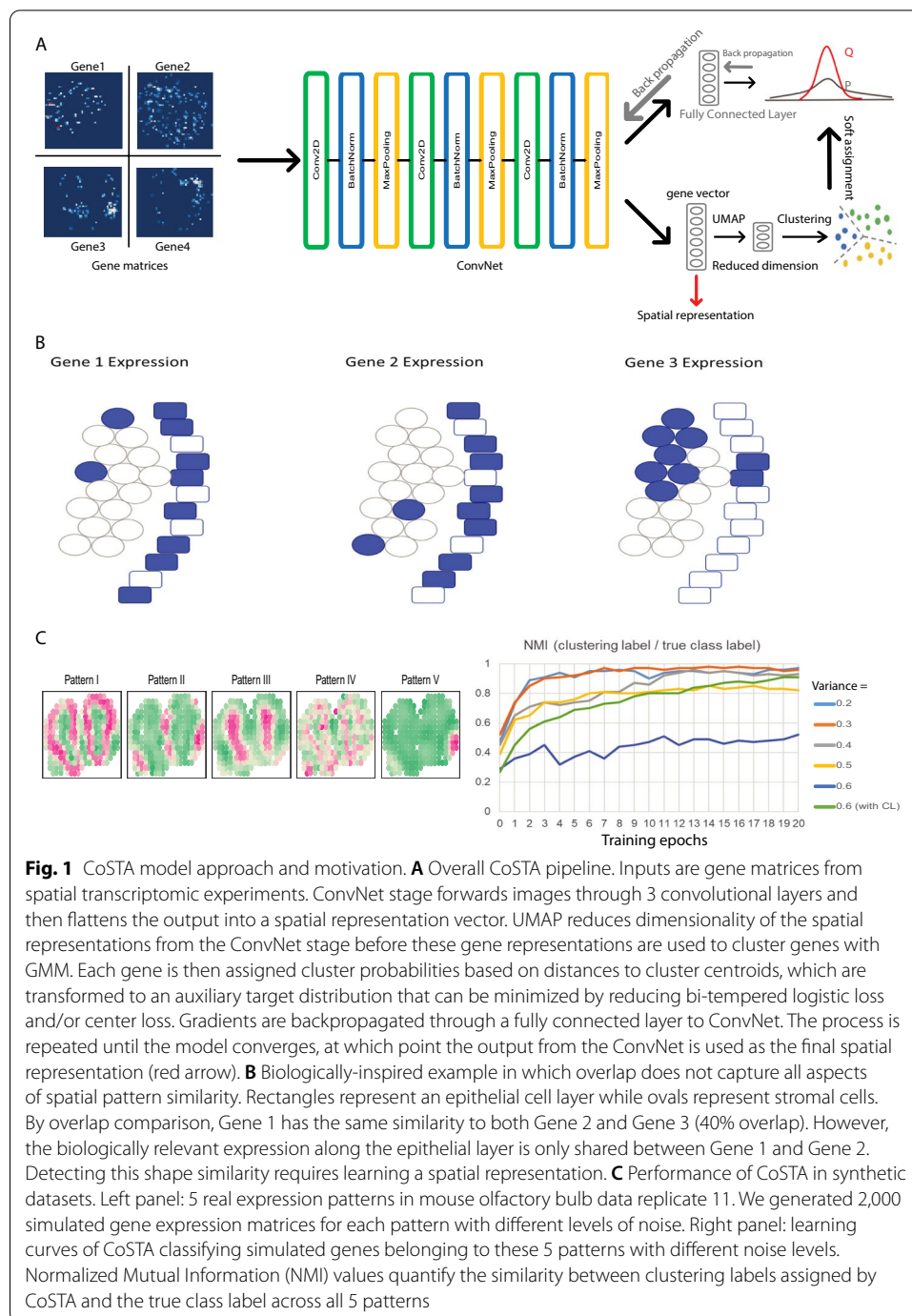
can fluorescently detect 10,000 mRNAs in situ at single cell resolution, and there are often groups of cells that have correlated gene expression with their neighbors to make up larger structures. However, the original report analyzed these expression patterns using PCA and hierarchical clustering, treating each cell as an independent feature, rather than preserving spatial positions of cell neighbors [2]. Slide-seq similarly produces high-throughput spatially resolved transcription information, using sequencing rather than fluorescence. Previous analyses of Slide-seq data first identified spatially non-random gene expression, but then looked for genes expressed in similar patterns using pixel-level overlap analysis rather than according to spatial features [3]. Existing algorithms for analysis of spatial transcriptomics are based on statistical modeling and primarily propose to distinguish spatially expressing or variable (SE or SV) genes from random spatial expression noise. For example, both SpatialDE and SPARK analysis approaches estimate how significant the spatial pattern of a gene is [4, 5]. SpatialDE further builds in an unsupervised pattern detection algorithm to cluster significant SE genes into different groups which have certain spatial patterns in collective. SPARK, in contrast, was designed only for finding SE genes. To examine spatial relationships between genes, this method still relies on hierarchical clustering that uses individual pixels as features. Therefore, even though SPARK can identify genes with significant spatial patterns, the latter part of the SPARK analysis decouples the expression from its original spatial context. Thus far, existing spatial transcriptomics analyses involve either multi-step complex feature engineering for spatial quantification or human-imposed rigid or statistical modeling-based screening of candidate SE genes. In the existing methods, the similarity of expression pattern between two genes is either binary- whether or not the genes cluster together- or is quantified based on pixel-level correlation.

In this work, we propose an approach inspired by computer vision and image classification to find relationships between spatial expression patterns of different genes while preserving the full spatial context (Fig. 1a). Our goal is to find quantitative comparisons between gene expression patterns in a way that preserves spatial relationships between neighboring cells and tissue regions. We aim for a method that will recognize an overall similar shape of expression, even if certain sets of pixels are not exactly overlapping. This is conceptually similar to image recognition in computer vision tasks. The use of convolutional neural networks brought a success of deep learning in computer vision and have demonstrated a wide range of applications, including image classification and object recognition. A few groups have proposed different approaches to use convolutional neural networks (ConvNet) in unsupervised learning [6–8]. Thus, here, we adopt an unsupervised ConvNet learning strategy for Spatial Transcriptomics Analysis (CoSTA). With simulated data, we show that CoSTA can correctly classify a variety of different spatial patterns and that the patterns CoSTA is detecting depend on spatial groupings rather than individual pixels. Then, we apply CoSTA to published MERFISH and Slide-seq data and show that CoSTA sometimes identifies smaller sets of genes with significant spatial relationships, but these identified relationships are biologically relevant.

Results

CoSTA architecture: training a ConvNet with pseudo-labels generated by GMM clustering

Though there are many unsupervised learning strategies, we chose to apply the workflow of DeepCluster, because it is straightforward and easy to implement [6]. Our CoSTA



approach consists of two main parts: clustering by Gaussian mixture model (GMM) and weight updating as commonly performed in training neural networks (Fig. 1a and see [Methods](#) for detailed description). Our inputs are sets of gene expression images, where each image is the matrix recording the expression levels of one gene at each position in space and all images belong to the same biological space. We first initialize a ConvNet randomly and then forward these gene expression matrices through the ConvNet. Our ConvNet consists of three convolutional layers, and each convolutional layer is followed

by a batch normalization layer and a max pooling layer. We flatten the matrix output from the last max pooling layer into a vector that captures the spatial features of the gene expression data. The size of this vector will vary depending on the image size from a given spatial transcriptomics technique. We then apply L2-normalization across features and reduce dimensionality by UMAP before we perform GMM clustering of genes. UMAP can preserve global and local structures during dimension reduction and previously showed better performance in image clustering than other dimension-reduction methods such as Isomap and t-SNE [7, 9]. The purpose of this clustering is to generate labels so that we can update the ConvNet as in other common supervised neural network training approaches. When the ConvNet is randomly initialized, features extracted by this ConvNet are weak. However, using them to generate labels can still guide the ConvNet to learn more discriminative features. Indeed, Caron et al. showed DeepCluster can learn from weak signal to bootstrap the discriminative power of a ConvNet [6]. Instead of giving each gene a single cluster label, we assign an auxiliary target distribution as a soft assignment. This approach emphasizes genes with high confidence in the clustering task and discounts noisy labels persisting from the random initialization of ConvNet. Doing this can also lead to more stable target values for training the neural network [8]. Finally, we use these soft assignments to train the ConvNet. We add a fully connected layer after the ConvNet that produces probabilities for each gene being assigned to each label. Thus, we can optimize the model by minimizing bi-tempered logistic loss based on Bregman Divergences between the soft assignments from GMM clustering and the probabilities from the fully connected layer [10]. In summary, the CoSTA approach uses a ConvNet clustering architecture which repeats (1) generating features by ConvNet, (2) generating soft assignments by GMM clustering, and (3) using soft assignments to update ConvNet. Once we finish training, we only retain the trained ConvNet for the purpose of feature extraction. Since the ConvNet mainly consists of convolutional layers, the final vector for each gene extracted by ConvNet should be a spatial representation. Using this spatial representation, we can then quantify the relationship between any two genes within one spatial transcriptomics dataset, visualize all SE genes in this dataset by UMAP, and assign patterns through common clustering algorithms. Further details about the rationale of this learning architecture can be found in [Methods](#).

Rationale for using spatial patterns rather than exact pixel overlap

To demonstrate the spatial information lost by overlap analysis and why a spatial representation approach such as CoSTA is useful, we present a simplified biologically-inspired conceptual example (Fig. 1b). In biological tissue sections, we commonly observe structures such as a tightly connected epithelial layer of cells (rectangles in the cartoon) adjacent to a collection of stromal cells (circles). In this example, the spatial expression patterns of three genes are shown. Comparing gene expression patterns by overlap only, we observe that Gene 1 and 2 have the same amount of overlap as Gene 1 and 3 (40%). Thus, an overlap approach to measure gene pattern similarity, like the one used in previous Slide-seq analysis, would report that Gene 1 is equally similar to both Gene 2 and Gene 3 [3]. However, biologically, it is relevant that Gene 1 and Gene 2 are expressed primarily in the epithelial layer while Gene 3 is expressed in the stroma. This biological

difference is not detected by strict overlap, but instead requires a spatial representation that would detect the vertical stripe of epithelial layer expression as a salient pattern. In computer vision, filters are commonly used to find this kind of local correlation, and the success of ConvNet in pattern recognition also relies in the use of filters for identifying local correlations. Using signals of how these 3 genes respond to the same filters, a ConvNet approach will identify Gene 1 and 2 as more similar and Gene 3 as less similar. Therefore, we are motivated to use our ConvNet clustering based CoSTA approach to prioritize similar shape more than overlap for biological cases where layers of cells and the overall patterns of groups of cells matter more than independent individual cell identities [11].

Tests on synthetic data show CoSTA's high specificity, reliance on spatial relationships, and ability to distinguish signal from noise

As a first test of CoSTA's ability to detect correlated spatial patterns in the absence of exact overlap, we use the MNIST handwritten digit image data [12]. When the aim is to find which digits have correlated handwritten patterns to the digit 3, CoSTA identifies only other instances of digit 3 as correlated (100% specificity). In contrast, overlap analysis finds some samples of all other digits as correlated digits of 3 (58% specificity) (Additional file 1: Fig. S1). Meanwhile, CoSTA identifies a smaller subset of the digit 3 as correlated (35% sensitivity) while overlap analysis captures more correlated digits overall (65% sensitivity) in its less specific set (Additional file 1: Fig. S1). As shown below, this increased specificity but possibly decreased sensitivity of CoSTA compared to other techniques appears to hold true in biological data as well.

Before applying CoSTA to real spatial transcriptomics data, we next tested its performance on 5 synthetic datasets, simulated based on real expression patterns from mouse olfactory bulb, following the simulation method in SPARK (Fig. 1c left panel) [5, 13]. We generated 2000 fake gene expression matrices for each pattern, to mimic data for 10,000 total genes. To simulate noise and variability for each gene, we added residual errors onto each spatial coordinate independently based on a normal distribution with mean of zero and variance ranging from 0.2 to 0.6. We then evaluated whether CoSTA could assign each simulated noisy gene to the correct pattern. To compare the CoSTA-derived cluster assignments to the true labels, we use the well-established cluster comparison metric Normalized Mutual Information (NMI) [14]. The NMI approaches a value of 1 as the assignment of genes to the 5 patterns becomes more and more accurate. When CoSTA was initialized, the NMI ranged from 0.27 to 0.57 (Fig. 1c right panel). As training proceeded, CoSTA learned discriminative features to distinguish the 5 patterns, eventually achieving NMIs from 0.85 to 0.98 against the true class label (Fig. 1c right panel, Additional file 13: Table S1). For the highest noise level (0.6) we found that combining both center loss (CL) and bi-tempered logistic loss during CoSTA training substantially improved CoSTA's accuracy (NMI increased from 0.52 to 0.91). However, CL pushes samples toward the 5 centroids and is only applicable when the final number of patterns is known. Thus, we do not include CL for true biological situations.

To demonstrate that CoSTA learns spatial rather than pixel-level patterns from these synthetic datasets, we shuffled the pixel positions in these synthetic datasets. Shuffling all the gene matrices exactly the same way keeps the pixelwise overlap

information identical while disrupting correlations between neighboring pixels, thus destroying the spatial pattern (see [Methods](#) for details). If a pattern detection method is successfully using spatial relationships between neighboring pixels, its ability to classify patterns should be disrupted by this kind of shuffling. Indeed, we found that CoSTA cannot distinguish the genes into correct pattern labels as well with shuffled data (NMI ranges from 0.32 to 0.89), demonstrating that CoSTA is detecting spatial features that depend on the positions of neighboring pixels, rather than features that can be captured by a set of single pixels (Additional file 2: Fig. S2 and Additional file 13: Table S1). When we applied the CoSTA model trained at 0.4 noise level to progressively more shuffled images, we found that the ability to classify genes into groups declined proportional to the amount of shuffling (Additional file 3: Figure S3A). We also tested SpatialDE on these true and shuffled synthetic datasets. SpatialDE performed very well on the true datasets, as expected. However, shuffling the data did not usually change the performance of SpatialDE (Additional file 13: Table S1), indicating an important difference between CoSTA and SpatialDE: SpatialDE is more likely to detect patterns of individual pixels while CoSTA emphasizes the spatial positions of these pixels relative to each other and overall shapes of patterns.

Using this same synthetic data, we next performed a disruption test to demonstrate a disadvantage of using individual pixels as features to analyze spatial transcriptomics data. For half of the simulated gene matrices, we masked a certain region of the pattern, and the masked region doesn't change expression pattern visually (Additional file 3: Fig. S3b). This mimics a situation in which a certain region is obscured or not sampled well for technical reasons experimentally. Using pixel overlap to identify patterns, in this case, assigns masked and unmasked genes into separate groups, even though they otherwise belong to the same pattern. In contrast, CoSTA is resistant to this disruption (Additional file 3: Fig. S3b).

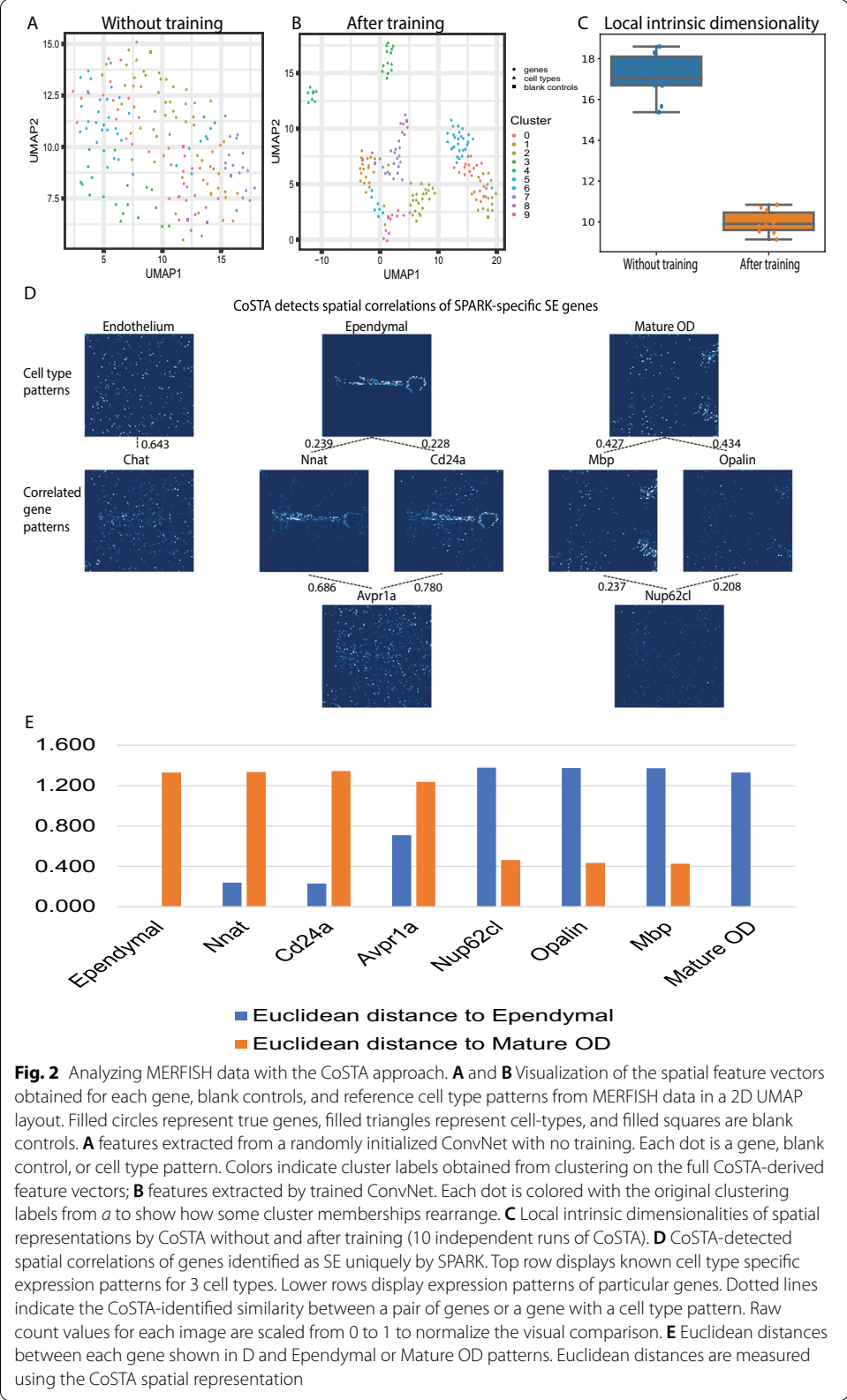
In real spatial transcriptomics data, not all genes will belong to a clear spatial pattern—some genes that are not relevant to the given tissue or condition may only yield random noise or be fairly uniformly expressed. To mimic this situation, we further followed the simulation approach in SPARK to generate synthetic datasets that have 5 spatial patterns and have mixed SE (spatially expressed) and non-SE genes (Additional file 4: Fig. S4). We trained CoSTA on these data with different ratios of SE and non-SE genes, from 90:10 to 10:90. We found that the representation of SE genes by CoSTA is distinct from non-SE genes, even when CoSTA was trained with a high percentage of non-SE genes. Meanwhile, CoSTA demonstrates the capacity to distinguish different patterns of SE genes even when non-SE genes exist (Additional file 4: Fig. S4). Further, CoSTA does not separate even a large number of non-SE genes into separate categories, showing that it does not create false signal out of noise. Here, we also note that a strength of CoSTA compared to methods like SpatialDE is that the output feature vector enables visualization, as is presented throughout these simulation results. While Spatial DE can classify genes into categories, it does not produce a result that can visualize how SE and non-SE genes are separated as we did here for CoSTA. Overall, the performance of CoSTA with synthetic data demonstrates that CoSTA can learn discriminating spatial features.

CoSTA classifies genes by cell type and identifies quantitative relationships between genes in MERFISH data

To extend the application of CoSTA to real spatial transcriptomics data, we first applied it to reanalyze a MERFISH dataset (see MERFISH Analysis in [Methods](#) for complete details) [15]. In order to compare with published analyses using the SPARK approach, we focused on the same slice of the mouse hypothalamus (Bregma + 0.11 mm from animal 18) [5]. The expression patterns of a set of 155 genes expected to be spatially variable were measured with MERFISH for this slice, along with 5 blank control genes. We first initialized a ConvNet and forwarded the MERFISH spatial gene expression matrices through it to obtain gene feature vectors. Then we clustered the 155 spatially variable genes with the 5 blank genes and with 9 cell type-specific expression patterns defined by the original publication through a combination of MERFISH and scRNA-seq data. We clustered these genes, controls, and cell type patterns into 10 groups and visualized them by UMAP. Without training, SE genes, control genes, and cell types are spread across the 2-dimensional UMAP space and boundaries between groups are not distinctively defined (Fig. 2a). Next, we trained the CoSTA model to obtain refined feature vectors. After training, SE genes, control genes and cell types formed distinct groups that have clearer boundaries in the 2D visualization (Fig. 2b) and refined cluster memberships that reproducibly and quantitatively form tighter clusters according to a linear intrinsic dimensionality (LID) estimator (Fig. 2c) [16].

From this MERFISH data, SPARK identified 145 SE genes including one blank control, and SpatialDE found 139 SE genes with one blank control [5]. CoSTA is designed primarily to detect similarities between spatial gene expression patterns, rather than to estimate spatial relevance (identify SE genes). So, to define which genes are called SE by CoSTA, we examined which genes CoSTA identified as highly correlated to one of the 9 pre-defined cell type specific expression patterns. We found a correlation threshold at which CoSTA identified 133 SE genes associated with one of the different cell type patterns, while none of the blank controls were called associated with a pattern (Additional file 14: Table S2). Thus, CoSTA's sensitivity is slightly lower than SPARK and SpatialDE, but with higher specificity (no blank controls detected). However, CoSTA's result is both more sensitive and more specific than the Trensceek approach, which only identified 108 SE genes and one blank control [17].

Three genes in this MERFISH dataset, *Avpr1a*, *Chat*, and *Nup62cl*, were highlighted by Sun et al., because they were only identified as SE by SPARK [5]. CoSTA is able to identify the spatial expression patterns of these genes, but also reveals by quantitative similarity that these genes are more distantly related to cell type expression patterns than other genes. We examined both the significantly similar groups determined by CoSTA and used the spatial representation learned by CoSTA to measure Euclidean distances of these genes to each other and to cell type expression patterns (Fig. 2D, E and Additional file 14: Table S2). For example, CoSTA identifies genes such as *Nnat* and *Cd24a* as significantly similar to the Ependymal cell type pattern (Dotted lines, Fig. 2D). *Avpr1a* is quantified as more distant from this Ependymal pattern (Fig. 2E), though it does show some similarity to *Nnat* and *Cd24a* (Fig. 2D). Similarly, *Mbp* and *Opalin* are significantly correlated to the Mature OD cell type pattern (Fig. 2D, E and Additional file 14: Table S2). *Nup62cl* is more distant from the Mature OD than *Opalin* and *Mbp*,

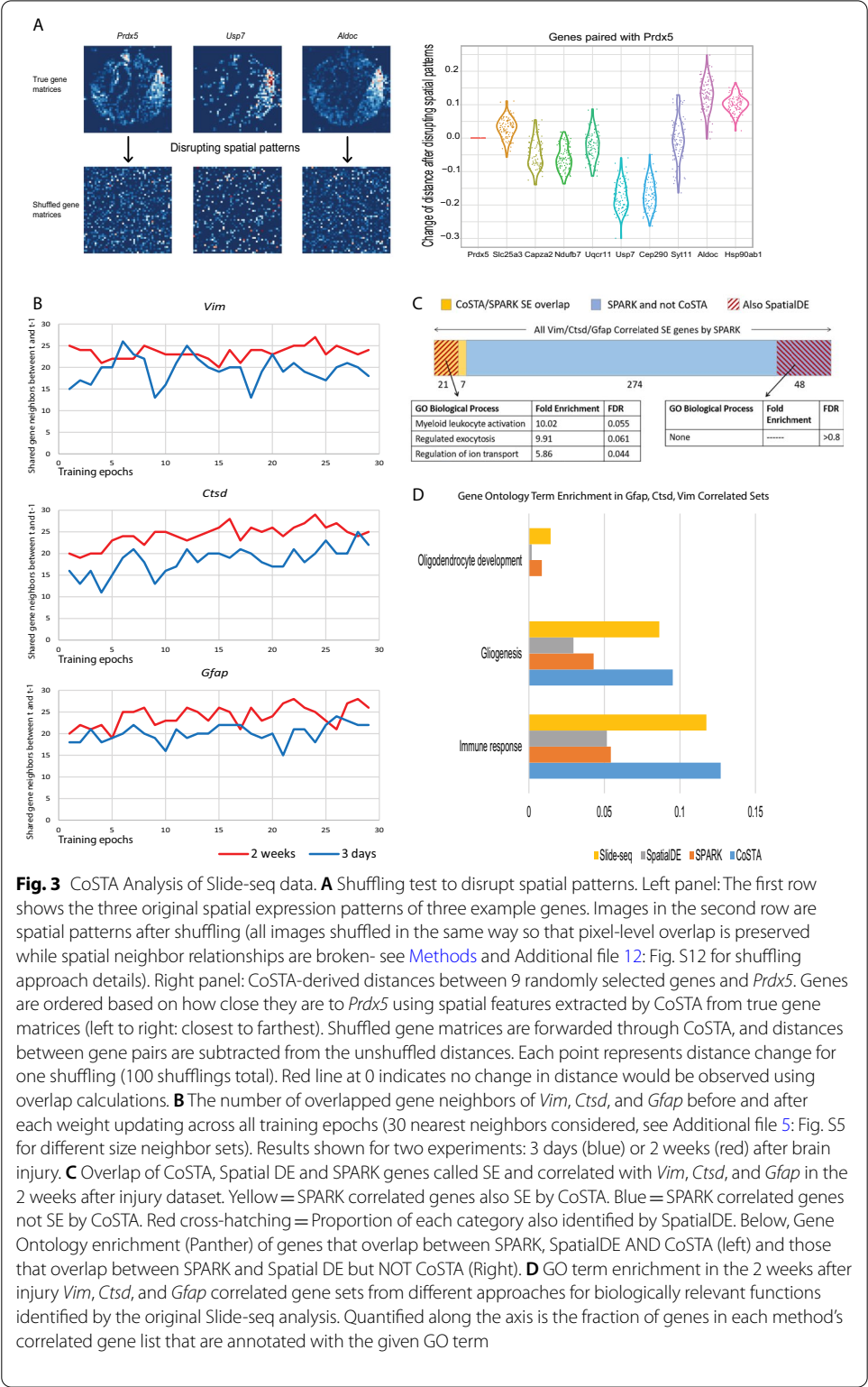


but is related to the expression patterns of *Mbp* and *Opalin*. Visual inspection of *Avpr1a* and *Nup62cl* confirms that these patterns are quite noisy and less similar to the key cell type pattern (Fig. 2d). Thus, by quantifying relationships between patterns rather than reporting uniform sets of SE genes, CoSTA clarifies that these genes are likely identified by SPARK and no other method because they are in fact less spatially similar to key cell type patterns. CoSTA's ability to quantify relationships between genes, rather than only categorizing genes, is important in biological situations, where there is often going to be a range of relative similarity that would be oversimplified by strict categorization.

CoSTA learns spatial pattern-dependent representations of Slide-seq data

We next expand our application of CoSTA to Slide-seq data. Slide-seq takes advantage of high-throughput single cell RNA sequencing and barcoding. Therefore, it enables spatial gene expression measurement for all genes in the genome [3]. As a first demonstration that CoSTA can be applied to this type of high-throughput spatial transcriptomics data, we performed an experiment-mixing test to evaluate whether CoSTA can separate different spatial patterns. Due to the unavailability of a “gold standard” for positive and negative spatial similarity of gene expression, we mixed gene matrices from four different spatial transcriptomics experiments by Slide-seq and tested the ability of CoSTA to deconvolve them [3]. Each overall experiment is performed on an independent brain slice of a different mouse, so the shapes and spatial features of each experimental sample overall constitute a large difference between experiments. Each gene within each experiment will have a somewhat different pattern (and it will be our next goal to distinguish those differences and similarities), but we first tested whether genes within the same experiment could be classified together based on their overall spatial features. We implemented training as above and then clustered the mixed experiment gene matrices into 4 clusters. The confusion matrix shows clustering labels are largely consistent with true experimental labels (Additional file 15: Table S3).

We next performed a shuffling test on gene matrices from one Slide-seq experiment, to break correlated patterns of neighboring regions in the way described for the shuffling of synthetic data above (see [Methods](#) for shuffling details). We trained a new model and examined model-reported similarity among expression patterns of ten random genes. If CoSTA successfully learned spatial features that distinguish the expression of these genes, the distances between two genes should change when spatial patterns and relationships between neighboring pixels are disrupted. We randomly selected *Prdx5* as the reference gene and calculated Euclidean distances of 9 other genes with it. We order these ten genes based on their distances to *Prdx5*. Then, we shuffled gene matrices 100 times, passed the shuffled matrices through the trained ConvNet, and recalculate paired distances with *Prdx5* (Fig. 3a). We find that in 5 of 9 comparisons, distances decreased upon shuffling, as the distinctive patterns captured by CoSTA were removed by shuffling, converting the matrices into generic, more similar patterns. In 4 of 9 comparisons, distances increased with shuffling, likely indicating that key similarities between the spatial patterns became disrupted during shuffling (Fig. 3b). In contrast, the similarity measured by overlap analysis would not change after shuffling since individual pixels were shuffled identically. This result again suggests, this time using real biological data, that the learned features by CoSTA are strongly tied to the spatial expression pattern.



Ensemble learning identifies stable relationships between spatial gene expression patterns

We next applied CoSTA to reanalyze two spatial transcriptomics datasets measured by Slide-seq [3]. These datasets are derived from two biological conditions: 3 days after

brain injury (“3 days”) and 2 weeks after brain injury (“2 weeks”). In the first investigation of these two datasets in Slide-seq, Rodriques et al. primarily focused on genes that were spatially correlated with *Vim*, *Ctsd* and *Gfap* at both 3 days and 2 weeks after brain injury [3]. For comparison, we also examined genes correlated with *Vim*, *Ctsd* and *Gfap* from our CoSTA results. One property of our approach is that features of each gene change every epoch when weights are updated. This may result in changes to the nearest neighbors of a gene during model training and can be used to infer how strong and stable the inferred spatial relationships are in a given condition. We measured the overlap between detected *Vim*, *Ctsd*, and *Gfap* neighbor genes before and after weight updating across training epochs, and we found neighbors tend to be more stable for the 2 weeks dataset than for the 3 days dataset (Fig. 3b and Additional file 5: Fig. S5). This may indicate that in the acute phase after injury, genes related to *Vim*, *Ctsd* and *Gfap* are more variable and less spatially patterned, but these patterns become stronger at 2-week time point after injury.

To screen significantly spatially patterned genes out from noise, we use ensemble learning. Briefly, we initialized 5 ConvNets and trained them separately. We then calculated the nearest neighbors for every gene in the same dataset, at neighbor set sizes of 5, 10, 15, 20, 25, 30, 40, 50, and 100. We use a broad range of neighboring levels because different genes may form different sizes of communities. Next, we calculated Jaccard similarities across the 5 CoSTA models and keep genes that have an averaged Jaccard similarity larger than 0.2 at least in one level. We call genes that pass the threshold “stable”, and genes that are filtered out as “unstable”. We propose that the percentage of stable vs. unstable genes represents the degree of spatial patterning in the experiment set. Overall, a smaller proportion of genes were considered stable at 3 days, consistent with the more variable gene neighbors observed for the 3-day condition above. These ‘stable’ genes can also be considered a CoSTA-derived set of ‘spatially expressed’ (SE) genes for comparison to SE genes identified by SPARK. The majority of CoSTA-SE genes are also called SE by SPARK (86% at 3 days and 78% at 2 weeks, Additional file 6: Fig. S6a). *Vim*, *Ctsd*, and *Gfap* are considered SE by CoSTA in the 2-week data but not in the 3-day dataset. Notably, *Vim*, *Ctsd*, and *Gfap* are also not present in the 3 days SE gene list identified by SPARK, and only *Ctsd* and *Gfap* were identified as SE genes by SPARK in the 2 weeks data. We note that less strongly patterned genes could reflect actively variable biological regulation (such as might happen during acute response to injury), not only technique noise. We are unable to definitively distinguish a weak spatial pattern from inherent noise, because of lack of “ground truth” for pattern matching. However, we can, as above, disrupt spatial patterns by shuffling the true datasets, maintaining pixelwise correlations between genes but removing spatial information (see [Methods](#) for shuffling approach details). We shuffled a whole set of gene matrices from 3 days and 2 weeks and applied CoSTA to these datasets. As when we shuffled simulated data in Additional file 2: Figure S2, we find that this shuffled dataset has overall lower NMI than its original dataset during training (Additional file 6: Fig. S6b; see [Methods](#) for details of NMI use). Further, substantially fewer SE genes are identified in the 2 week randomized data as compared to the real data (Additional file 6: Fig. S6c). This again demonstrates that CoSTA captures spatial features that are distinct from individual pixel information. For true

3-day and shuffled 3-day data, there is not a clear difference in the number of identified SE genes (Additional file 6: Fig. S6c). This again suggests that the spatial patterns are much less strong within the 3-day dataset. Indeed, few patterns are visually obvious for example gene matrices from 3 days (Additional file 7: Fig. S7a).

CoSTA identifies smaller, but specific and biologically relevant, sets of spatially correlated genes compared to SPARK and SpatialDE

We focused our further analysis on the 2-week data. We applied SpatialDE and SPARK to this dataset for comparison to CoSTA. The original Slide-seq publication previously identified 843 genes that are correlated with *Vim*, *Ctsd*, and *Gfap* via overlap analysis [3]. However, CoSTA, with a rigid neighbor similarity stability threshold, identified many fewer correlated genes (63 with z-scores < -2.325), and only 19 genes matched the original Slide-seq set (Additional file 8: Fig. S8a). SPARK first identified 1294 significantly SE genes and then clustered them into 10 groups by hierarchical clustering with individual pixels as features. Our CoSTA correlated gene list only has 5 gene overlaps with genes that are grouped with *Vim*, *Ctsd*, and *Gfap* by SPARK. We also used SpatialDE to find significant SE genes. Surprisingly, the whole dataset passed the SpatialDE test for significant spatial expression. Then, we applied the unsupervised pattern detection algorithm built in SpatialDE to cluster genes into 10 groups. This resulted in a large number of genes grouped with *Vim*, *Ctsd*, and *Gfap*. A majority of our CoSTA set (41 genes) overlaps with genes identified by SpatialDE (Additional file 8: Fig. S8a). The set of correlated genes identified by CoSTA is much smaller than sets identified by the other 3 methods. This is in part because CoSTA requires stable relationships between neighboring genes to be classified as an SE gene at all, and only SE genes can then be identified as highly similar to the genes of interest. Indeed, out of 350 genes identified by SPARK as correlated with *Vim*, *Ctsd*, and *Gfap*, only 28 are even classified as SE genes by CoSTA. However, we observe evidence that this CoSTA-identified SE subset is reliable and meaningful. First, 75% of these CoSTA-SE genes identified by SPARK are also identified as correlated by SpatialDE, while in the remaining non CoSTA-SE set, there is only a 15% overlap between SPARK and SpatialDE (Fig. 3c). Further, the genes overlapping between SPARK, SpatialDE, and CoSTA have biologically relevant function enrichment (such as ion transport and exocytosis) while the genes overlapping between SPARK and SpatialDE but not identified as SE by CoSTA show no function enrichment at all (Fig. 3c). We also observe visible evidence of spatial pattern similarity to the 3 genes of interest among genes considered SE and highly correlated by CoSTA and less evidence of similarity for genes identified only by SPARK (Additional file 8: Fig. S8b).

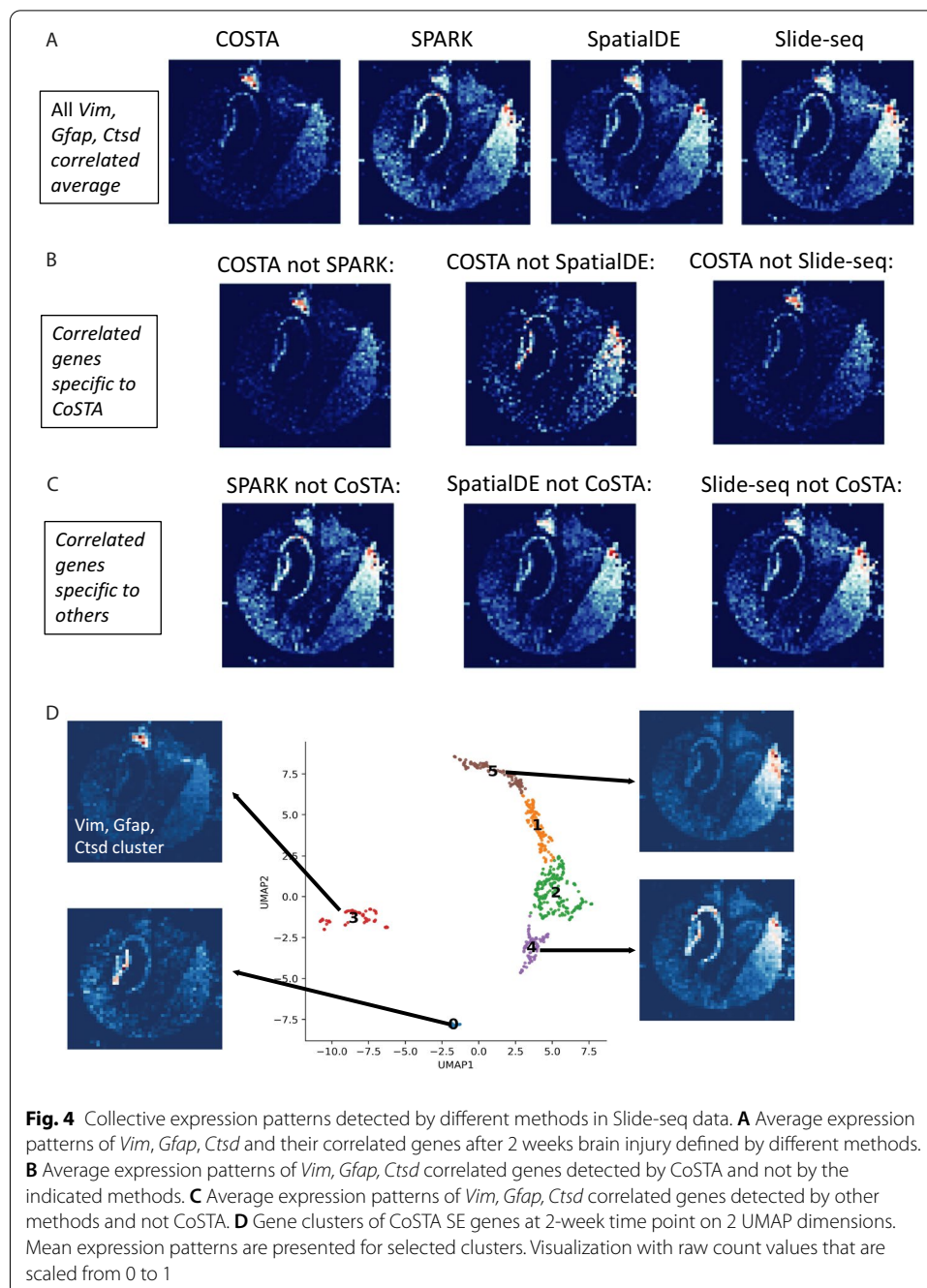
Further, we find that the 63 genes identified by CoSTA as significantly correlated to *Vim*, *Gfap*, and *Ctsd* are highly enriched for meaningful biological function. In the original study, Rodriques et al. highlighted that genes correlated with *Vim*, *Ctsd*, and *Gfap* are enriched for functions in immune response, gliogenesis and oligodendrocyte development—all functions that are biologically expected in response to injury [3]. We found that the correlated genes identified by CoSTA have higher enrichment in immune response and gliogenesis than the genes identified by SpatialDE, SPARK and this original Slide-seq report (Fig. 3d). However, none of genes fall into category of oligodendrocyte

development. When we visually inspected expression patterns of genes in the category of oligodendrocyte development, their individual and collective patterns do not have similarities to expression patterns of *Vim*, *Ctsd*, and *Gfap*. They are either noisy or expressed globally (Additional file 7: Fig. S7b). From results above, we conclude that CoSTA returns a reduced, stringent set of correlated genes that are more enriched for biological significance than the larger sets returned by other methods.

As noted earlier, one key difference between CoSTA and these other methods is that CoSTA provides not only sets of similar genes, but also quantitative pairwise comparisons between all genes. Thus, we can extract from CoSTA a ranked list of how similar each CoSTA-SE gene is to *Vim*, *Ctsd*, and *Gfap* (Additional file 16: Table S4). This enables us to search for enriched biological functions using similarity rankings rather than an arbitrary cutoff using the GOrilla enrichment tool [18]. Using the whole ranked list, we find novel enriched functions such as collagen metabolism, astrocyte differentiation, and vascular endothelial growth factor signaling that may be relevant to damage repair (Additional file 8: Fig. S8c). Pixelwise correlation can also be used to create a ranked similarity list. When these two approaches are compared, we observe some highly ranked genes shared by both approaches with clear pattern overlap to the query genes. Where the two approaches differ widely in gene ranking, CoSTA-specific genes tend to have the key patterns of expression as contiguous patterns superimposed on a generic weak background while pixel-specific genes tend to have isolated pixels overlapping the key areas (Additional file 8: Fig. S8d).

To avoid observation bias by only looking at a few example genes that are classified differently by different approaches, we next globally compared the types of spatial patterns detected uniquely by CoSTA and other previous methods. For each method (CoSTA, SpatialDE, SPARK, and the original Slide-seq overlap approach), we consider the list of genes classified as spatially correlated with *Vim*, *Ctsd*, and *Gfap* as described above. The average expression pattern of genes detected as correlated to these query genes varies somewhat according to approach. Notably, the average CoSTA pattern is more localized to the upper right region, where the damage was induced (Fig. 4a). In contrast, SPARK, SpatialDE, and Slide-seq each identify so many correlated genes that their average pattern looks very much like the average pattern of all genes in the dataset (compare Fig. 4a with Additional file 9: Fig. S9) rather than distinctive. This again emphasizes the smaller, but perhaps more specific, set of genes identified as correlated by CoSTA. When we compare genes identified as correlated by CoSTA and not certain other techniques, we can see that CoSTA-unique genes have certain local patterns that were not captured as much by other methods (Fig. 4b). In contrast, again, genes detected by other methods and not CoSTA look more similar on average to the average gene expression of the entire gene set (Fig. 4c).

Finally, rather than using a significant correlation threshold, we clustered all CoSTA-determined SE genes at the 2-week time point into 6 groups using the learned spatial representation. The cluster that contains *Vim*, *Ctsd*, and *Gfap* (cluster 3) is composed of 89 genes expressed in a distinct pattern (Fig. 4d and Additional file 17: Table S5). Other clusters also successfully identify distinctive spatial patterns of expression (Fig. 4d and Additional file 9: Fig. S9). We also used SpatialDE to cluster



SE genes identified by CoSTA into 6 clusters. We found that the two methods share many commonalities in detecting patterns, with some disagreements (Additional file 9: Fig. S9). Notably, when only the narrower set of SE genes identified by CoSTA is used, the cluster of genes identified by SpatialDE containing *Vim*, *Gfap*, and *Ctsd* (cluster 2, Additional file 9: Fig. S9) has a much more specific, localized pattern than when SpatialDE default settings are used to classify all genes. This again suggests that CoSTA provides a meaningful increase in specificity by identifying genes with stable spatial relationships.

Discussion

We have shown that our CoSTA approach can successfully implement deep learning ideas from computer vision to infer spatial gene expression relationships. This approach can be applied to any technology that outputs an image-type matrix of gene expression information for each gene, including not only Slide-seq [3, 19] and MERFISH [20] explored here, but also STARmap [21], 10 × Visium (10 × Genomics), and HDST [22] approaches. Identifying spatial patterns from high-throughput spatial transcriptomics data is still challenging, however. We often do not have a clear ground truth answer for what should be detected as a pattern vs. noise and what similarities in patterns are most biologically relevant. Different approaches will have different strengths and weaknesses depending on the types of patterns and relationships to be detected. The very first step in any approach to analyzing spatial transcriptomics data is estimating significant SE genes. To identify SE genes, SpatialDE relies on the assumption that spatial expression of a given gene follows a multivariate normal distribution across spatial coordinates [4]. However, this assumption leads all genes in a Slide-seq dataset to be identified as SE genes by SpatialDE. This may occur because noisy signals generated by the Slide-seq experiment may also follow or are confounded within the multivariate normal distribution. Therefore, a multivariate normal model will not be able to distinguish spatial patterns from noise in certain types of experimental data. Different from SpatialDE, both SPARK and CoSTA make use of kernels to identify SE genes. SPARK defined 5 periodic and 5 gaussian kernels to cover a range of possible spatial patterns that the authors believe are observed in common biological datasets [5]. Therefore, identifying SE genes involves a statistical evaluation of how well kernels match spatial patterns of interest. This SPARK approach is very valuable if an experimental dataset is accompanied by prior knowledge about relevant spatial patterns. Kernels in CoSTA also serve a similar purpose but are not predefined. Instead, kernels in CoSTA are learned through training a neural network. To identify SE genes, we rely on the idea that a true spatial pattern should be collective, which means a group of genes should share a spatial pattern. Therefore, when we apply kernels learned independently from 5 ConvNets, genes in the same group should have similar responses to these kernels. Conversely, a noisy gene expression pattern would respond to the 5 sets of ConvNet kernels differently, clustered with different groups of genes each time. Indeed, we showed that this kernel approach helps identify a more focused set SE genes in Slide-seq data without requiring an a priori definition of relevant patterns that SPARK requires. We have shown by various measures that the SE genes identified by CoSTA are a much smaller set, but with high enrichment for meaningful biological function, and more likely to be also detected by multiple other methods, increasing confidence in this set.

Identification of SE genes is just the beginning of extracting biological meaning from spatial gene expression. Careful analysis of the spatial relationships between genes is also necessary. Often, as in overlap analysis, studying gene relationships is based on vectorizing gene expression patterns and measuring their similarities in a latent space without considering spatial information such as the position of neighboring datapoints. One key motivation for CoSTA, therefore, is to preserve a spatial and shape representation of gene expression patterns. In comparison, SPARK does not have a pattern detection function, but can be combined with hierarchical clustering with pixels as features to

assign each gene a pattern label. SpatialDE implements a clustering model based on a spatial Gaussian-process-based (GP) prior [4]. This clustering model is an extension of GMM with the addition of a spatial prior on cluster centroids. Therefore, pattern detection by SpatialDE goes beyond the pixel level. In our method, we define the key goal as learning a spatial representation for each gene. We have demonstrated that features learned by CoSTA are not isolated to individual pixels, while SpatialDE responds more to individual pixel information in our simulations. Because of use of convolutional layers, spatial features learned by our method represent **local patterns and multiple local patterns together form the global pattern for the gene matrix**. Finally, vectorizing gene matrices allows us not only to find different spatial patterns within a dataset by clustering but also to **study spatial relationships of pairs of genes**. Such a pairwise examination, in contrast, is not implemented in SpatialDE.

Not only in detection of a narrower set of SE genes, but also in identifying relationships between genes, our results consistently suggest that **CoSTA provides more specific though less sensitive results than other methods**. Throughout our analyses, we find that overlap approaches, as well as SPARK and SpatialDE tend to group together larger sets of genes that are **more distant in their spatial pattern relationships, while CoSTA captures a narrower and more specific set of genes**. This was observed in our analysis of digit image data as well as in applications to Slide-Seq and, to a lesser extent, MERFISH. This difference in outcomes again demonstrates the different advantages and disadvantages of different approaches. CoSTA would likely be more useful in a case where users want to **narrow their set of candidate related genes for future experiments**. We also note throughout the Methods section alterations to parameters of CoSTA that could allow for detection of more general patterns.

Again, depending on the biological reality underlying the data, different approaches will have different advantages. The CoSTA approach will have advantages in cases where **overall pattern shape is important**, while direct overlap calculations may perform better when exact cell to cell correlation is more biologically relevant. The CoSTA approach may also have future applications to datasets in which images of different genes are not from the identical biological section, but instead from neighboring tissue slices, as is common in traditional histology. If a pattern or shape of expression is maintained while exact overlap is lost, as we demonstrated with our simulated masking approach, CoSTA could still detect such a pattern similarity where an overlap approach would not.

Conclusions

In this study, we demonstrated that our deep learning CoSTA approach provides a different angle to spatial transcriptomics analysis by focusing on the shape of expression patterns. CoSTA includes more information about the positions of neighboring pixels than does an overlap or individual pixel correlation approach. CoSTA can be applied to any form of spatial transcriptomics data that are represented in matrix form to find genes expressed in similar patterns as well as to evaluate the strength of the spatial patterning of each gene. We find that CoSTA captures more focused groups of spatially related genes while still detecting the biological function information found by other approaches that report larger sets of related genes.

Methods

Resizing gene images and normalization

The Slide-seq data was obtained from: https://portals.broadinstitute.org/single_cell/study/slide-seq-study. The raw images of Slide-seq consist of over 1,000,000 pixels, which makes computation difficult. Therefore, we first binned 100 pixels into one pixel and resized matrices from different experiments into the same 48X48 image size. This results in a lower resolution, which may obscure small-scale fine details, but large scale global features of expression patterns of genes are preserved. CoSTA can be applied to any spatial transcriptomics dataset at any resolution, as long as the user has sufficient computational resources available. To avoid extreme computational burden, we recommend that users interested in high resolution features zoom into regions of interest and crop images in that region to efficiently apply CoSTA to their data. After binning, we normalized gene matrices as described in Svensson et al. [4]. This normalization involves finding the total gene expression counts for each pixel across all gene matrices and then normalizing each pixel of each matrix by the log total counts across all matrices for this pixel. If this normalization is not performed, the expression of a gene could be over or undercounted at certain spatial locations where expression levels were systematically high or low for all genes. Normalization by total counts at each pixel ensures that our approach captures the spatial covariance for each gene beyond this potential artifactual effect. For visualization of expression patterns, we instead use averaged raw count values, and scale values from 0 to 1 divided by the maximum value. Thus, expression images in all figures are in 0 to 1 scale. This allows a more direct visual inspection of the raw data.

CoSTA architecture

1. ConvNet

The ConvNet stage of CoSTA consists of 3 convolutional layers for Slide-seq and MERFISH analysis. Inputs are sets of spatial gene expression images (matrices) as described above. We first initialize a ConvNet randomly and then forward these gene expression matrices through the ConvNet. All weights in convolutional layers are initialized on a Xavier uniform distribution. Each convolutional layer is activated by a rectified linear unit function and is followed by a batch normalization layer and a max pooling layer to reduce the size of the output. To produce a feature vector for each gene, we flatten the matrix output from the last max pooling layer by concatenating all matrix columns into a single column. One fully connected layer is added to the model after the last max pooling layer with a customized softmax activation to produce outputs as probabilities (See **4. Loss Function**). The fully connected layer is only used during training, when we need gradients to pass backwards through the model. Once trained, this fully connected layer will be discarded, and we use L2-normalized outputs as the spatial representations. Specific parameters used in ConvNet, such as the number and size of filters in each convolutional layer, can be found in python code. We note that different numbers of convolutional layers have been used for

different image classification tasks. We recommend that users start with a 3-convolutional-layer network for initial data exploration. However, if a dataset has a larger size of gene matrices, outputs from the 3-convolutional-layer network will be very long vectors. Therefore, users can increase the number of convolutional layers to decrease the dimensions of outputs if needed.

2. UMAP and clustering

The flattened spatial representation vector output from the three convolutional layers is reduced by UMAP before GMM clustering. We implemented UMAP using the original python source code [9]. We set up “n_neighbors=20” and “min_dist=0”, while using UMAP for dimension reduction. To cluster samples into N clusters, a user can reduce dimensions to N UMAP-dimensions. In this study, we reduce all samples to 30 UMAP-dimensions and cluster all samples into 30 clusters by GMM. While 30 clusters are used here for the model training purpose, once the model is trained, the user can use the final output vector of spatial features to cluster genes into any number of groups desired. To test the influence of the initial choice of number of clusters, we tested 10, 20, and 30, 50, 75, and 100 clusters in 2-week Slide-seq data. Using larger numbers of clusters leads to the identification of fewer SE genes (Additional file 10: Fig. S10a). Our model can converge no matter how many clusters are used for training (Additional file 10: Fig. S10b). For a purpose of comparison, we called the 15 nearest genes of *Vim*, *Ctsd*, and *Gfap* individually, and total 45 genes in one test as correlated genes were used for comparing effects of the number of clusters. The choice of the number of clusters will influence the scale of correlated expression pattern detected (Additional file 10: Fig. S10c). More global pattern differences will be detected using smaller numbers of clusters while finer scale pattern distinctions are detected with larger numbers of clusters (Additional file 10: Fig. S10c). Increasing the number of clusters will also bring a disadvantage of larger computational cost and longer training time (Additional file 18: Table S6). In this case, 30 clusters show good specificity, and the detected spatial pattern is not further refined with increasing cluster numbers (Additional file 10: Fig. S10c). Without ground truth for a dataset, the number of clusters must be chosen based on the scale of patterns desired to be detected for a particular biological application and the results inspected visually.

3. Auxiliary target distribution as soft assignment

After clustering, we calculate centroids by averaging samples in the same cluster (Eq. 1).

$$c_i = \frac{1}{M_i} \sum_{j=1}^{M_i} x_{ij} \quad (1)$$

where c_i is the centroid for the i th cluster, M_i is the total number of samples in this cluster, and $x_{i,j}$ is a reduced UMAP vector for the j th sample in the i th cluster.

Then, each sample is assigned probabilities based on Euclidean distances to cluster centroids (Eq. 2).

$$p(y = i|x) = \frac{e^{1/d_i}}{\sum_{i=1}^N e^{1/d_i}} \quad (2)$$

where d_i is the Euclidean distance of sample x to the centroid c_i , and N is the total number of clusters.

Next, we transform probabilities of each sample to an auxiliary target distribution using Eq. (3).

$$q_{ij} = \frac{p_{ij}^2/f_i}{\sum_{i=1}^N (p_{ij}^2/f_i)} \quad (3)$$

where $f_i = \sum_{j=1}^M p_{ij}$. i denotes the i th cluster and j denotes the j th sample, p_{ij} is probability that the j th sample belongs to the i th that we get through Eq. (2). q_{ij} is the auxiliary target probability that the j th sample belongs to the i th cluster. This transformation was proposed by Xie et al., which is raising p_{ij} to the second power and then normalizing by frequency per cluster [23]. The use of power 2 is to highlight samples that have high confidence in the clustering task and discount samples for which the model is uncertain about cluster assignment.

4. Loss function

To optimize the neural network, we use bi-tempered logistic loss based on Bregman Divergences as the primary loss function. Bi-tempered logistic loss was proposed by Amid et al. and showed advantage of making supervised learning robust to noise [10]. To achieve the robustness, they devised tempered softmax function and tempered logistic loss and gave detailed mathematical reasons behind (Eq. 4, 5). We reason that training CoSTA also faces the problem of unknown noise within the data, because clustering will assign wrong labels to samples. This is even true when clustering is based on the ConvNet that is randomly initialized. Therefore, use of bi-tempered logistic loss is to deal with wrong or uncertain labels generated by clustering. In the equation below, t_1 and t_2 are two temperature parameters proposed in the original work. t_1 controls the log-transformation of input values, while t_2 controls the exponential function of activated input values. When both t_1 and t_2 are equal to 1, bi-tempered logistic loss is the common KL-divergence loss with softmax activation.

$$L = y_i (\log_{t_1} y_i - \log_{t_1} \hat{y}_i) - \frac{1}{2 - t_1} (y_i^{2-t_1} - \hat{y}_i^{2-t_1}) \quad (4)$$

where $\log_{t_1}(x)$ can approximate to $\frac{1}{1-t_1} (x^{1-t_1} - 1)$. y_i is the target value and \hat{y}_i is the predicted value out of the fully connected layer.

$$\hat{y}_i = \exp_{t_2}(\hat{\alpha}_i - \lambda_{t_2}(\hat{\alpha})) \quad (5)$$

where $\hat{\alpha}_i$ is linear activation of output of the fully connected layer for the i th cluster, and $\lambda_{t_2}(\hat{\alpha}) \in \mathbb{R}$ is s.t. $\sum_{j=1}^k \exp_{t_2}(\hat{\alpha}_j - \lambda_{t_2}(\hat{\alpha})) = 1$.

Center loss is an optional setting in our model. Center loss was first proposed to assist models to learn discriminative representations in supervised learning [24]. Optimizing models with center loss is equal to minimizing intra-class variation defined by Eq. (6).

$$L_c = \frac{1}{2} \sum_{j=1}^{M_i} \|x_i - c_j\|^2 \quad (6)$$

where c_i is the centroid of i th cluster, and x_j is the hidden features of j th sample in this cluster.

Because lowering center loss will push samples closer to the cluster center, the learned representations will be more discriminative in the hidden space. Though we did not use center loss to train models for Slide-seq data, we found that adding center loss during training can substantially improve accuracy in Fashion image data (Additional file 11: Fig. S11) and the synthetic data with variance as 0.6. If a user has a biological dataset with some degree of known ground truth for comparison, initial data exploration should explore whether combining center loss and bi-tempered logistic loss is more appropriate to capture the known spatial features of the data.

5. Normalized mutual information

Unlike supervised learning, we do not have ground truth for training in the CoSTA approach. To monitor how well training proceeds, we use normalized mutual information (NMI) to compare clustering labels before and after weight updating across training epochs. Increase of NMI during training indicates a decreased changing of clustering labels and thus suggests convergence of model. We cannot hold aside a validation set during CoSTA training. Therefore, NMI also serves as a metric of overfitting. Once we do not observe a large jump of NMI in consecutive epochs, we consider that the model has converged. For tests with synthetic data, we also use NMI to quantify how well the CoSTA-assigned labels match the true labels.

6. Experiments with common image datasets

MNIST handwritten, USPS-digit, and Fashion image datasets were downloaded from: <http://yann.lecun.com/exdb/mnist/>, <https://www.kaggle.com/bistaumanga/usps-dataset> and <https://www.kaggle.com/zalando-research/fashionmnist>. These datasets come with true labels, and we noticed that the CoSTA approach can learn to predict more true labels than the model that is just initialized and exceeds UMAP + GMM with pixel values as features (Additional file 11: Fig. S11). For the Fashion image dataset, CoSTA was greatly improved after we add center loss with bi-tempered logistic loss as a whole loss function. However, the learning ability of CoSTA with these datasets is less than with supervised learning approaches (typically > 95% accuracy). The highest accuracy we got is 0.961 (MNIST handwritten), 0.931 (USPS-digit) and 0.686 (Fashion), as measured by NMI between the clustering label and true class label. NMIs achieved with CoSTA applied to the MNIST and Fashion datasets are higher than for all other deep learning clustering methods, and the CoSTA NMI for USPS scores second in the ranking of deep learning approach performance [7]. We also tested whether SpatialDE can identify

patterns in these three image datasets. We used the automatic histology pattern detection implemented in SpatialDE to cluster images in MNIST handwritten, USPS-digit, and Fashion into 10 groups, and SpatialDE achieved 0.532 (MNIST handwritten), 0.658 (USPS-digit), and 0.568 (Fashion) NMIs, which are even lower than UMAP+GMM clustering with pixels (Additional file 11: Fig. S11).

Simulation datasets

1. Simulation of synthetic data with 5 patterns

We followed the simulation approach in SPARK to generate 10,000 fake genes that can be assigned into 5 distinct patterns [5]. We added residual errors onto each spatial coordinate independently based on a normal distribution with mean of zero and variance ranging from 0.2 to 0.6, resulting in 5 synthetic datasets with different noise levels. The simulation code can be found at <https://github.com/xzhoulab/SPARK>.

2. Synthetic data with mask

We selected synthetic data with variance as 0.4 for this test. We arbitrarily selected a region to mask, a region in the blue circle of Additional file 3: Fig. S3b. Though it is an arbitrary selection, we intentionally avoid any region that is crucial for each expression pattern. Therefore, the mask region will not disrupt spatial patterns visually. For each pattern, we randomly chose half of the simulated genes and added the mask by suppressing expression in that region to zero. The other half of the simulated genes remain intact. Thus, we have 5,000 genes each with and without masks.

3. Mimic real spatial transcriptomic data by mixing SE and non-SE genes

We still focus on synthetic data with variance of 0.4. We further introduced non-SE genes to build more synthetic datasets that have different ratios of SE and non-SE genes. Code to generate non-SE genes can also be found at <https://github.com/xzhoulab/SPARK>. In this test, we generated 5 synthetic datasets with SE and non-SE ratios ranging from 90:10 to 10:90.

Shuffling approach to disrupt spatial information but not pixel correlation

To evaluate the degree to which CoSTA and other methods detect patterns in space vs. pixelwise information, we used a data shuffling approach that would preserve the pixel-by-pixel correlation between different gene image matrices but disrupt the neighboring pixel spatial pattern (Additional file 12: Fig. S12). Each gene matrix (expression image) was flattened into a single vector by concatenating all rows of the matrix into a single row. Then, the positions of individual elements in these vectors were shuffled identically for all gene expression images. That is, for all images, the data at position 2 might now be at position 10, while position 10 would now be at position 3, etc. Then, the vectors were reformed into a matrix and these identically shuffled gene matrices were passed through a given analysis tool. Thus, any methods (such as pixel overlap or correlation) that depend only on the one-to-one relationship between pixels between two gene images

will perform exactly the same on the shuffled and original data. In contrast, methods that capture information about shared patterns between neighboring pixels (such as a broader patch of high gene expression) would perform differently on the shuffled data when this previously existing broader pattern is disrupted.

SE gene calling

To call out SE genes, we use an approach of ensemble learning. Simply put, we train 5 CoSTA models independently. We then calculate a set of nearest neighbors for every gene in the same dataset, using neighbor set sizes of 5, 10, 15, 20, 25, 30, 40, 50, and 100. This is because different genes with their neighbors may form a community with different sizes. Using a broad range of neighboring set sizes can enable us to include SE genes that only form a small community with a few genes as well as SE genes that fall into a large gene group. Next, we calculated Jaccard similarities across the 5 ConvNets and keep genes that have averaged Jaccard similarity larger than 0.2 at least in one level of neighbor set sizes: 5, 10, 15, 20, 25, 30, 40, 50, or 100.

Correlated gene calling

To find significant correlated genes, we use the learned features from one of 5 CoSTA models to calculate Euclidean distance pairwise between all genes. For example, to get significant correlated genes with *Vim*, we calculated distances of all other SE genes to *Vim* based on the learned features. Then, we used these distances to create a null distribution. Distances that have Z-scores lower than -2.323 ($p < 0.01$) are considered significant, and genes that have significant distances would be called out as correlated genes to *Vim*. Because we trained 5 independent models, we obtain 5 sets of correlated genes for each SE gene in the data. Then, we keep correlated genes that show up in at least in 3 models.

MERFISH analysis

We obtained the MERFISH dataset collected on the mouse preoptic region of the hypothalamus from Dryad [15] (<https://datadryad.org/stash/dataset/doi:10.5061/dryad.8t8s248>), and we used the slice at Bregma +0.11 mm from animal 18 for analysis as used for SPARK analysis [5]. We reduced the image resolution tenfold and resized images to 85X85 matrices. Next, we directly applied a customized CoSTA model to the MERFISH dataset. This customized approach has the same general architecture that defines CoSTA, as described above. The customized ConvNet also has three convolutional layers but each convolutional layer has a larger filter, to reduce the overall size of the output. To compare with results from SPARK, we created null distributions for correlated gene calling by permuting images 100 times. Permuted images are forwarded through CoSTA to get permuted spatial features. Then we calculated their Euclidean distances with the spatial features of the true image, and these distances serve as the null distribution. Because the 9 defined cell type expression patterns are known, significantly correlated genes to these 9 expression patterns were called SE genes. For each gene in this MERFISH dataset, including the 5 blank controls, we calculated its Euclidean distances and its 100-time shuffled distances to the 9 expression patterns. If the true Euclidean distance of one gene to one cell type pattern are lower than Z-score -2.323 , we call this gene an SE gene that

is correlated to the expression pattern typical of this particular cell type. To visualize the training process, we project the feature vectors of each gene onto the first two UMAP dimensions and label each gene according to clusters defined using the whole feature vector. We use a linear intrinsic dimensionality (LID) estimator to quantify the change in cluster distinctness before and after training. This estimator mainly measures a ratio between distance of each datapoint to its the second closest datapoint and distance to its closest datapoint. Ratios are ordered from low to high and it fits a line that crosses the origin. The slope of this line represents the LID of this data in the latent space. Simply put, the lower LID, the more clustered datapoints are in the latent space. Indeed, among 10 different runs, spatial representations after training show lower LIDs than without training.

Analysis of slide-seq with SPARK and SpatialDE

Analysis of Slide-seq with SPARK and SpatialDE follows the standard analysis pipelines proposed by these two methods, with default parameters. Code of analysis can be found at the GitHub repository (<https://github.com/rpmccordlab/CoSTA>).

Abbreviations

ConvNet: Convolutional neural network; SE or SV gene: Spatially expressed or spatially variable gene; CoSTA: Unsupervised ConvNet learning strategy for spatial transcriptomics analysis.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12859-021-04314-1>.

Additional file 1. Supplementary Fig. 1: Comparison of CoSTA and overlap analysis performance in finding correlated digits to digit 3. 1000 images are sampled from the full MNIST dataset, and each digit contains 100 samples. CoSTA (red bars) uniquely calls samples of digit 3 as correlated to digit 3. However, overlap analysis (blue bars) identifies some instances of all digits as showing some overlap with digit 3. CoSTA is more specific, but less sensitive: CoSTA reports a smaller number of correlated digit 3 images (bottom right) while overlap analysis reports a greater number of correlated digits overall.

Additional file 2. Supplementary Fig. 2: Learning curves of CoSTA using true and shuffled synthetic datasets. 2,000 simulated gene matrices were used for each pattern, as in Fig. 1, with different levels of noise added ("variance"). Shuffling for each pattern and each simulated gene was performed identically so that pixelwise correlations were preserved but spatial relationships between neighbors were disrupted. (see Methods and Fig. S12 for shuffling approach details) NMI compares the clustering labels generated by CoSTA against the true class label.

Additional file 3. Supplementary Fig. 3: Performance of CoSTA using synthetic datasets with perturbations. A) Left: the same initial spatial patterns as in Fig. 1 were used. CoSTA was applied to classify 2,000 simulated genes for each pattern from the original patterns (top), half shuffled (middle), and fully shuffled (bottom) patterns. Applying trained CoSTA representations of simulated genes are visualized by using spatial representation in 2D UMAP. Genes are colored based on the true synthetic pattern from which they are derived. Silhouette scores quantify how well the representation distinguishes different patterns. (Closer to 1 = more distinguishable patterns are recovered). B) Disruption test through masking. Half of the simulated genes from each pattern have a masked region, simulating experimental missing data. The masked region is circled in blue in the upper panel. Representation of simulated genes based on pixelwise values (left) and features extracted by CoSTA (right) are visualized in 2D UMAP, and genes are colored based on pattern type from which they were generated (upper panel) or according to whether they belonged to the masked or unmasked set (lower panel).

Additional file 4. Supplementary Fig. 4: Training CoSTA with different ratios of SE and non-SE genes. (A) To simulate non-SE genes, five patterns without clear spatial features were used. (B) Simulated non-SE genes were mixed with SE genes simulated from the 5 patterns in Figure 1 in different ratios from 90:10 to 10:90. CoSTA representations of these gene mixtures are visualized in 2D UMAP. Genes are colored based on pattern membership (top) or SE type (bottom). (C) Silhouette scores quantify how well the representation distinguishes different patterns for SE and non-SE genes across different mixture ratios.

Additional file 5. Supplementary Fig. 5: The number of overlapped neighbors of Vim, Cttd, and Gfap before and after each weight updating across all epochs, considering either 10 nearest neighbors (left), 20 nearest neighbors (center), or 50 nearest neighbors (right).

Additional file 6. Supplementary Fig. 6: The number of SE genes after 3 days and 2 weeks brain injury. (A) Overlap of SE genes identified by SPARK or CoSTA. (B) Learning curve of CoSTA with original and shuffled data. (see Methods and Fig. S12 for shuffling approach details) Y-axis shows NMI calculated between cluster labels at training epoch t and cluster labels at previous epoch $t-1$. X-axis shows training epoch t . (C) Percent of all measured genes that are called SE genes by the 3 approaches.

Additional file 7. Supplementary Fig. 7: Expression patterns of Vim, Ctsd, and Gfap 3 days and 2 weeks after brain injury. (A) Expression patterns of Vim, Ctsd, and Gfap 3 days after brain injury. (B) Expression patterns of Vim, Ctsd, Gfap and genes involved in oligodendrocyte development (bottom row) 2 weeks after brain injury. Patterns that are visibly similar between Vim, Gfap, and Ctsd (small red boxes) are not strikingly visible in oligodendrocyte development genes.

Additional file 8. Supplementary Fig. 8: Comparison of SPARK, SpatialDE, CoSTA, and pixel overlap results. (A) Overlap of gene lists correlated with Vim, Ctsd, and Gfap at 2 weeks after injury identified by CoSTA, SPARK, SpatialDE, and overlap analysis ("Slide-seq"). (B) Examples of gene expression images for genes detected as similar to Vim, Gfap, and Ctsd by SPARK and also by CoSTA (left) or SPARK and not CoSTA (right). Numbers below images indicate the rank of the given gene in the list of correlated genes. See Figure S7 for expression patterns of the query genes. All images are scaled between 0 and 1 for visualization purposes. Key visible regions of high expression in Vim, Gfap, and Ctsd are circled in red for cross comparison of all images. (C) Gene Ontology term enrichment evaluated by Gorilla using the ranked correlated gene list produced by CoSTA (see Table S4). (D) Examples of gene expression images for genes highly ranked by CoSTA only (left), pixel only (middle), and both (right) as similar to Vim, Gfap, and Ctsd. Annotations next to gene names indicate rankings in CoSTA "C" and Pixel "P".

Additional file 9. Supplementary Fig. 9: Expression patterns of SE genes identified by CoSTA 2 weeks after brain injury. SE genes were clustered into 6 groups by SpatialDE and CoSTA. CoSTA cluster numbers correspond to Figure 4d and the most similar SpatialDE cluster is placed below the most closely corresponding CoSTA cluster when possible. The SpatialDE cluster containing Vim, Gfap, and Ctsd is cluster 2. Average expression pattern in 3rd row shows the overall pattern of all genes combined in the 2-week dataset.

Additional file 10. Supplementary Fig. 10: Effect of cluster number on CoSTA results with 2-week post injury Slide-seq data. a, SE genes identified by CoSTA with 10-100 clusters. b, CoSTA learning curve with 10-100 clusters. Y-axis shows NMI calculated between cluster labels at training epoch t and cluster labels at previous epoch $t-1$. X-axis shows training epoch t . c, Mean expression pattern of genes found to be correlated with Vim, Gfap and Ctsd identified by CoSTA with cluster numbers ranging from 10-100. Raw count values are scaled from 0 to 1 for these visualizations.

Additional file 11. Supplementary Fig. 11: CoSTA approach applied to clustering USPS, MNIST and Fashion datasets. Left panels: Models were trained for 10 epochs. After each weight updating, we clustered images into 10 clusters and directly compared them to true class labels through NMI. The grey line indicates clustering by UMAP+GMM with pixel values as features. The black line indicates clustering by SpatialDE. The orange line represents learning with combined center loss and bi-tempered logistic loss in Fashion dataset. Right panels: NMIs between clustering at the t th updating and the previous ($t-1$)th updating.

Additional file 12. Supplementary Fig. 12: Shuffling approach to preserve pixel correlation but disrupt spatial information. (A) Cartoon representing shuffling approach. Right: two initial 4x4 gene matrices (dimensions are small for example purposes). Each matrix shows a certain pattern of expression where a cluster of neighboring pixels show similar gene expression. Pixels are numbered so their positions can be tracked through the shuffling process. Middle: the 4x4 matrix is flattened into a single vector and then the positions of pixels are shuffled in the same way for Gene1 and Gene2 (orange arrows show a few example pixel rearrangements). The pixel ordering within each image is disrupted but each gene shares the same pixel ordering with other genes. This preserves individual pixel correlations across images from different genes but disrupts the spatial ordering and relationships between neighboring pixels. Right: shuffled vectors are reformed into a 4x4 matrix. (B) Example of shuffling result for 2 example Slide-seq gene matrices. Left: original gene expression image matrices. Shuffling is applied identically to the two genes as shown in A. Right: Resulting shuffled matrices. Visible spatial patterns are gone, but the pixel correlation of the two images would remain the same.

Additional file 13. Supplementary Table 1: Comparison of CoSTA and SpatialDE classification of 10,000 simulated genes belonging to 5 spatial patterns (see Fig. 1). Normalized Mutual Information is used to measure the similarity between CoSTA or SpatialDE-derived cluster assignments and true cluster assignments (values closer to 1 indicate a higher concordance between true and predicted cluster memberships). At noise level 0.6, CoSTA performs better with the addition of center loss (0.91 vs. 0.52). For shuffled data, each gene matrix was identically shuffled as described in Fig. S12 and in the Methods. This shuffling preserves pixel correlations between genes but disrupts overall spatial patterns, allowing an evaluation of whether the tested analysis detects pixelwise or spatial information.

Additional file 14. Supplementary Table 2: Clusters of SE genes identified by CoSTA in the MERFISH dataset (cell type patterns are included in clusters).

Additional file 15. Supplementary Table 3: CoSTA was applied to gene images from 4 different Slide-seq experiments and evaluated for whether it could separate gene images correctly into which original tissue slice (overall pattern) they came from. The table shows the confusion matrix of clustering labels derived from CoSTA results compared to the original known experimental label.

Additional file 16. Supplementary Table 4: All CoSTA SE genes in Slide-seq data 2 days after injury ranked according to their similarity to query genes Vim, Ctsd, and Gfap. Ranking according to CoSTA feature vectors as well as by pixel correlation are shown. Columns 4 and 5 indicate whether each gene was found to be correlated to these query genes by SPARK or SpatialDE.

Additional file 17. Supplementary Table 5: List of genes in each cluster derived by CoSTA from the 2-week Slide-seq data, as shown in Figure 4D.

Additional file 18. Supplementary Table 6: Runtime of CoSTA for 3-day and 2-week Slide-seq data. Runtimes are measured in minutes and under different numbers of clusters being assigned during training.

Acknowledgements

We thank Tian Hong, Tongye Shen, and Amir Sadovnik for insightful discussion.

Authors' contributions

YX conceived the project, developed the computational approach, and performed analyses. RPM supervised the project, performed some analyses and prepared some figures, and YX and RPM wrote the manuscript. All authors have read and approved the final manuscript.

Funding

This research was supported in part by NIH NIGMS grant R35GM133557 to R.P.M. The funding body played no role in the design of the study nor in the collection, analysis, and interpretation of data, nor in writing the manuscript.

Availability of data and materials

The processed Slide-seq datasets were retrieved from https://singlecell.broadinstitute.org/single_cell/study/SCP354/slide-seq-study. We also deposited processed MERFISH and Slide-seq data and scripts for all analyses in this study at the GitHub repository (<https://github.com/rpmccordlab/CoSTA>) under an Open Source Initiative compliant MIT license. The version of the code used in the manuscript is available at <https://doi.org/10.5281/zenodo.3948711>.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent to publish

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹UT-ORNL Graduate School of Genome Science and Technology, University of Tennessee, Knoxville, TN, USA. ²Department of Biochemistry & Cellular and Molecular Biology, University of Tennessee, Knoxville, TN, USA.

Received: 14 June 2021 Accepted: 2 August 2021

Published online: 09 August 2021

References

- Burgess DJ. Spatial transcriptomics coming of age. *Nat Rev Genet.* 2019;20(6):317.
- Eng CHL, Lawson M, Zhu Q, Dries R, Koulina N, Takei Y, Yun J, Cronin C, Karp C, Yuan G-C, et al. Transcriptome-scale super-resolved imaging in tissues by RNA seqFISH. *Nature.* 2019;568(7751):235.
- Rodrigues SG, Stickels RR, Goeva A, Martin CA, Murray E, Vanderburg CR, Welch J, Chen LM, Chen F, Macosko EZ. Slide-seq: a scalable technology for measuring genome-wide expression at high spatial resolution. *Science.* 2019;363(6434):1463–7.
- Valentine S, Sarah AT, Oliver S. SpatialDE: identification of spatially variable genes. *Nat Methods.* 2018;15(5):343–6.
- Sun S, Zhu J, Zhou X. Statistical analysis of spatial expression patterns for spatially resolved transcriptomic studies. *Nat Methods.* 2020;17(2):193–200.
- Caron M, Bojanowski P, Joulin A, Douze M. Deep clustering for unsupervised learning of visual features. 2018. [arXiv:1807.05520v2](https://arxiv.org/abs/1807.05520v2).
- McConville R, Santos-Rodriguez R, Piechocki RJ, Craddock I. N2D: (Not Too) deep clustering via clustering the local manifold of an autoencoded embedding. 2019. [arXiv:1908.05968v6](https://arxiv.org/abs/1908.05968v6).
- Xie J, Girshick R, Farhadi A. Unsupervised deep embedding for clustering analysis. In: Proceedings of the 33rd International Conference on Machine Learning. Vol 48, 2016; pp. 478–87. New York, NY, USA.
- McInnes L, Healy J, Melville J. UMAP: uniform manifold approximation and projection for dimension reduction. 2018. [arXiv:1802.03426v2](https://arxiv.org/abs/1802.03426v2).
- Amid E, Warmuth MK, Anil R, Koren T. Robust bi-tempered logistic loss based on bregman divergences. 2019. [arXiv:1906.03361v3](https://arxiv.org/abs/1906.03361v3).
- Addison M, Xu Q, Cayuso J, Wilkinson DG. Cell identity switching regulated by retinoic acid signaling maintains homogeneous segments in the hindbrain. *Dev Cell.* 2018;45(5):606–620.e603.
- Li D. The MNIST database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Process Mag.* 2012;29(6):141–2.
- Ståhl PL, Salmén F, Vickovic S, Lundmark A, Navarro JF, Magnusson J, Giacomello S, Asp M, Westholm JO, Huss M, et al. Visualization and analysis of gene expression in tissue sections by spatial transcriptomics. *Science (Am Assoc Adv Sci).* 2016;353(6294):78–82.

14. Vinh NX, Epps J, Bailey J. Information theoretic measures for clusterings comparison: variants, properties, normalization and correction for chance. *J Mach Learn Res*. 2010;11:2837–54.
15. Moffitt JR, Bambach-Mukku D, Eichhorn SW, Vaughn E, Shekhar K, Perez JD, Rubinstein ND, Hao J, Regev A, Dulac C, et al. Molecular, spatial, and functional single-cell profiling of the hypothalamic preoptic region. *Science*. 2018;362(6416):eaau5324.
16. Facco E, d'Errico M, Rodriguez A, Laio A. Estimating the intrinsic dimension of datasets by a minimal neighborhood information. *Sci Rep*. 2017;7(1):12140–8.
17. Edsgård D, Johnsson P, Sandberg R. Identification of spatial expression trends in single-cell gene expression data. *Nat Methods*. 2018;15(5):339–42.
18. Eden E, Navon R, Steinfeld I, Lipson D, Yakhini Z. GOrilla: a tool for discovery and visualization of enriched GO terms in ranked gene lists. *BMC Bioinform*. 2009;10(1):48–48.
19. Stickels RR, Murray E, Kumar P, Li J, Marshall JL, Di Bella DJ, Arlotta P, Macosko EZ, Chen F. Highly sensitive spatial transcriptomics at near-cellular resolution with Slide-seqV2. *Nat Biotechnol*. 2021;39(3):313–9.
20. Chen KH, Boettiger AN, Moffitt JR, Wang S, Zhuang X. Spatially resolved, highly multiplexed RNA profiling in single cells. *Science*. 2015;348(6233):aaa6090.
21. Wang X, Allen WE, Wright MA, Sylwestrak EL, Samusik N, Vesuna S, Evans K, Liu C, Ramakrishnan C, Liu J, et al. Three-dimensional intact-tissue sequencing of single-cell transcriptional states. *Science*. 2018;361(6400):eaat5691.
22. Vickovic S, Eraslan G, Salmén F, Klughammer J, Stenbeck L, Schapiro D, Åijö T, Bonneau R, Bergenstråhle L, Navarro JF, et al. High-definition spatial transcriptomics for in situ tissue profiling. *Nat Methods*. 2019;16(10):987–90.
23. Yang J, Parikh D, Batra D. Joint unsupervised learning of deep representations and image clusters. 2016. [arxiv:1604.03628v3](https://arxiv.org/abs/1604.03628v3).
24. Wen Y, Zhang K, Li Z, Qiao Y. A discriminative feature learning approach for deep face recognition. In: *Computer vision—ECCV 2016*. 2016. Cham: Springer; 2016. pp. 499–515.

Publisher's Note

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

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

