

# Interactive Visualization of Metric Distortion in Nonlinear Data Embeddings using the `distortions` Package

Kris Sankaran<sup>1</sup>, Shuzhen Zhang<sup>2</sup>, Chenab<sup>2</sup>, Marina Meilă<sup>3</sup>

<sup>1</sup>Department of Statistics, University of Wisconsin–Madison, Madison, 53706, WI, USA.

<sup>2</sup>Department of Statistics, University of Washington, Seattle, 98195, WA, USA.

<sup>3</sup>Department of Computer Science, University of Waterloo, Waterloo, N2L 3G1, Canada.

## Abstract

Nonlinear dimensionality reduction methods like UMAP and *t*-SNE can help to organize high-dimensional genomics data into manageable low-dimensional representations, like cell types or differentiation trajectories. Such reductions can be powerful, but inevitably introduce distortion. A growing body of work has demonstrated that this distortion can have serious consequences for downstream interpretation, for example, suggesting clusters that do not exist in the original data. Motivated by these developments, we implemented a software package, `distortions`, which builds on state-of-the-art methods for measuring local distortion and displays them in an intuitive and interactive way. Through case studies on simulated and real data, we find that the visualizations can help flag fragmented neighborhoods, support hyperparameter tuning, and enable method selection. We believe that this extra layer of information will help practitioners use nonlinear dimensionality reduction methods more confidently. The package documentation and notebooks reproducing all case studies are available online at <https://krisrs1128.github.io/distortions/site/>.

## 1 Background

Nonlinear dimensionality reduction methods like UMAP and *t*-SNE are central data visualization tools in modern biology. By projecting high-dimensional molecular profiles into lower dimensions, they reveal salient biological variation across cells. These methods support diverse applications, including developmental trajectory analysis, reference atlas construction, and disease characterization. They are included in widely used data analysis workflows like Scanpy [62] and Seurat [55] and have been popular in practice, reflecting their utility in modern biological research. Nonetheless, these methods have been controversial [5, 25, 32], because they can introduce distortions and artefacts. These shortcomings include exaggerating cluster differences, failing to capture density variation, and suggesting non-existent trajectories [10, 7, 60, 59, 17], which can complicate and cast doubts on the biological interpretation of the observed patterns, potentially leading to false discoveries.

Although alternative dimensionality reduction methods have been proposed that are arguably more principled, their adoption remains limited. For this reason, recent research has focused on wrapper methods designed to prevent artefacts and to support accurate embedding interpretation. These include improved method initialization [22], adaptations for visualization faithfulness [38, 29], automatic hyperparameter selection [63, 26], and statistical tests to flag problematic embedding regions [63, 26]. These methods provide valuable guidance for creating embedding visualizations. However, their static nature limits the amount of contextual information they can display. Nonlinear embedding distortions are local, direction-dependent –

055 stretching in some directions while contracting in others – and spatially variable, changing gradually from  
056 point to point or abruptly between clusters. This complexity makes it difficult for static visualizations to faithfully  
057 represent distortion context without inducing information overload. Moreover, while existing diagnostic  
058 methods can highlight problematic regions, the reasons underlying their selection (e.g., warped neighbor  
059 distances) must remain hidden to avoid visual clutter. Further, existing methods vary in their capacity to  
060 remove distortions or provide quantitative measures of the associated improvements.

061 To address these limitations, we introduce the `distortions` package, which uses interactive visualization  
062 to display the sources of distortion in nonlinear embeddings. We adopt a mathematically rigorous definition  
063 of local metric distortion, which is not tied to any particular data embedding algorithm [35], available in  
064 our package as `local_distortions()`. Our paper applies this measure to biological data for the first time.  
065 To render the rich information returned by `local_distortions()`, we introduce a version of the *focus-plus-context*  
066 principle [15, 46, 11], supporting the progressive and user-controlled disclosure of sources of  
067 distortion (like fragmented neighborhoods, defined below) based on user interaction, while maintaining the  
068 overall query context. This approach helps users interactively flag algorithmic artefacts and answer questions  
069 about them that are impossible to answer in full detail with static visualizations. Further, by introducing  
070 a new method for interactively isometrizing an embedding, we make it possible to obtain a distortion-free  
071 view of the underlying data's intrinsic geometry in the vicinity of the region of interest.

072 In summary, this paper makes the following contributions:

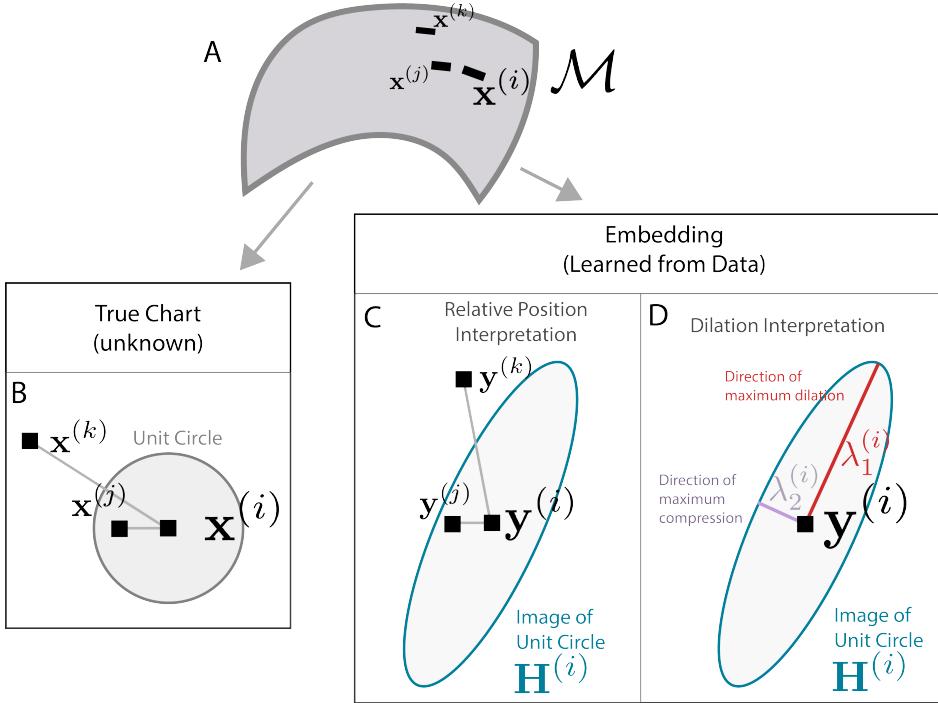
- 073
- 074 1. Applying state-of-the-art measures of local distortion from the manifold learning literature [44, 34, 35]  
075 to single-cell data for the first time. These methods reveal systematic differences in the interpretation of  
076 embedding distances across cell types and highlight contiguous neighborhoods that become fragmented  
077 during dimensionality reduction.
  - 078 2. Demonstrating the practical utility of distortion measures in choosing between algorithms and hyper-  
079 parameters. We find that these metrics support objective comparison of embedding results, and the  
080 accompanying visualizations provide insight into qualitatively different types of distortion.
  - 081 3. Developing interactive visualizations that highlight distorted regions and enable local corrections. We  
082 introduce an isometrization method that allows users to interactively correct distortions locally within  
083 regions of interest. Additional focus-plus-context approaches reveal distorted neighborhoods based on user  
084 queries of summary visualizations.

085 We validate this functionality using data with known low-dimensional structure, then apply the package to  
086 three single-cell datasets, showing the potential for improved biological interpretation and nonlinear embed-  
087 ding method application. The package is hosted at <https://pypi.org/project/distortions/> and documented  
088 at <https://krisrs1128.github.io/distortions/site/>.

## 089 1.1 Distortion estimation

090 To set up our results, we briefly review distortion estimation. Embedding methods aim to learn a low-  
091 dimensional, potentially nonlinear manifold on which the data lie. This manifold hypothesis is motivated by  
092 the fact that only certain patterns of gene expression are plausible, due to regulatory constraints. Geometrically,  
093 every point on the manifold can be mapped to a local coordinate system, called a *chart*. Biologically,  
094 the local coordinates are directions of shifting activity of latent biological processes. An ideal embedding  
095 method would perfectly recover these intrinsic charts, ensuring that distances on the biological manifold  $\mathcal{M}$   
096 are reflected in the embedding. Such a distance-preserving manifold embedding is called an *isometry*.

097 Even in linear dimensionality reduction, distances require careful interpretation. For example, in principal  
098 component analysis (PCA) plots, it is recommended that the axes be rescaled to reflect the relative  
099 variances explained by each component [39]. This issue becomes more difficult in nonlinear settings, where  
100 the interpretation of relative distances can vary locally across regions of the visualization [44]. Practical  
101 algorithms inevitably introduce distortion, systematically dilating some directions while compressing others.  
102 Depending on the direction of movement and the starting point, traveling the same distance in the  
103 embedding space might correspond to different distances along the manifold. Though we may not be able  
104 to avoid distortion, we can at least estimate it. Here we will call this estimate RMetric (Section 4.3 explains  
105 this name) and can be represented in various equivalent ways, as shown in Fig 1. For instance, the function  
106 `local_distortions()` returns RMetric as a matrix  $\mathbf{H}^{(i)}$ . The matrix  $\mathbf{H}^{(i)}$  gives a quantitative measure of



**Fig. 1:** Interpreting the matrices  $\mathbf{H}^{(i)}$  generated by the RMetric algorithm. A. Three points on a hypothetical manifold  $\mathcal{M}$ . B. The three points from panel A arranged on one of the charts that defines  $\mathcal{M}$ . A unit circle with respect to the intrinsic metric around  $\mathbf{x}^{(i)}$  is overlaid. C. The embedding algorithm distorts the unit circle from A. Though the distances and angles between samples have changed, the ratios of their distances to the unit circle have not. Though the true distortion around  $\mathbf{y}^{(i)}$  is unknown, it can be estimated using  $\mathbf{H}^{(i)}$  (blue ellipse). D. The same  $\mathbf{H}^{(i)}$  as panel C, but emphasizing the directions and degree of maximum dilation and compression.

local distortion induced by an embedding method. A mathematical treatment is provided in the Methods section, and we refer to [36] for an in-depth discussion.

We visually encode the local distortions  $\mathbf{H}^{(i)}$  with ellipses, displayed at each embedded point. Ellipses with circular shapes reflect regions where the embedding approximates an isometry. Thus the size and orientation of ellipses gives the principal directions of stretch/compression around point  $i$ , and the ellipse itself can be seen as a polar plot of the stretch (or compression) associated to each direction from point  $i$  (Fig 1). Specifically, larger ellipses appear when distances have been inflated and the major axes appear in the direction of most extreme dilation. This approach generalizes Tissot's indicatrix from cartography [23] to high-dimensional embedding algorithms.

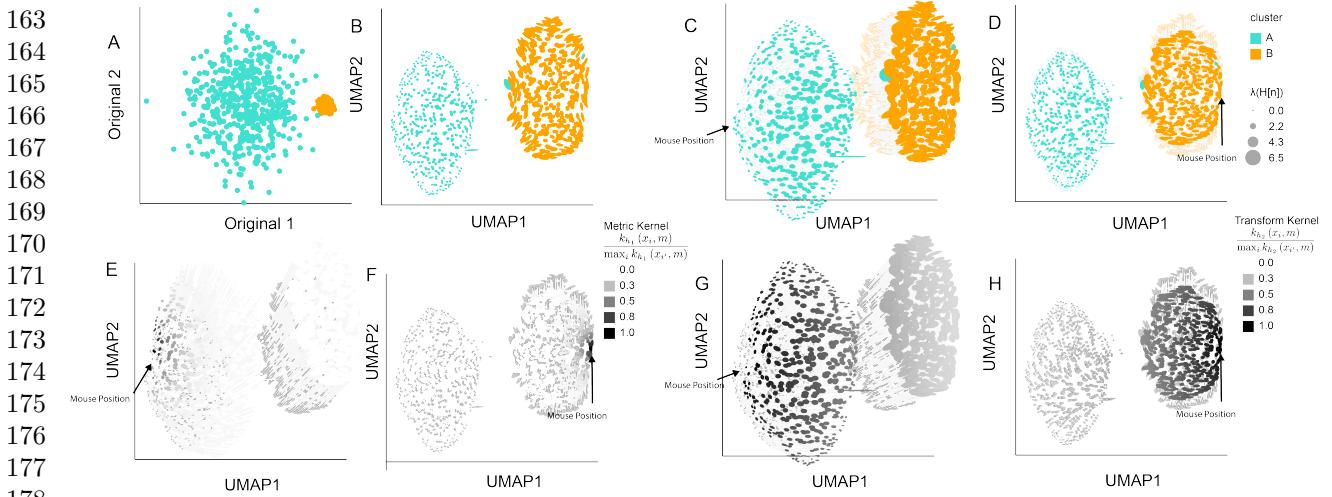
## 2 Results

### 2.1 Detecting cluster-specific differences in local metrics

This section gives two examples where local metric visualization highlights systematic differences in embedding interpretation across clusters.

**2.1.1 Gaussian mixtures with different variances** We evaluate the recovery of intrinsic geometry in an embedding of a mixture of two Gaussians. We sampled 500 points each from two components:  $\mathcal{N}(\mu_k, \sigma_k^2 I_2)$  where  $\mu_A = (0, 0)^\top$ ,  $\mu_B = (30, 0)^\top$ ,  $\sigma_A = 10$ , and  $\sigma_B = 1$ . The resulting mixture is shown in Fig 2A. We applied UMAP with 50 neighbors and a minimum distance of 0.5. Despite the large differences in variance, UMAP returned clusters with comparable sizes and densities (Fig 2B). We applied the RMetric algorithm with a geometric graph Laplacian constructed from the 50-nearest neighbor graph and rescaling  $\epsilon = 1$ . The affinity kernel radius was set to the mean of the original data distances between neighbors on this

109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162

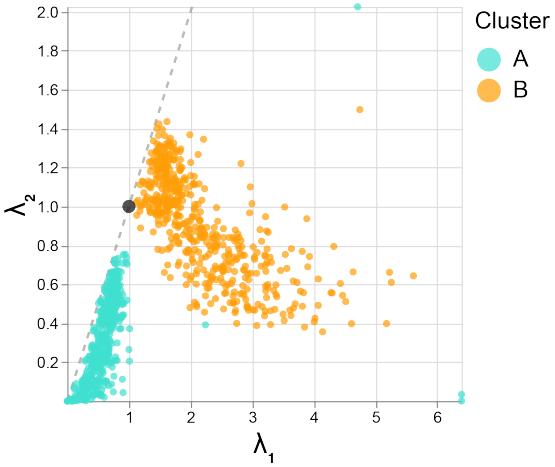


**Fig. 2:** Interactive isometrization partially restores density differences in an Gaussian mixture embedding. A. Original simulated data. Cluster B has smaller variance compared to Cluster A. B. Ellipse orientation and sizes encode differences in local metrics in the UMAP embedding. Smaller ellipses mean that the same distance in the embedding space corresponds to larger distances in the original data space. C. The isometrization interaction updates ellipse size and positions to reflect the local metric in the hovered-over region. This partially restores the difference between cluster variances that were lost in the initial embedding. D. The analogous isometrization when hovering over Cluster B (orange). Cluster B slightly shrinks, while Cluster A (blue) remains at its original size. E. The normalized kernel similarities defining the contribution of each  $\mathbf{H}^{(i)}$  to the  $\mathbf{H}^*$  used in the isometrization from panel B, as given in equation (4). F. The analog of panel E for the mouse interaction in Panel D. G. The normalized kernel similarities describing the extent to which each point is moved from its original position, as given in equation (5). The analog of panel G for the mouse interaction in panel D.

graph. To prevent samples with outlying  $(\lambda_1^{(i)}, \lambda_2^{(i)})$  from obscuring variation among the remaining points, we truncated  $\lambda^{(i)}$  at a maximum value of 5; this affects 6 samples. To ensure that isometrization does not uniformly contract or expand neighborhoods across the visualization, we further divided all  $\mathbf{H}^{(i)}$  by a scaling factor  $\frac{1}{4N} \sum_{i'} \sum_{k,k'} \mathbf{H}_{kk'}^{(i')}$ .

The resulting local metrics  $\mathbf{H}^{(i)}$  are overlaid as ellipses in Fig 2B-H. Fig 2B shows that Cluster A has smaller ellipses than Cluster B, correctly reflecting the differences in cluster variance lost by the UMAP embedding. Fig 3 shows the coordinates of the truncated  $\lambda^{(i)}$  plotted against one another. The clear separation in singular values across clusters reinforces the qualitative differences in ellipse sizes from Fig 2A. Fig 2C-D show the isometrized versions of Fig 2B when hovering over samples in Cluster A and B, respectively. These interactions recalculate the embedding locations and ellipse sizes to bring the local metrics  $\mathbf{H}^{(i)}$  near the viewer's mouse position closer to the identity  $I_2$ , resulting in more circular ellipses. Thin grey lines connect the isometrized and the original embedding coordinates. When hovering over Cluster A, the samples in that cluster become spread further apart, while those in Cluster B are translated to the right but remain at their original density. In contrast, when hovering over Cluster B, the samples in that cluster contract while those in Cluster A remain close to their original positions.

More precisely, Fig 2C calculates a “local” metric  $\mathbf{H}^*$  based on the weights in Fig 2E, which are high (darker) near the viewer's mouse position. An exact isometrization with respect to the current region of interest would update embedding coordinates  $\mathbf{y}^{(i)}$  to  $(\mathbf{H}^*)^{-\frac{1}{2}} \mathbf{y}^{(i)}$  [44] across the entire visualization. We instead restrict the transformation to areas close to the viewer's current interaction region. Informally, the darker points in Fig 2G are allowed to be updated more aggressively than the lighter points; the formal transformation is detailed in equation (5). The analogs of Fig 2E and G for the interaction in Fig 2D are given in Fig 2F and H.



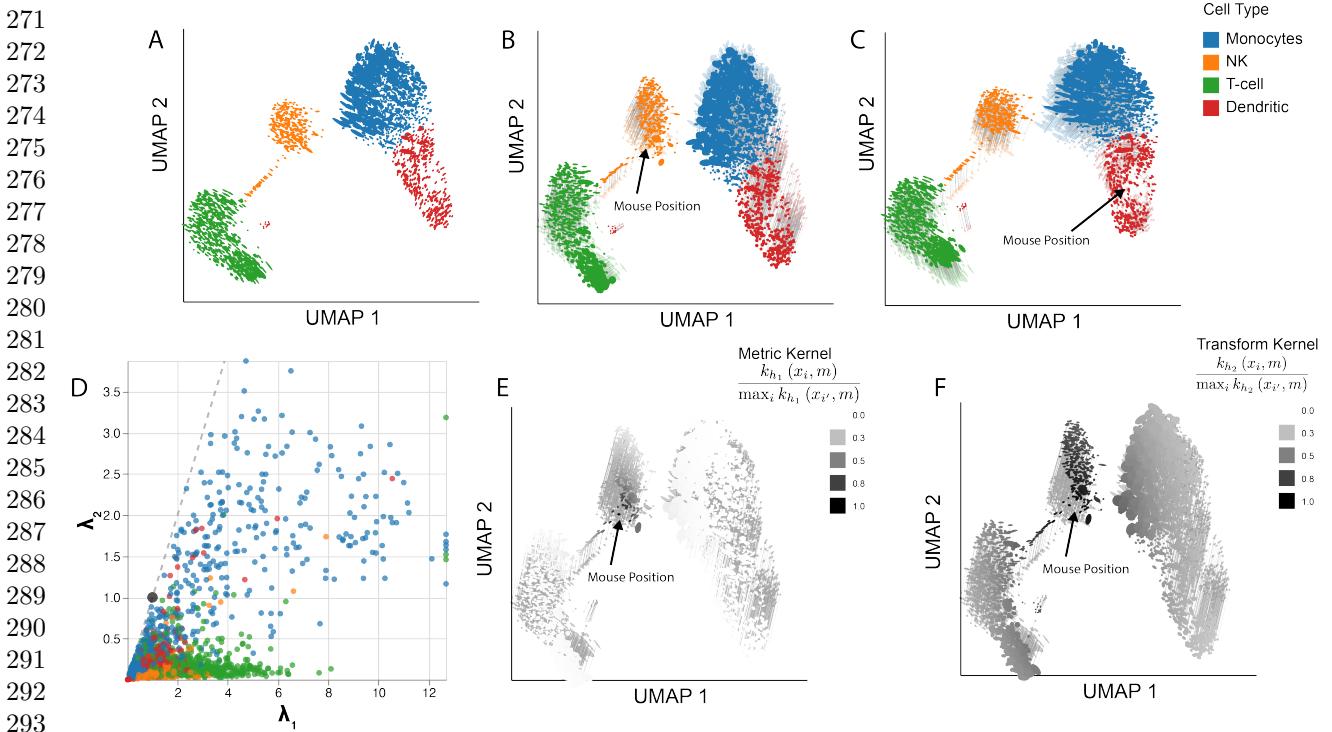
**Fig. 3:** Singular values  $\lambda_1^{(i)}$  and  $\lambda_2^{(i)}$  of the  $\mathbf{H}^{(i)}$  estimated in Figure 2. The larger  $(\lambda_1^{(i)}, \lambda_2^{(i)})$  in Cluster B results in the larger ellipse sizes for that cluster, indicating that embedding distances for this cluster have been “spread out” relative to original, pre-embedding distances. This effect is consistent with the data shown in Figure 2A.

**2.1.2 Local metrics vary across cell types in a PBMC atlas** We next analyze peripheral mononuclear blood cell (PBMC) single-cell genomics data (2683 cells, 1838 genes) from 10X Genomics, with default processing from scanpy [50, 62], which included total sum scaling (TSS), a  $\log(1+x)$  transformation, and highly variable gene filtering. We then applied UMAP (50 neighbors, minimum distance 0.5) to the PCA-denoised data (top 40 components). Cell types were identified with Leiden clustering and canonical marker genes CD79A and MS4A1 (B cells), FCER1A and CST3 (dendritic cells), GNLY and NKG7 (NK cells), FCGR3A (monocytes), IGJ (plasma cells), and CD3D (T cells). To estimate the data’s intrinsic geometry, we applied RMetric using the geometric graph Laplacian constructed from the 50-nearest neighbor graph and rescaling  $\epsilon = 5$ . The affinity kernel radius was set to three times the mean of the original data distances within this graph. To prevent a few highly skewed ellipses from influencing the remaining points, we truncated the singular values  $(\lambda_1^{(i)}, \lambda_2^{(i)})$  from above at 2.5. We divide by the same  $\frac{1}{4N} \sum_{i'} \sum_{k,k'} \mathbf{H}_{kk'}^{(i')}$  scaling factor as in Gaussian mixture example above.

The resulting ellipse-enriched embedding Fig 4A reveals systematic metric differences across cell types. T cells ellipses are oriented with major axes in the northwest/southeast direction, suggesting that distances orthogonal to this direction compressed in the embedding. In contrast, dendritic cells are generally oriented in the southwest/northeast direction, suggesting greater spread away from the monocytes than the embedding alone indicates. Fig 4D displays the truncated and rescaled singular values  $(\lambda_1^{(i)}, \lambda_2^{(i)})$ . Points closer to the  $x$ -axis correspond to ellipses that are more eccentric than those near the center of the plot. Cell types differ systematically in this view as well, reinforcing our conclusion that local metrics are associated with cell type. The panel also draws attention to the high condition numbers among subsets of the T and NK cells. In contrast, many monocytes lie in the middle of the panel; these are the more circular embeddings in Fig 4A. Further, zooming into Fig 4D, Appendix Fig A1 reveals a second monocyte subset with smaller singular values. This pattern matches the bimodality in monocyte size distribution in Figure 4A. The UMAP embedding appears to have collapsed the two monocyte subsets, compressing distances for smaller points and dilating them for larger points. Though changing the visual markers from circles to ellipses is a small difference, the associated local metrics reveals valuable context about how UMAP warps intrinsic geometry across the visualization.

Fig 4 illustrates isometrizations for two cell types. Fig 4B - C show the embedding after placing the mouse over the cluster of NK cells (Fig 4B) and dendritic cells (Fig 4C), respectively. In both panels, the solid ellipses represent the updated embedding, while transparent ellipses and thin lines indicate the original positions and metrics. Isometrization over the NK cells expands the main NK cluster and increases the distance from the main cluster and the subset bridging T cells and NK cells, consistent with the northwest/southeast

217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270



**Fig. 4:** Isometrization of the PBMC UMAP embeddings. A. Ellipse orientation and sizes vary systematically across regions of the embedding, indicating differences in local metrics within and between cell types. B. An updated version of panel A when the mouse is positioned over a subset of NK cells. Transparent ellipses mark the cells' original positions, and lines connect the original and updated locations. C. Isometrization when hovering over dendritic cells. D. The windsorized singular values  $\lambda_1^{(i)}, \lambda_2^{(i)}$  associated with  $\mathbf{H}^{(i)}$  across cells. Ellipse size is determined by  $\lambda_1^{(i)}\lambda_2^{(i)}$  and eccentricity by  $\lambda_1^{(i)}/\lambda_2^{(i)}$ . A version that zooms into the region near the origin is given in Appendix Fig A1. E. The normalized kernel similarities defining the local metric  $\mathbf{H}^*$  (equation (4)) when the mouse is placed as in panel C. F. The analog of panel E when the mouse is placed as in panel G. The normalized kernel similarities defining the region of transformation (equation (5)) when the mouse is placed as in panel C. H. The analog of panel G when the mouse is placed as in panel D.

304  
305  
306

307 orientation of ellipses in Fig 4A. Figs 4E - F display the normalized kernel similarities used to define the local  
308 metric  $\mathbf{H}^*$  and the regions of transformation. The interaction over the dendritic cells (Fig 4C) increases the  
309 spread of cells close to the mouse position. Monocyte orientations are shifted slightly, but other cell types  
310 remain largely unchanged. Together, this suggests that the distortion of NK cells is more severe than that  
311 of dendritic cells in this choice of UMAP embedding.  
312

## 313 2.2 Identifying fragmented neighborhoods

314

315 There is emerging evidence that nonlinear dimensionality reduction methods can introduce embedding  
316 discontinuities [26], meaning that some points that are nearby in the original space end up being embedded  
317 as far from one another as those that are originally very different. In particular, points that lie in the same  
318 neighborhood in the original space may be fragmented into different embedding regions, complicating the  
319 interpretation of between-cluster relationships. To address this, *distortions* provides metrics for quantifying  
320 fragmentation at the neighborhood and pair levels, based on the relationship between observed vs. embedding  
321 distances among nearby points in the original data space, as detailed in the Methods (“Identifying fragmented  
322 neighborhoods”). A focus-plus-context visualization approach [15, 24, 46] then allows viewer interactions to  
323 progressively reveal the extent of fragmentation within different embedding regions. We provide examples  
324 below.

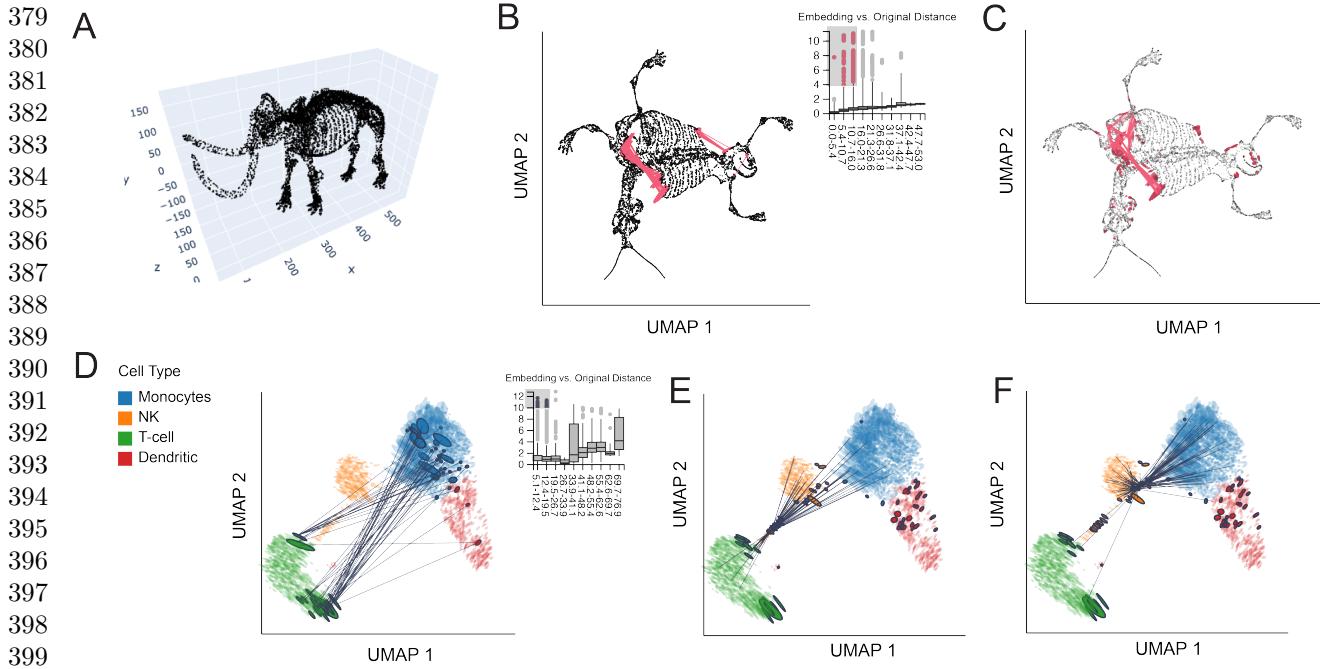
**2.2.1 Mammoth skeleton** We evaluate this strategy on UMAP embeddings of a three-dimensional mammoth skeleton point cloud (Fig 5A) generated by the Smithsonian Museums in an effort to digitize their collection [40]. This dataset has been used to study the artifacts introduced by UMAP [7]. It has the advantage of being directly visualizable in three dimensions (Fig 5A). Further, the data exhibit patterns at both global and local scales. For example, a successful dimensionality reduction method must preserve global relationships, like the relative positions of tusk, skull, and legs, and also fine-scale differences, like the distinction between bones in the rib cage. We applied UMAP (50 neighbors, minimum distance 0.5) to embed the 10,000 samples available in these data. The RMetric algorithm was applied using a geometric graph Laplacian constructed from the 50-nearest neighbor graph and a rescaling  $\epsilon = 5$ . The affinity kernel radius was set to 3 times the mean of the original data distances between neighbors on this graph. The resulting  $\mathbf{H}^{(i)}$  are directly encoded using ellipse dimensions without any post-processing. To identify distorted pairs, we used the boxplot display with outlier threshold set to  $10 \times \text{IQR}$ . To identify distorted neighborhoods, we applied the bin-based screening metric (see Methods) with  $\kappa = 0.1$  and  $\sigma = 3$ , requiring that at least 10% of neighbor distances be poorly preserved. This flags 425 potentially fragmented neighborhoods.

Fig 5B shows the boxplot widget overlaid on the UMAP. Reassuringly, the median embedding distances increase monotonically as the distances in the original space increase. However, within each bin, the distribution of embedding distances is skewed, especially for small distances in the original data space. Many pairs of points within these bins appear much further apart in the embedding than expected. In the current display, the viewer has selected the outliers within the three leftmost bins, highlighted in pink. The corresponding pairs are linked together in the main embedding view. These pairs include points on the left and right shoulders of the mammoth. These points are close to one another in three dimensions, but have been spread apart by the embedding. UMAP appears to reflect geodesic rather than Euclidean distance, effectively “flattening” the mammoth skeleton. In addition to the left and right shoulder pairs, the highlighted outliers include neighbors where one point lies on the last right-side rib bones and the other on the right side of the pelvis. UMAP embeds these adjacent bones further apart than appropriate, another distortion of the original structure.

The fragmented neighborhoods displays in Fig 5C confirms these findings. For example, the flattening of the shoulder is evident in the chain of fragmented neighborhoods in this region. Points further along the rib on the right and pelvis are also highlighted, as in the boxplot view. In this case, the viewer’s mouse lies over the right shoulder of the mammoth. Unlike the boxplot view, this allows us to view all the neighbors of distorted points near the mouse, showing the neighbors along the chest and arm whose distances are not outlying in the boxplot. This view also reveals more isolated fragmentations, for example on the skull, arms, and tail. Hovering over these points shows that a large fraction of their neighbors have also been spread apart (e.g., left and right hand sides of the skull), even if their absolute embedding distances are not large enough to stand out in the boxplot interactions.

**2.2.2 PBMC gene expression** We next identify distorted embedding pairs in the PBMC example. The boxplot widget again reveals outliers when the original distances are small (Fig 5D). The viewer has brushed the top outliers within the two leftmost bins. This selection highlights pairs of T cells and monocytes that are close despite the apparent embedding separation. This view distinguishes between two regions with distorted pairs within the T cell cluster. Unlike the NK cells, which visibly cluster into two subtypes, these two subtypes of distorted T cells do not stand out from the main T cell cluster. Nonetheless, both subtypes are near the boundary of the overall cluster. This suggests that in high dimensions, the T cell cluster may be curved in a way that allows these subgroup to be closer to the monocytes than is visible in the embedding.

We next identify fragmented neighborhoods using the bin-based strategy ( $L = 10$ ,  $\kappa = 0.2$ ,  $\sigma = 2$  threshold), resulting in 72 cells with fragmented neighborhoods. Fig 5E highlights fragmented NK cell neighborhoods that are separated from the main NK cluster. These distorted neighborhoods are centered on cells that are connected to both T cells and monocytes. These cells appear to bridge several cell types. Fig 5F shows a different subgroup of distorted NK cells. These cells lie along the NK cluster periphery, with many links to monocytes but few to T cells. This suggests that the sharp separation between the NK cells and monocytes may be an embedding artifact. A subset of dendritic cell are flagged as distorted, and hovering over them shows that they are neighbors with distant monocytes. Two subtypes of T cells are flagged as distorted; these largely overlap with those highlighted by the boxplot visualization in Fig 5D. Only four monocytes with fragmented neighborhoods are flagged, suggesting that this cluster does not suffer from



**Fig. 5:** Fragmented neighborhoods and links. A. The original mammoth point cloud, before applying any dimensionality reduction. B. Pairs with poorly preserved distances in the mammoth data. The viewer has selected pairs of points that are close to one another in the original space, but which are far apart in the embedding. C. Analogous poorly fragmented neighborhoods defined using the quantile smoothing criteria. D. Pairs with poorly preserved distances in the PBMC data. Distances between NK cells, T cells, and monocytes have been exaggerated by the UMAP embedding. E. A subgroup of NK cells with poorly preserved neighborhoods in the PBMC data. Cells close to the viewer's mouse interaction have neighbors spread across T cells, NK cells, and monocytes. F. A different subgroup of NK cells with fragmented neighborhoods. These cells often have neighbors lying on the far boundary of monocytes.

409

410

411 fragmentation as severely as the others. Interactive distortion visualization can reveal different degrees and  
412 types of distortion across and within cell types.  
413

414

415

### 2.3 Guiding method selection and tuning

416

417

In addition to interpreting individual embedding visualizations, distortion metrics can be used to compare different embedding methods and hyperparameter choices. They give a quantitative way to judge how well competing embeddings preserve the original data's structure. Further, the interactivity implemented in `distortions` makes it possible to explore why distortions arise, without overwhelming viewers with all contextual information at once. In this section we use three example datasets to illustrate how distortion visualization can guide method selection and tuning.

422

423

**2.3.1 Clarifying how hyperparameter choice impacts distortion** Hyperparameters in nonlinear dimensionality reduction methods like UMAP and *t*-SNE can substantially influence results [22, 3, 60]. Distortion visualization can reveal the trade-offs imposed by specific choices. We evaluate this using the hydra cellular differentiation data from [53]. This study used single-cell RNA sequencing to measure gene expression of a developing hydra polyp, an organism notable for its regenerative ability. Fig 1 of their paper is a *t*-SNE that clarifies the cellular composition of hydra tissue as well as the differentiation paths from stem and progenitor cells to specialized cell types. To create a setting with greater statistical instability and where hyperparameters may play a more important role, we take a random sample of 2000 of the original 24,985 cells. As in the analysis of the PBMC data, we apply TSS normalization, a  $\log(1 + x)$  transformation, and filter to the top 1000 highly variable genes.

We apply  $t$ -SNE to the PCA denoised data (30 components) with perplexity values of either 80 or 500. To estimate distortion we use the RMetric algorithm with a geometric graph Laplacian with 50 nearest neighbors and a rescaling  $\epsilon = 5$ . The radius for the affinity kernel was set to three times the average original data distance in the 50-nearest neighbor graph. Both the boxplot widget and the neighborhood fragmentation visualizations suggest qualitatively different types of distortion across the two perplexity settings. At a perplexity of 80, the fragmented neighborhoods occur in the gaps between cell type clusters (Fig 6A). Hovering over these neighborhoods reveals connections to adjacent cell types (Fig 6C), suggesting that some transitions in gene expression programs between cell types may in fact be more gradual. These blurrier transitions are captured at a perplexity of 500 (Fig 6B). However, at this hyperparameter value, many fragmented neighborhoods appear along the top and bottom boundaries of the embedding. Interacting with the display reveals that at this hyperparameter choice, the embedding fails to preserve distances between peripheral neighborhoods. For example, the viewer's selection in Fig 6D highlights neighbors that have been split across opposite sides of the visualization.

This reveals a basic trade-off: higher perplexity better reflects distances between main cell types but arbitrarily places rarer types, while lower perplexity correctly places these rare clusters correctly at the cost of inflating distances between common cell types. This additional context gives confidence in the conclusions drawn within specific regions of separate visualizations. These conclusions can still be reliable even when no single view preserves all relevant properties of the original high-dimensional data. Further, though the qualitative differences between hyperparameter choices would be difficult to obtain through manual inspection of the distances within the embedding output, the interactive display allows the differences to pop out naturally.

**2.3.2 Comparing initialization strategies using distortion metrics** Nonlinear dimensionality reduction methods can be sensitive to initialization strategies. Indeed, most single cell analysis packages use a preliminary dimensionality reduction step, like PCA or Laplacian eigenmaps [1], to initialize the optimization [22, 48]. We next study whether distortion metrics can detect issues arising due to poor initialization. To this end, we rerun the UMAP analysis of the PBMC data and consider a random, rather than the default spectral, initialization. All other dimensionality reduction and visualization hyperparameters remain as before. Fig 7 presents the results. Compared to Fig 4A, Fig 7A separates the NK cells into distinct groups falling on opposite sides of a main cluster of dendritic cells and monocytes. Brushing outlying neighbor distance pairs in the boxplot in Fig 7C highlights the fact that these two groups share many neighbors, and that the gap is artificial: many NK cells are neighbors with T cells despite lying on opposite sides of the plot. This suggests that the spectral initialization, which places T cells and NK cells adjacent to one another, better preserves their neighborhood relationships.

Fig 7D displays the fragmented neighborhoods, analogous to Fig 5E. Though some T-cell-adjacent NK cells had been flagged in the spectrally-initialized embedding, a larger number are distorted in the random initialization, including many with neighbors in the monocyte cluster. Further, the reduced  $y$ -axis range in Fig 7B relative to Fig 4D draws attention the greater eccentricity of ellipses in the random initialization, indicating larger distortion of local metrics. Importantly, none of these issues with the random initialization are detectable from the embedding coordinates alone. Both ellipse eccentricity and interaction with distortion summary metrics add context for understanding the importance of effective UMAP initialization.

**2.3.3 Analyzing density preservation in a *Caenorhabditis elegans* cell atlas** We applied our package to a single-cell atlas of *Caenorhabditis elegans* development [41], originally gathered to characterize the gene programs activated during different phases of embryogenesis in the *C. elegans* model system. These data include measurements on 86,024 cells, of which 93% have been manually annotated with cell types by the authors. Nonlinear embeddings applied to this dataset are known to obscure meaningful differences in local density, causing biologically meaningful cell types to appear sparser or denser than appropriate [38]. Therefore, we compared UMAP with the density-preserving algorithm DensMAP and use `distortions` to evaluate the improvement in local metric preservation. By utilizing our package's distortion summaries, we can highlight neighborhoods that are artificially fragmented in the embeddings and quantify the reduction in distortion made possible through the DensMAP algorithm.

Before applying either method, we applied a PCA denoising step which reduced the data to 100 dimensions. We applied both UMAP and DensMAP with 10 neighbors and a minimum distance of 0.5. To simplify the distortion analysis, we considered a random sample of 5000 cell embeddings from each of 10 randomly

487 chosen cell types (arcade, glia, pharyngeal neuron, intestinal and rectal muscle, M, excretory duct and pore,  
488 hypodermis, intestine, and rectal gland). This restriction is analogous to focusing on a subset of cell types  
489 when testing whether putative cell types are truly distinct or a visualization artifact [54, 12]. We used RMetric-  
490 ric to estimate local metric distortion using a geometric graph Laplacian based on the 10-nearest neighbor  
491 graph and affinity kernel radius set to three times the average original neighbor distance in this graph. To  
492 identify distorted neighborhoods associated with each method, we apply the bin-based strategy with  
493  $L = 10, \kappa = 0.4, \sigma = 3$  to flag points where a fraction of at least 40% of neighbors have embedding distance  
494 at least  $3 \times \text{IQR}$  away from the median within the corresponding bin of original distances.

495 Fig 8A-B shows the resulting fragmented neighborhoods. In both embeddings, the degree of fragmentation  
496 varies by cell type. For example, pharyngeal neuron neighborhoods are often fragmented by both algorithms,  
497 while few fragmented neighborhoods are centered on hypodermis cells. Qualitatively, the DensMAP embed-  
498 ding is less compressed into tight clusters than UMAP, suggesting that UMAP may artificially inflate the  
499 embedding space densities. Despite using the same graphical encoding scales, the UMAP ellipses also appear  
500 to be less uniformly sized. The more compact “hair” plots reinforce this conclusion (Fig 8D-E). Each seg-  
501 ment corresponds to one ellipse in panels A - B. The segments are oriented along the minor axis of the  
502 ellipses, and their lengths encode condition number  $\lambda_1^{(i)} / \lambda_2^{(i)}$ . We note that these hair-like graphical marks  
503 can be substituted for ellipses in all visualizations and interactions discussed above, including the boxplot  
504 and isometrization displays.

505 Further, the distortion metrics provide quantitative support of DensMAP’s ability to preserve intrinsic  
506 geometric information. For example, the histogram in Fig 8C shows that the UMAP resulted in systematically  
507 larger metric condition numbers, suggesting more systematic metric distortion. Further, Fig 8F shows that  
508 across choices of  $\kappa$ , the DensMAP results in fewer fragmented neighborhoods than UMAP. In this case, the  
509 distortion metrics led to a stable conclusion across hyperparameters. However, more generally, whether a  
510 neighborhood is flagged as fragmented can be dependent on the choices ( $L, \kappa, \sigma$ ), and local metric estimates  
511 like  $\lambda^{(i)}$  can depend on the graph Laplacian neighborhood and radius hyperparameter choices. We have  
512 followed the recommendations discussed in [44], but it is worth considering that the estimated degree of  
513 distortion can be dependent on hyperparameter choices.  
514

515

## 516 2.4 Software architecture and extensibility

517

518 Considering that no single definition of distortion exists for nonlinear dimensionality reduction, the  
519 `distortions` package adopts a “loosely coupled” design to ensure extensibility [14]. Each visualization  
520 accepts viewer-provided specifications of fragmented neighborhoods or links. Alternative distortion metrics  
521 can be implemented in independent functions as long as their formats are consistent. Similarly, visualizations  
522 can composed from viewer-specified graphical marks and interactions, similar in spirit to `ggplot2` [61] and  
523 `altair` [57]. For example, consider the interactions with fragmented neighborhoods of the PBMC data in Fig  
524 5E - F. If the distorted neighborhoods are stored in a dictionary `N`, then the interactive plot can be created  
525 with  
526

```
527 dplot(embedding) \
  .mapping(x="embedding_0", y="embedding_1", color="cell_type") \
  .geom_ellipse() \
  .inter_edge_link(N=N)
```

528

529 and the result will appear in a jupyter notebook cell.

530 This loose coupling also simplifies the application of our visualization to other distortion summarization  
531 approaches. We illustrate this by using the scDEED algorithm [63] to flag dubious cells in a t-SNE embedding  
532 of the PBMC data (Fig 9A). This figure was generated by applying the default scDEED workflow to the  
533 PBMC data, adding local metrics  $\mathbf{H}^{(i)}$  to the resulting embeddings, and replacing `embedding` and `N` in the  
534 call with the corresponding scDEED output. The resulting visualization is consistent with expectations –  
535 by hovering the mouse close to the scDEED flagged cells, we see these cells often have neighbors in the  
536 original space that are placed far apart in the embedding space (Fig 9B). We note that no such interactive  
537 display has previously been available for output from the scDEED package.  
538

We can also customize the graphical marks, styling, and labels in a format familiar to ggplot2 and altair users. For example, the visualization from the code block above can be customized using

```
dplot(embedding, width=440, height=340)\# custom plot size  
  .mapping(x="embedding_0", y="embedding_1", color="cell_type")\#  
  .geom_ellipse(radiusMax=15, radiusMin=1)\# custom point size  
  .inter_edge_link(N=N, threshold=.1, strokeWidth=0.4)\# narrower interaction window  
  .scale_color(legendTextSize=15)\# increase legend size  
  .labs(x="UMAP 1", y="UMAP 2")# custom labels
```

Further, we can switch from the fragmented neighborhood to the boxplot interaction (Fig 5D) by simply substituting the `inter_edge_link` call with `inter_boxplot`. This modular approach also enables the specification of new graphical marks and interactivity. For example, for large datasets, ellipses can be replaced with line segments, as in Fig 8D - E. This more compact encoding of local distortions is accomplished by substituting the `geom_ellipse` mark with `geom_hair`.

### 3 Summary and conclusions

Nonlinear embedding visualizations have been essential to progress in high-throughput biology, offering visual overviews that have guided advances in diverse applications like cell atlas construction [9, 45], cell differentiation trajectories [4, 13], and functional diversity mapping in metagenomes [51]. However, their potential for misinterpretation is well-documented [10, 60, 5, 22]. The community has made significant progress in characterizing and minimizing distortion [38, 63, 26, 29, 18], and the `distortions` package offers an interactive visualization toolbox that draws from manifold learning concepts and complements these advances.

Moderate distortions are accurately characterized by the RMetric algorithm, whose results can be graphically encoded in ellipse or hair plots, while more severe distortions are flagged via fragmented neighborhood plots. RMetric emphasizes how the intrinsic geometry is warped across different regions of the embedding space, alerting analysts to failures in density preservation and compression/dilation in certain embedding directions (Fig 2). Further, by flagging fragmented neighborhoods, we could identify clusters that are more closely related in the original data space than the embedding suggests. In the PBMC fragmentation example (Fig 5D-F), we found that a subset of T cells had many neighbors coming from the monocyte cluster, despite these clusters appearing on the opposite sides of the embedding visualization. Further, interaction with distortion metrics highlighted trade-offs between the types of distortion introduced by different hyperparameter choices (Fig 6). Finally, local isometrization offers the scientist a kind of magnifying glass into the local geometry of the original data, making it possible to zoom in and query low-level sample relationships that can be lost in global reductions.

We acknowledge limitations in our approach. For instance, our summaries depend on viewer-specified hyperparameters, like the number of bins  $L$  or neighborhood fraction  $\kappa$  in the bin-based fragmented neighborhood definition. Regarding distortion, while the local distortion RMetric is a well-defined differential geometric quantity, measuring distortions at larger scales is open to subjective preferences. For example, the fragmented neighborhood definition relies on distances in the original high-dimensional space, while alternatives could consider geodesic distances along the original manifold embedded in the high-dimensional spaces. Pairs of points could be flagged by outlierness, or by constructing a neighborhood graph in the embedding space and comparing the two.

Future work should ensure that a variety of long-range distortion measures are covered. Our modular software architecture will support straightforward extensions to new definitions and visualization layers. We expect that continued effort in this space will result in visualization techniques that can allow analysts to gain valuable insights from exploratory overviews while contextualizing their inherent limitations. While nonlinear dimensionality reduction methods cannot fully preserve all metric properties from the original data space, these exploratory views can guide more appropriate interpretation, allowing scientists to communicate results confidently and avoid the pitfalls of false discoveries due to algorithmic artefacts. By overlaying quantitative summaries of the distortion introduced by embedding algorithms, the `distortions` package aids researcher

595 intuition and facilitates critical evaluation of the embedding visualizations that have become standard in  
 596 modern biological analysis.

597

## 598 4 Methods

599

### 600 4.1 Notation

601

In the following, matrices will be denoted in bold uppercase letters, e.g.  $\mathbf{A}$ , vectors in bold lowercase, e.g.  $\mathbf{v}$ , vector and matrix elements by additional subscripts, e.g.  $\mathbf{A}_{ii'}$ , and other scalars by unbolded Latin and Greek letters. The index  $i$  will be reserved for denoting the  $i^{th}$  data point, and it will be used as a superscript on vectors and matrices associated with it. Thus, the original data is  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i)}, \dots, \mathbf{x}^{(n)} \in \mathbb{R}^D$ , where  $n$  is the sample size and  $D$  is the dimension of the data. The embedded data points are denoted  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)} \in \mathbb{R}^d$ , where  $d \leq D$  is the *embedding dimension*. Here we formally define the variables that underlie the algorithms in the `distortions` package. For more background on the statistical and mathematical basis of embedding algorithms, the reader is referred to the [35] review.

602

### 603 4.2 Neighborhood graph

604

Embedding algorithms such as UMAP [33], Isomap [56], t-SNE [30], DiffusionMaps [8], or LTSA [58] each output different embeddings  $\mathbf{y}^{(1:n)}$ , but they all start from the same data representation, which is the *neighborhood graph*. Specifically, the first step in embedding data as well as in analyzing an embedding is to find neighbors of each data point  $\mathbf{x}^{(i)}$ . This leads to the construction of the neighborhood graph as follows. Every data point  $\mathbf{x}^{(i)}$  represents a node in this graph, and two nodes are connected by an edge if their corresponding data points are neighbors. We use  $\mathcal{N}_i$  to denote the neighbors of  $\mathbf{x}^{(i)}$  and  $k_i = |\mathcal{N}_i|$  be the number of neighbors of  $\mathbf{x}^{(i)}$ . This graph, with suitable weights that summarize the local geometric and topological information in the data, is the typical input to a nonlinear dimension reduction algorithm.

605

There are two usual ways to define neighbors. In the *k-nearest neighbor (k-NN) graph*,  $\mathbf{x}^{(i')}$  is the neighbor of  $\mathbf{x}^{(i)}$  iff  $\mathbf{x}^{(i')}$  is among the closest  $k$  points to  $\mathbf{x}^{(i)}$ . In a *radius-neighbor graph*,  $\mathbf{x}^{(i')}$  is a neighbor of  $\mathbf{x}^{(i)}$  iff  $\|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\| \leq r$ , with  $r$  a parameter that defines the neighborhood scale. The k-NN graph has many computational advantages since it is connected for any  $k > 1$  and each node has between  $k$  and  $2k - 1$  neighbors (including itself). Many software packages are available to construct (approximate) k-NN graphs fast for large data [16, 6, 37].

606

The distances between neighbors are stored in the distance matrix  $\mathbf{A}$ , with  $\mathbf{A}_{ii'}$  being the distance  $\|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\|$  if  $\mathbf{x}^{(i')} \in \mathcal{N}_i$ , and infinity if  $\mathbf{x}^{(i')}$  is not a neighbor of  $\mathbf{x}^{(i)}$ . For biological data analysis, specialized distance functions can replace the generic Euclidean distance [27, 21, 64, 19]. From  $\mathbf{A}$ , another data representation is calculated, in the form of an  $n \times n$  matrix of weights that are decreasing with distances. This is called the *similarity matrix*. The weights are given by a *kernel function* [52], for example, the Gaussian kernel, defined as

$$631 \quad \mathbf{K}_{ii'} := \begin{cases} \exp\left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\|^2}{h^2}\right), & \mathbf{x}^{(i')} \in \mathcal{N}_i, \\ 632 \quad 0, & \text{otherwise.} \end{cases} \quad (1)$$

633

In the above,  $h$ , the kernel width, is another hyperparameter that must be tuned. Note that, even if  $\mathcal{N}_i$  would trivially contain all the data points, the similarity  $\mathbf{K}_{ii'}$  would be vanishingly small for faraway data points. Therefore, (1) effectively defines a radius-neighbor graph with  $r \propto h$ . Hence, a rule of thumb is to select  $r$  to be a small multiple of  $h$  (e.g.,  $r \approx 3h-10h$ ) [35]. <sup>1</sup>

634

The neighborhood graph augmented with the distance matrix  $\mathbf{A}$  or with similarity matrix  $\mathbf{K}$  has many uses:

635

1. As stated above, it serves as a starting point for embedding algorithms.
2. In this paper,  $\mathbf{K}$  is used to calculate the local distortion.
3. In this paper,  $\mathbf{A}$  is used to detect the fragmented neighborhoods.

636

---

<sup>1</sup>Sometimes, the simple similarity

$$637 \quad \mathbf{K}_{ii'} := \begin{cases} 1, & \mathbf{x}^{(i')} \in \mathcal{N}_i, \\ 638 \quad 0, & \text{otherwise} \end{cases} \quad (2)$$

639

is used. This similarity matrix  $\mathbf{K}$  is the unweighted adjacency matrix of the neighborhood graph, and completely ignores the distances.

4. Neighborhood graphs are also used in estimating the intrinsic dimension, in Topological Data Analysis,  
namely in finding the loops and hollows in the data, as well as in other Geometric Data Analysis tasks.

While most embedding algorithms can take as input both types of neighborhood graphs (or resulting distance or similarity matrices), the embeddings obtained will be influenced by the type of graph and by the hyperparameter value used with it. For other uses, one type of graph or another may be optimal. In particular, for the purpose of estimating distortion, it is necessary to use the radius-neighbor graph, as this guarantees the distortion estimated is unbiased.

### 4.3 Distortion estimation with the Dual Pushforward Riemannian Metric

The distortion estimation function `local_distortions()` implements the algorithm introduced by [43]. Given an embedding  $\mathbf{y}^{(1:n)}$  of data  $\mathbf{x}^{(1:n)}$  with similarity matrix  $\mathbf{K}$  computed from radius-neighbor graph, `local_distortions()` outputs for each embedding point  $\mathbf{y}^{(i)}$  a  $d \times d$  matrix  $\mathbf{V}^{(i)}$  whose column  $\mathbf{v}_1^{(i)}, \dots, \mathbf{v}_d^{(i)}$  represent the *principal directions* of distortion at data point  $i$ . The stretch in direction  $\mathbf{v}_j^{(i)}$  is given by  $\lambda_j^{(i)}$ . When  $\lambda_j^{(i)} = 1$  there is no stretch, for  $\lambda_j^{(i)} > 1$  the embedding stretches the data in direction  $\mathbf{v}_j^{(i)}$ , and for  $0 < \lambda_j^{(i)} < 1$  the embedding shrinks the data along this direction. Thus, the principal directions are orthogonal directions in the embedding where the algorithm induces pure stretch. Intuitively, the values  $\lambda_j^{(i)}$  represents the local unit of length in direction  $\mathbf{v}_j^{(i)}$ .

The principal directions and stretch values result from the eigendecomposition of the symmetric, positive definite matrix  $\mathbf{H}^{(i)} = \mathbf{V}^{(i)} \text{diag}\{\lambda_1^{(i)}, \dots, \lambda_m^{(i)}\} \mathbf{V}^{(i)\top}$ . For an embedding with no distortion, namely an *isometric* embedding,  $\mathbf{H}^{(i)} = \mathbf{I}_d$  the unit matrix.

The local correction at  $\mathbf{y}^{(i)}$  is the inverse  $\mathbf{G}^{(i)}$  of  $\mathbf{H}^{(i)}$ ; in technical terms  $\mathbf{G}^{(i)}$  is known as the *embedding (push-forward) Riemannian metric*. Obviously, the eigendecomposition of  $\mathbf{G}^{(i)}$  is given by  $\mathbf{V}^{(i)}$  and  $1/\lambda_1^{(i)}, \dots, 1/\lambda_m^{(i)}$ . Thus, to correct the distortion in direction  $\mathbf{y}^{(i')} - \mathbf{y}^{(i)}$ , one calculates  $\mathbf{G}^{(i)}(\mathbf{y}^{(i')} - \mathbf{y}^{(i)})$ . The orientation and length of this vector with origin in  $\mathbf{y}^{(i)}$  are the corrected direction and distance to nearby point  $\mathbf{y}'$ .

Hence, for any data embedding, it is sufficient to estimate, at all points  $\mathbf{y}^{(1:n)}$ , the matrices  $\mathbf{G}^{(1:n)}$ , which represent the auxiliary information enabling correct distance computations, as if working with the original data, even though the embedding may not have preserved them. The same  $\mathbf{G}^{(1:n)}$  can be used to preserve not only geodesic distances but also other geometric quantities such as angles between curves in  $\mathcal{M}$  or volumes of subsets of  $\mathcal{M}$ . Further uses of the distortion and correction matrices are described in [35, 43, 20], and here we present a corrected visualization based on  $\mathbf{G}^{(i)}$ .

### 4.4 Computational complexity of RMetric

The complexity of the RMetric computation is dominated by the construction of the neighborhood graph. Since this graph is already computed for the purpose of embedding the data, we will only consider the overhead. Obtaining the similarity  $\mathbf{K}$  involves a fixed set of operations per graph edge (i.e. calculating the kernel value), hence order  $m$  operations total, where  $m$  is the number of edges in the neighborhood graph. Further computations also are proportional to  $m$ . Computing the RMetric at point  $i$  requires  $\sim k_i d^2$  operations, where  $k_i$  is, as above, the number of neighbors of  $i$ . Hence, obtaining the RMetric at all points requires  $\sim md^2$  operations.<sup>2</sup> Further eigendecompositions and inversion of  $\mathbf{H}^{(i)}$  are order  $d^3$  per data point, hence  $nd^3$  total.

Since the optimal neighborhood graph is a sparse graph (since it should only capture distance to nearby points and ignore the distances to far-away points),  $m$  is much smaller than the maximum value  $n(n-1)/2$ . In practice, on large data sets, we have always found that computing the RMetric is much faster than computing the embedding itself. The same is true for the isometrization algorithm, in which the overhead after RMetric computation is to apply a simple transformation to every embedded point.

<sup>2</sup>Since  $\sum_i k_i = 2m$ .

#### 703 4.5 Selecting the hyperparameter $h$

704 We recommend [35] for a tutorial on the choice of parameters  $k$  and/or  $h$  (with  $r$  being a small multiple  
705 of  $h$ ). An automatic method for choosing these parameters, reminiscent of cross-validation, was introduced  
706 by [43] and can be found in the `megaman` package <https://mmp2.github.io/megaman/>.  
707

708 As a general rule of thumb, if a neighborhood graph results in a good embedding, then the neighborhood  
709 scale is the appropriate one for the RMetric as well. Hence, if the embedding is obtained via a radius-neighbor  
710 graph, then the same graph, or same  $\mathbf{K}$  matrix should be used for `local_distortions()`. If a  $k$ -NN graph  
711 was used, then we recommend selecting  $h$  so that the row sums of  $\mathbf{K}$  average  $k$ , the neighborhood parameter  
712 of the  $k$ -NN graph.  
713

#### 713 4.6 Identifying fragmented neighborhoods

714 To compare distances across the original and embedding space, let:

$$716 \quad \mathcal{D} := \bigcup_{i=1}^N \left\{ \left( \|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\|, \|\mathbf{y}^{(i)} - \mathbf{y}^{(i')}\| \right) \in \mathbb{R}^2 \text{ for } i' \in \mathcal{N}_i \right\}$$

717 The `distortions` package supports two strategies for flagging neighbors with poorly preserved distances,  
718 which form the basis for defining fragmented neighborhoods.  
719

##### 721 Bin-based strategy

722 This approach partitions the original space distances into  $L$  evenly-sized bins and detects outliers in the  
723 embedding distances within each bin. Let  $\pi_O(\mathcal{D})$  and  $\pi_E(\mathcal{D})$  extract the original and embedding distances  
724 from  $\mathcal{D}$ , respectively. With  $d_{\min} = \inf \pi_O(\mathcal{D})$  and  $d_{\max} = \sup \pi_O(\mathcal{D})$ , set the binwidth  $w = \frac{1}{L} (d_{\max} - d_{\min})$   
725 and partition the original data distances into intervals  $I_l = [d_{\min} + w(l-1), d_{\min} + wl]$ . The embedding  
726 distances within bin  $l$  are,  
727

$$728 \quad \mathcal{D}_l := \{d \in \mathcal{D} : \pi_O(d) \in I_l\}$$

729 where we have abused notation and applied the projection  $\pi_O$  to an individual distance tuple  $d$ . For each  
730 bin, we compute the interquartile range of associated embedding distances,  
731

$$733 \quad \text{IQR}_l = Q_{0.75}(\pi_E(\mathcal{D}_l)) - Q_{0.25}(\pi_E(\mathcal{D}_l))$$

735 where  $Q_\alpha$  extracts the  $\alpha$ -quantile. A distance tuple  $d \in \mathcal{D}$  is considered outlying if,

$$737 \quad \pi_E(d) \notin [Q_{0.5}(\pi_E(\mathcal{D}_l)) - \sigma \text{IQR}_l, Q_{0.5}(\pi_E(\mathcal{D}_l)) + \sigma \text{IQR}_l]$$

739 where  $\sigma$  controls the outlier threshold. Note that neighborhood distances can be considered outlying for  
740 two qualitatively different reasons. The embedding distance may be either too large, where truly neighboring  
741 points may be artificially spread apart. This is labeled  $\mathcal{O}_l^+$  in Fig 10. Alternatively, they may be too small,  
742 where distant points are inappropriately collapsed on top of one another ( $\mathcal{O}_l^-$  in Fig 10). All bin- $l$  outliers  
743 are collected into the set  $\mathcal{O}_l = \mathcal{O}_l^- \cup \mathcal{O}_l^+$ .  
744

745 We define fragmented neighborhoods using the outlier sets  $\mathcal{O}_l$ . We consider  $\mathbf{y}^{(i)}$  to be the center of a  
746 fragmented neighborhood if,

$$747 \quad \frac{\left| \left\{ i' : \left( \|\mathbf{x}^{(i)} - \mathbf{x}^{(i')}\|, \|\mathbf{y}^{(i)} - \mathbf{y}^{(i')}\| \right) \in \cup \mathcal{O}_l \right\} \cap \mathcal{N}_i \right|}{|\mathcal{N}_i|} \geq \kappa, \quad (3)$$

750 that is, if at least a fraction  $\kappa$  of the distances to its neighbors belong to at least one outlier set  $\mathcal{O}_l$ . This  
751 procedure is illustrated graphically in Fig 10.  
752

##### 753 Window-based strategy

754 The window-based strategy parallels the bin-based approach but uses running windows centered at each  
755 point. For each  $d_0 \in \mathcal{D}$ , we define a window  $\mathcal{D}_{\text{win}}$  of the  $\Delta$  nearest points with respect to  $\pi_O(d_0)$ . Within

each window, we compute the interquartile range (IQR) of the embedding distances and flag  $d \in \mathcal{D}$  as an outlier if its embedding distance is more than  $\sigma$  IQRs from the median embedding distance in the window,

$$\pi_E(d) \notin [Q_{0.5}(\pi_E(\mathcal{D}_{\text{win}}(d_0))) - \sigma \text{IQR}(d_0), Q_{0.5}(\pi_E(\mathcal{D}_{\text{win}}(d_0))) + \sigma \text{IQR}(d_0)]$$

where  $\text{IQR}(d)$  is the interquartile range of  $\pi_E(\mathcal{D}_{\text{win}}(d))$ . As with the bin-based strategy, a neighborhood is fragmented if at least a fraction  $\kappa$  of its neighbor pairs are flagged as outliers. This approach leads to smoother IQR boundaries compared to the bin-based approach, but is more computationally involved.

#### 4.7 Focus-plus-context visual interaction

Adding distortion information to standard nonlinear embedding visualizations is challenging because the additional context can overwhelm an already complex visualization, making them even more difficult to understand. The `distortions` package addresses this challenge through the focus-plus-context principle [15, 24, 46]. This approach displays distortion information locally (“focus”) while maintaining the broader visual overview (“context”). The region within which to display additional information is set by the viewer’s interactions. We implement three forms of focus-plus-context interactivity, adapted to visualize fragmented neighborhoods, distance preservation, and local isometries, respectively.

**4.7.1 Mouseover interactions to reveal fragmented neighborhoods** This visualization supplements the original embedding overview by highlighting fragmented neighborhoods when their centers are hovered over. The centers may be defined using either the bin-based or window-based strategies described above. Before interaction, the fragmented neighborhood centers are highlighted with a distinctive stroke and color, guiding attention to regions of the embedding that are enriched with fragmentation. When the viewer’s mouse is moved to a location  $m \in \mathbb{R}^2$ , all fragmented neighborhoods with centers within a distance  $\delta$  of  $m$  are highlighted. Specifically, an edge is drawn between  $\mathbf{y} \in \mathbb{R}^2$  and  $\mathbf{y}' \in \mathbb{R}^2$  if:

1.  $\|\mathbf{y} - m\| \leq \delta$ .
2.  $\mathbf{y}$  satisfies the fragmented neigorhood criterion (Equation 3).
3.  $\mathbf{y}'$  is one of the top  $k$  neighbors of  $\mathbf{y}$  in the original data space.

The neighbors  $\mathbf{y}'$  are highlighted when their corresponding edge links are visible. The hyperparameters  $\delta$  and  $k$  must be specified by the viewer. We default to the  $k$  used in the original embedding. Since the neighborhoods are fragmented, the associated edge links typically span large regions of the embedding space, making interactive updates necessary to prevent occlusion from overlapping edges.

**4.7.2 Brush interactions to visualize distance preservation** The focus-plus-context principle supports visualization of individual edges with poorly preserved distances, rather than entire neighborhoods. A brushable widget is placed alongside the main embedding visualization and displays boxplots that compare binned distances in the original data space ( $x$ -axis) with the distances in the embedding space ( $y$ -axis). This boxplot overview builds on the static approach of [22]. Boxplot whiskers are capped at  $\sigma$  times the IQR, with outliers beyond this range drawn as distinct points. The number of bins and  $\sigma$  are user-specified hyperparameters. As the brush is moved, the embedding visualization updates to highlight edges between neighbors with brushed and outlying embedding distances. The coordinated display allows viewers to focus on specific distorted neighbor pairs within the context given by the overview boxplots.

**4.7.3 Mouseover interactions to update local isometries** The `distortions` package supports interactions that provide an intuitive understanding of local metric differences induced by the embedding. In this view, the mouse’s position is used to isometrize neighborhoods centered around it, providing an interactive, local version of the isometrization algorithm from [43]. Rather than modifying the entire embedding to induce an isometry around a selected point, this view updates the region around the mouse position. To isometrize the embedding with respect to sample  $i'$ , [43] suggest the transformation,

$$\mathbf{y}^{(i)} \rightarrow \left( \mathbf{H}^{(i')} \right)^{-1} \mathbf{y}^{(i)}$$

811 For focus-plus-context interaction, we isometrize only samples near the mouse position  $m$  and smoothly  
812 interpolate the transformation as the mouse moves between samples. We implement,  
813

$$\mathbf{y}^{(i)} \rightarrow k_{h_1}(\mathbf{y}^{(i)}, m) \tilde{\mathbf{y}}^{(i)} + (1 - k_{h_1}(\mathbf{y}^{(i)}, m)) \mathbf{y}^{(i)} \quad (4)$$

816 where,  
817

$$\tilde{\mathbf{y}}^{(i)} := (\mathbf{H}^*)^{-1}(\mathbf{y}^{(i)} - m) + m \quad (5)$$

$$\mathbf{H}^* = \sum_{j=1}^N \left[ \frac{k_{h_2}(\mathbf{y}^{(j)}, m)}{\sum_{j'=1}^N k_{h_2}(\mathbf{y}^{(j')}, m)} \right] \mathbf{H}^{(j)}. \quad (6)$$

824 and  $k_{h_g}$  denotes the Gaussian kernel with bandwidth  $h_g$  and  $\mathbf{H}^*$  represents a local average of  $\mathbf{H}^{(i)}$ . The  
825 parameter  $h_1$  controls the size of the region affected by isometrization, and  $h_2$  controls the the region defining  
826  $\mathbf{H}^*$ . This interactive coordinate system update is related to fisheye distortion [47], where local geometries  
827 are deliberately distorted to focus on specific samples.  
828

#### 829 4.8 Package software architecture

830 The `distortions` software architecture must support low-level graphical marks, like ellipses, and  
831 interactions, like updating fragmented neighborhood links on mouseover, that are unavailable in existing  
832 visualization software. These customizations cannot come at the cost of support for higher-level data struc-  
833 tures from modern computational biology software. To this end, we have defined a standalone javascript  
834 package (`distortions-js`, <https://www.npmjs.com/package/distortions>) for visual components and interac-  
835 tions, and a separate python package (`distortions`, <https://pypi.org/project/distortions/>) for higher-level  
836 algorithms and distortion computation. The javascript implementation is built around a `DistortionPlot`  
837 class, which exports a `mapping` method to encode dataset fields in the visual channels from each `geom*` ele-  
838 ment, as well as methods for each interaction type. All graphical marks are rendered as SVG elements on  
839 a parent canvas. This is necessary, as standard javascript plotting libraries like `vega` [49] and `observable`  
840 `plot` [42] do not support ellipse visualization. Brush events are implemented using the d3-brush library [2],  
841 and legends are drawn using d3-legend [28].  
842

843 The python package connects to `distortions-js` through the `anywidget` [31] package, allowing interac-  
844 tive javascript execution within Jupyter and Quarto notebooks. This approach converts python dictionary  
845 objects storing data and plot specifications into javascript data structures for visualization in the browser or  
846 notebook cell. The embeddings can be passed in through an `AnnData` experimental object [62]. For intrinsic  
847 geometry estimation, we use the `megaman` package [34], which is designed for scalable nonlinear dimensional-  
848 ity reduction and supports estimation of the local metrics  $\mathbf{H}^{(i)}$  for each sample  $i$ . Open source code can be  
849 found at <https://github.com/krisrs1128/distortions> and <https://github.com/krisrs1128/distortions-js>, docu-  
850 documentation is given at <https://krisrs1128.github.io/distortions/site/>. We note that the packages can be used  
851 independently.  
852

853 **Acknowledgements.** MM gratefully acknowledges the DataShape Group at INRIA Saclay and the long  
854 program on “Data driven materials informatics” at the Institute for Mathematical and Statistical Innovation  
855 (IMSI) for hosting her, and last but not least, the University of Washington Department of Statistics, where  
856 most of this research was conducted. KS and MM acknowledge the University of Wisconsin-Madison SILO  
857 seminar series, where this collaboration was initiated.  
858

## 858 A Supplementary Figures

## 859 References

- 860 [1] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps and spectral techniques for embedding and  
861 clustering. *Advances in neural information processing systems*, 14, 2001.  
862 [2] Michael Bostock, Vadim Ogievetsky, and Jeffrey Heer. D<sup>3</sup> data-driven documents. *IEEE transactions*  
863 *on visualization and computer graphics*, 17(12):2301–2309, 2011.  
864

- [3] T Tony Cai and Rong Ma. Theoretical foundations of t-sne for visualizing high-dimensional clustered data. *Journal of Machine Learning Research*, 23(301):1–54, 2022. 865  
866
- [4] Junyue Cao, Malte Spielmann, Xiaojie Qiu, Xingfan Huang, Daniel M. Ibrahim, Andrew J. Hill, Fan Zhang, Stefan Mundlos, Lena Christiansen, Frank J. Steemers, Cole Trapnell, and Jay Shendure. The single-cell transcriptional landscape of mammalian organogenesis. *Nature*, 566(7745):496–502, February 2019. 867  
868  
869  
870
- [5] Tara Chari and Lior Pachter. The specious art of single-cell genomics. *PLOS Computational Biology*, 19(8):e1011288, August 2023. 871  
872
- [6] Jie Chen, Haw ren Fang, and Yousef Saad. Fast approximate knn graph construction for high dimensional data via recursive lanczos bisection. *Journal of Machine Learning Research*, 10(69):1989–2012, 2009. 873  
874
- [7] Andy Coenen and Adam Pearce. Understanding UMAP — pair-code.github.io. <https://pair-code.github.io/understanding-umap/>, 2019. [Accessed 06-06-2025]. 875  
876
- [8] Ronald R Coifman and Stéphane Lafon. Diffusion maps. *Applied and computational harmonic analysis*, 21(1):5–30, 2006. 877  
878
- [9] The Tabula Muris Consortium. Single-cell transcriptomics of 20 mouse organs creates a tabula muris. *Nature*, 562(7727):367–372, October 2018. 879  
880
- [10] Persi Diaconis, Sharad Goel, and Susan Holmes. Horseshoes in multidimensional scaling and local kernel methods. 2008. 881  
882
- [11] Julia Fukuyama, Kris Sankaran, and Laura Symul. Multiscale analysis of count data through topic alignment. *Biostatistics*, 24(4):1045–1065, June 2022. 883  
884
- [12] F Richard Guo and Rajen D Shah. Rank-transformed subsampling: inference for multiple data splitting and exchangeable p-values. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 885  
886  
887
- [13] Luke T.G. Harland, Tim Lohoff, Noushin Koulena, Nico Pierson, Constantin Pape, Farhan Ameen, Jonathan Griffiths, Bart Theeuwes, Nicola K. Wilson, Anna Kreshuk, Wolf Reik, Jennifer Nichols, Long Cai, John C. Marioni, Berthold Göttgens, and Shila Ghazanfar. A spatiotemporal atlas of mouse gastrulation and early organogenesis to explore axial patterning and project in vitro models onto in vivo space. *Cell Reports*, 44(8):116047, August 2025. 888  
889  
890  
891  
892
- [14] Jeffrey Heer and Maneesh Agrawala. Software design patterns for information visualization. *IEEE transactions on visualization and computer graphics*, 12(5):853–860, 2006. 893  
894
- [15] Jeffrey Heer and Stuart K Card. Doitrees revisited: scalable, space-constrained visualization of hierarchical data. In *Proceedings of the working conference on Advanced visual interfaces*, pages 421–424, 2004. 895  
896  
897
- [16] Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, pages 604–613, 1998. 898  
899  
900
- [17] Rafael Irizarry. Simply Statistics: Biologists, stop putting UMAP plots in your papers — simplystatistics.org. <https://simplystatistics.org/posts/2024-12-23-biologists-stop-including-umap-plots-in-your-papers/>, 2024. [Accessed 17-07-2025]. 901  
902  
903
- [18] So Won Jeong and Claire Donnat. Lobstur: A local bootstrap framework for tuning unsupervised representations in graph neural networks, 2025. 904  
905
- [19] Yuge Ji, Tessa D Green, Stefan Peidli, Mojtaba Bahrami, Meiqi Liu, Luke Zappia, Karin Hrovatin, Chris Sander, and Fabian J Theis. Optimal distance metrics for single-cell rna-seq populations. *BioRxiv*, pages 2023–12, 2023. 906  
907  
908
- [20] Dominique Joncas, Marina Meila, and James McQueen. Improved graph laplacian via geometric Self-Consistency. In I Guyon, U V Luxburg, S Bengio, H Wallach, R Fergus, S Vishwanathan, and R Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4457–4466. Curran Associates, Inc., 2017. 909  
910  
911  
912
- [21] Taiyun Kim, Irene Rui Chen, Yingxin Lin, Andy Yi-Yang Wang, Jean Yee Hwa Yang, and Pengyi Yang. Impact of similarity metrics on single-cell rna-seq data clustering. *Briefings in bioinformatics*, 20(6):2316–2326, 2019. 913  
914  
915
- [22] Dmitry Kobak and George C. Linderman. Initialization is critical for preserving global data structure in both t-sne and umap. *Nature Biotechnology*, 39(2):156–157, February 2021. 916  
917  
918

- 919 [23] Piotr H Laskowski. The traditional and modern look at Tissot's indicatrix. *The American Cartographer*,  
920 16(2):123–133, 1989.
- 921 [24] Keith Lau, Ronald A Rensink, and Tamara Munzner. Perceptual invariance of nonlinear focus+ con-  
922 text transformations. In *Proceedings of the 1st Symposium on Applied perception in graphics and*  
923 *visualization*, pages 65–72, 2004.
- 924 [25] Jan Lause, Philipp Berens, and Dmitry Kobak. The art of seeing the elephant in the room: 2d  
925 embeddings of single-cell data do make sense. *PLOS Computational Biology*, 20(10):e1012403, October  
926 2024.
- 927 [26] Zhexuan Liu, Rong Ma, and Yiqiao Zhong. Assessing and improving reliability of neighbor embedding  
928 methods: a map-continuity perspective. *Nature Communications*, 16(1), May 2025.
- 929 [27] Catherine Lozupone and Rob Knight. Unifrac: a new phylogenetic method for comparing microbial  
930 communities. *Applied and environmental microbiology*, 71(12):8228–8235, 2005.
- 931 [28] Susie Lu. d3-legend.susielu.com. <https://d3-legend.susielu.com/>, 2016. [Accessed 11-08-2025].
- 932 [29] Rong Ma, Xi Li, Jingyuan Hu, and Bin Yu. Uncovering smooth structures in single-cell data with  
933 pcs-guided neighbor embeddings. *biorXiv*, July 2025.
- 934 [30] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning*  
935 *research*, 9(Nov):2579–2605, 2008.
- 936 [31] Trevor Manz, Nezar Abdennur, and Nils Gehlenborg. anywidget: reusable widgets for interactive analysis  
937 and visualization in computational notebooks. *Journal of Open Source Software*, 9(102):6939, 2024.
- 938 [32] Vivien Marx. Seeing data as t-sne and umap do. *Nature Methods*, 21(6):930–933, May 2024.
- 939 [33] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and  
940 projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- 941 [34] James McQueen, Marina Meilă, Jacob VanderPlas, and Zhongyue Zhang. Megaman: Scalable manifold  
942 learning in python. *Journal of Machine Learning Research*, 17(148):1–5, 2016.
- 943 [35] Marina Meilă and Hanyu Zhang. Manifold learning: What, how, and why. *Annual Review of Statistics*  
944 *and Its Application*, 11(1):393–417, 2024.
- 945 [36] Marina Meilă and Hanyu Zhang. Manifold learning: what, how, and why. *Annual Reviews in Statistics*  
946 *and its Applications*, (accepted, 2024).
- 947 [37] Marius Muja and David G Lowe. Fast approximate nearest neighbors with automatic algorithm  
948 configuration. *VISAPP (1)*, 2(331-340):2, 2009.
- 949 [38] Ashwin Narayan, Bonnie Berger, and Hyunghoon Cho. Assessing single-cell transcriptomic variability  
950 through density-preserving data visualization. *Nature biotechnology*, 39(6):765–774, 2021.
- 951 [39] Lan Huong Nguyen and Susan Holmes. Ten quick tips for effective dimensionality reduction. *PLoS*  
952 *computational biology*, 15(6):e1006907, 2019.
- 953 [40] Max Noichl. Max Noichl — Flattening Mammoths — maxnoichl.eu. <https://www.maxnoichl.eu/projects/mammoth/>, 2019. [Accessed 08-07-2025].
- 954 [41] Jonathan S Packer, Qin Zhu, Chau Huynh, Priya Sivaramakrishnan, Elicia Preston, Hannah Dueck,  
955 Derek Stefanik, Kai Tan, Cole Trapnell, Junhyong Kim, et al. A lineage-resolved molecular atlas of c.  
956 elegans embryogenesis at single-cell resolution. *Science*, 365(6459):eaax1971, 2019.
- 957 [42] Jeffrey M. Perkel. Reactive, reproducible, collaborative: computational notebooks evolve. *Nature*,  
958 593(7857):156–157, May 2021.
- 959 [43] Dominique Perrault-Joncas and Marina Meila. Non-linear dimensionality reduction: Riemannian metric  
960 estimation and the problem of geometric discovery, May 2013.
- 961 [44] Dominique Perrault-Joncas and Marina Meilă. Riemannian learning of manifolds. <https://sites.stat.washington.edu/mmp/Papers/RMetric.pdf>, 2013. Unpublished manuscript; available at the authors'  
962 website (accessed 2025-08).
- 963 [45] Regev, Aviv, Teichmann, Sarah A, Lander, Eric S, Amit, Ido, Benoist, Christophe, Birney, Ewan,  
964 Bodenmiller, Bernd, and et al. The human cell atlas. 2017.
- 965 [46] Kris Sankaran and Susan Holmes. Interactive visualization of hierarchically structured data. *J. Comput.*  
966 *Graph. Stat.*, 27(3):553–563, 2018.
- 967 [47] Manojit Sarkar and Marc H. Brown. Graphical fisheye views of graphs. In *Proceedings of the SIGCHI*  
968 *conference on Human factors in computing systems - CHI '92*, CHI '92, pages 83–91. ACM Press, 1992.
- 969 [48] Rahul Satija, Jeffrey A Farrell, David Gennert, Alexander F Schier, and Aviv Regev. Spatial  
970 971 972

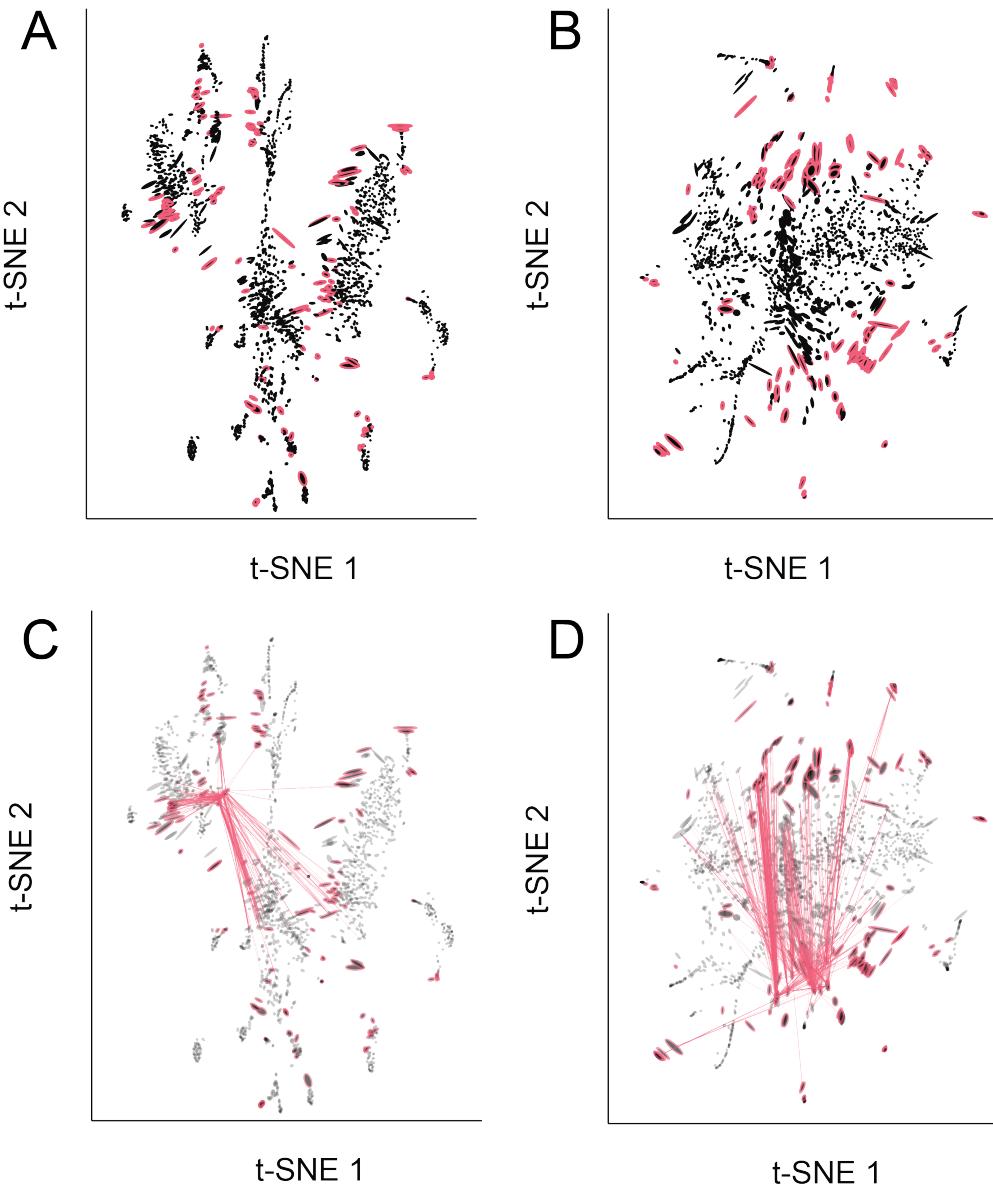
- reconstruction of single-cell gene expression data. *Nature biotechnology*, 33(5):495–502, 2015. 973
- [49] Arvind Satyanarayan, Ryan Russell, Jane Hoffswell, and Jeffrey Heer. Reactive vega: A streaming 974 dataflow architecture for declarative interactive visualization. *IEEE Transactions on Visualization and 975 Computer Graphics*, 22(1):659–668, 2016. 976
- [50] Scanpy Development Team. scanpy.datasets.pbmc3k\_processed — scanpy.readthedocs.io. [https://scanpy.readthedocs.io/en/stable/generated/scanpy.datasets.pbmc3k\\_processed.html](https://scanpy.readthedocs.io/en/stable/generated/scanpy.datasets.pbmc3k_processed.html), 2025. [Accessed 977 08-07-2025]. 978
- [51] Jonas Schluter, Ana Djukovic, Bradford P. Taylor, Jinyuan Yan, Caichen Duan, Grant A. Hussey, Chen 979 Liao, Sneh Sharma, Emily Fontana, Luigi A. Moretti, Roberta J. Wright, Anqi Dai, Jonathan U. Peled, Ying Taur, Miguel-Angel Perales, Benjamin A. Siranosian, Ami S. Bhatt, Marcel R.M. van den Brink, Eric G. Pamer, and Joao B. Xavier. The taxumap atlas: Efficient display of large clinical 980 microbiome data reveals ecological competition in protection against bacteremia. *Cell Host and Microbe*, 981 31(7):1126–1139.e6, July 2023. 982
- [52] Bernhard Schölkopf, Koji Tsuda, and Jean-Philippe Vert. *Kernel methods in computational biology*. 983 MIT press, 2004. 984
- [53] Stefan Siebert, Jeffrey A. Farrell, Jack F. Cazet, Yashodara Abeykoon, Abby S. Primack, Christine E. Schnitzler, and Celina E. Juliano. Stem cell differentiation trajectories in hydra resolved at single-cell 985 resolution. *Science*, 365(6451), July 2019. 986
- [54] Dongyuan Song, Kexin Li, Xinzhou Ge, and Jingyi Jessica Li. Clusterde: a post-clustering differential 987 expression (de) method robust to false-positive inflation caused by double dipping. *Research Square*, 988 pages rs-3, 2023. 989
- [55] Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalexi, William M 990 Mauck, Yuhan Hao, Marlon Stoeckius, Peter Smibert, and Rahul Satija. Comprehensive integration of 991 single-cell data. *cell*, 177(7):1888–1902, 2019. 992
- [56] Joshua B Tenenbaum, Vin de Silva, and John C Langford. A global geometric framework for nonlinear 993 dimensionality reduction. *science*, 290(5500):2319–2323, 2000. 994
- [57] Jacob VanderPlas, Brian Granger, Jeffrey Heer, Dominik Moritz, Kanit Wongsuphasawat, Arvind 995 Satyanarayan, Eitan Lees, Ilia Timofeev, Ben Welsh, and Scott Sievert. Altair: Interactive statistical 996 visualizations for python. *Journal of Open Source Software*, 3(32):1057, December 2018. 997
- [58] Jianzhong Wang. Local tangent space alignment. In *Geometric Structure of High-Dimensional Data 998 and Dimensionality Reduction*, pages 221–234. Springer, 2012. 999
- [59] Shu Wang, Eduardo D. Sontag, and Douglas A. Lauffenburger. What cannot be seen correctly in 2d 1000 visualizations of single-cell ‘omics data? *Cell Systems*, 14(9):723–731, September 2023. 1001
- [60] Martin Wattenberg, Fernanda Viégas, and Ian Johnson. How to use t-sne effectively. *Distill*, 2016. 1002
- [61] Hadley Wickham and Carson Sievert. *ggplot2: elegant graphics for data analysis*, volume 10. Springer 1003 New York, 2009. 1004
- [62] F Alexander Wolf, Philipp Angerer, and Fabian J Theis. Scanpy: large-scale single-cell gene expression 1005 data analysis. *Genome biology*, 19:1–5, 2018. 1006
- [63] Lucy Xia, Christy Lee, and Jingyi Jessica Li. Statistical method scdeed for detecting dubious 2d single- 1007 cell embeddings and optimizing t-sne and umap hyperparameters. *Nature Communications*, 15(1):1753, 1008 2024. 1009
- [64] Hongxuan Zhai and Julia Fukuyama. A convenient correspondence between k-mer-based metagenomic 1010 distances and phylogenetically-informed  $\beta$ -diversity measures. *PLOS Computational Biology*, 19(1):e1010821, 2023. 1011
- 1012
- 1013
- 1014
- 1015
- 1016
- 1017
- 1018
- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025
- 1026

1027  
1028  
1029  
1030

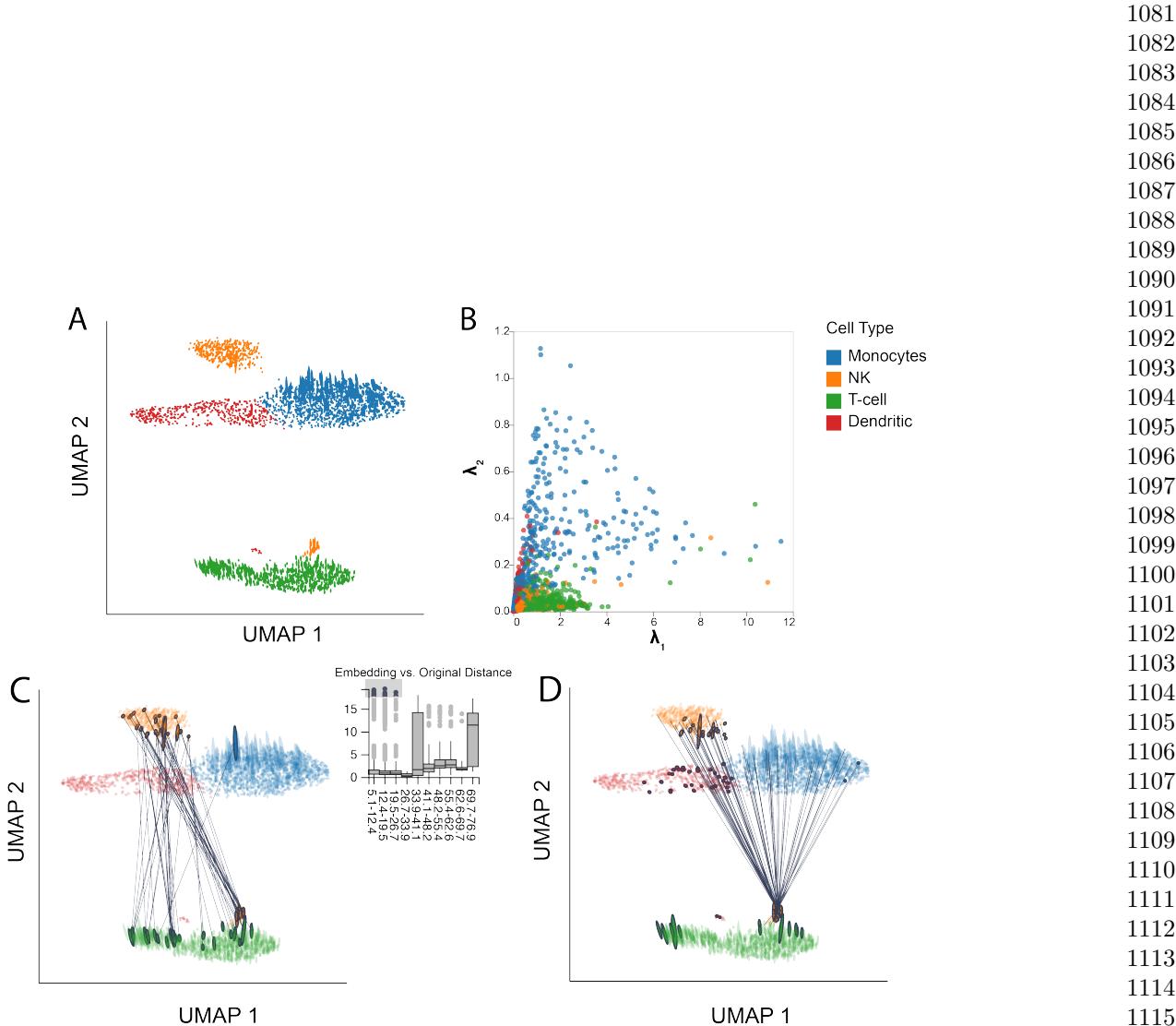
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047

1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066

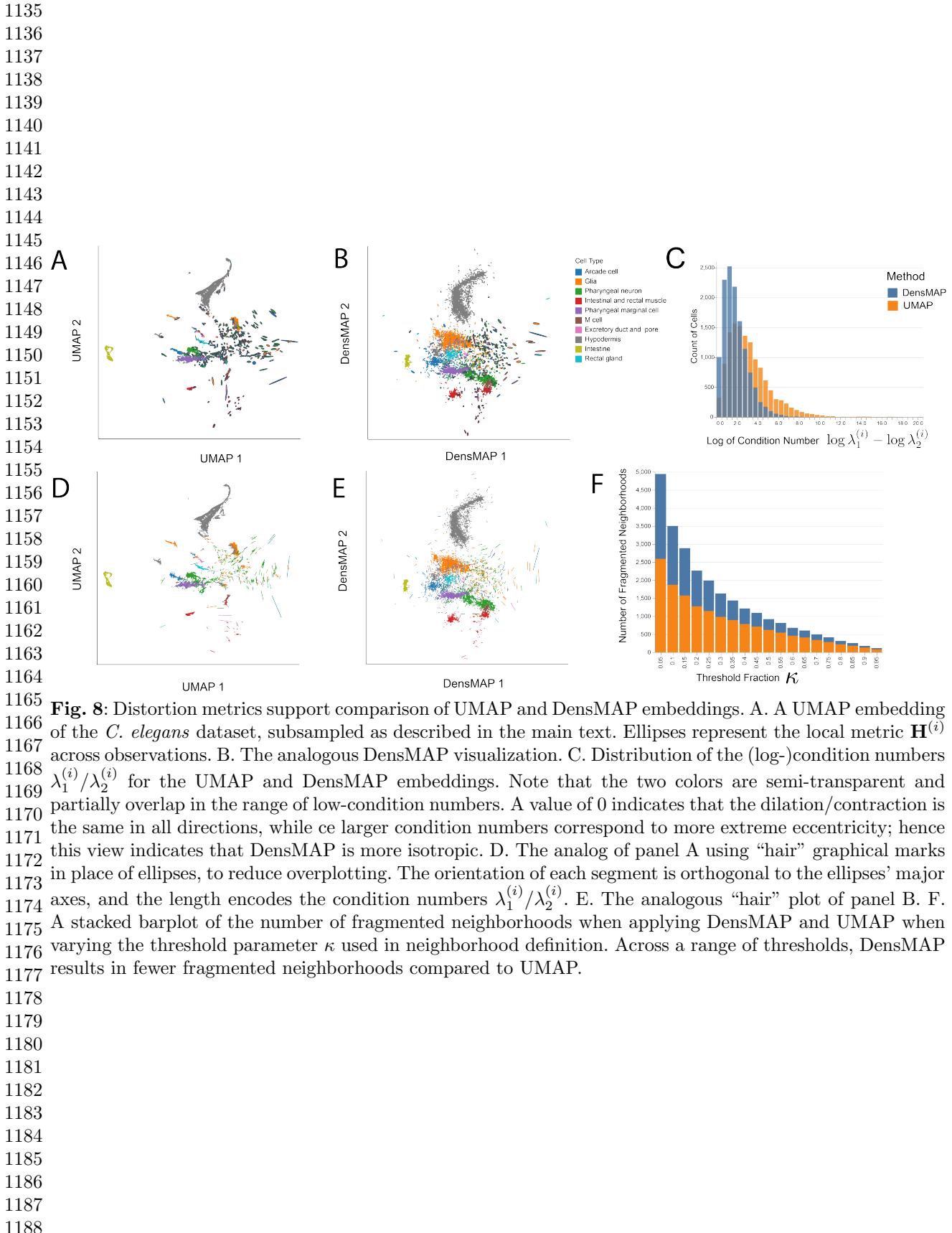
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080

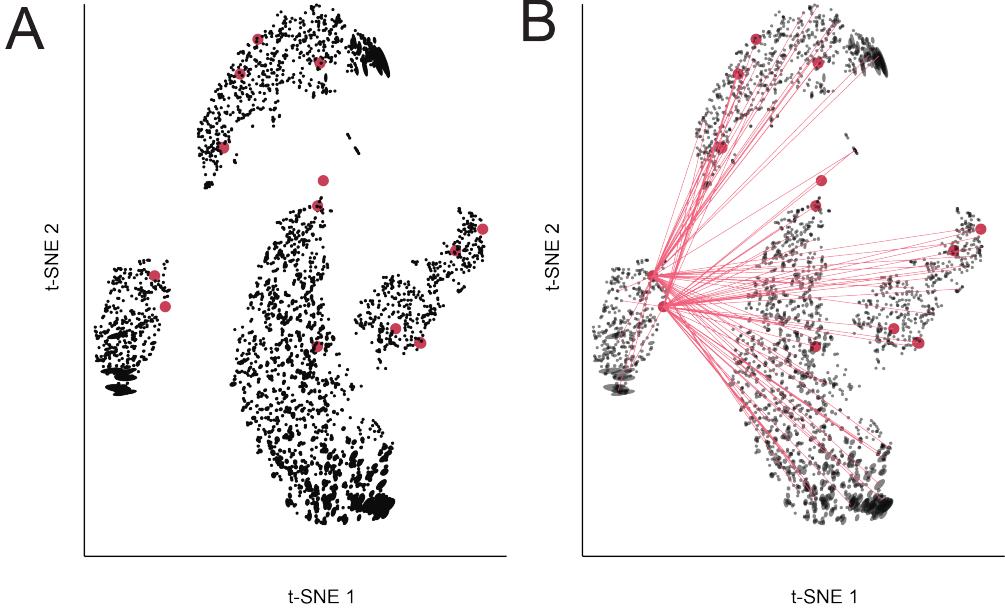


1068 **Fig. 6:** Salient characteristics of distortion vary across hyperparameter settings. A. The *t*-SNE embedding  
1069 of the hydra cell atlas dataset when perplexity hyperparameter is set to 80. This embedding exaggerates the  
1070 distinction between cell type clusters. B. The analogous view when the *t*-SNE perplexity is set to 500. At this  
1071 hyperparameter value, the main clusters are now more overlapping, but the distances along the periphery  
1072 of the embedding are less well preserved. C. Hovering over fragmented neighborhoods near the bottom-left  
1073 of the embedding in panel A shows that neighbors are often shared between clusters. D. Hovering over a  
1074 fragmented neighborhood in Panel B shows that points near the periphery can be neighbors with points  
1075 spread throughout the visualization.

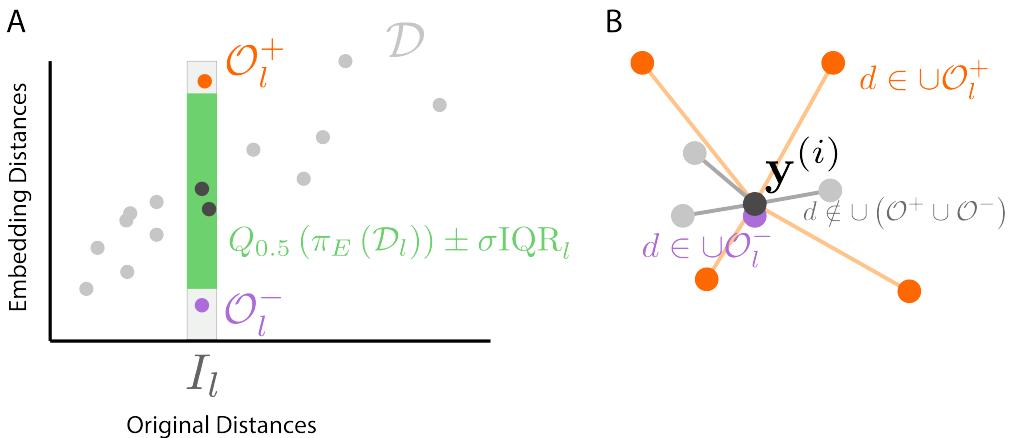


**Fig. 7:** Distortion visualizations highlight problems with randomly initialized UMAP. A. UMAP embedding of the PBMC data when applying random initialization. B. The analog of Fig 4D in the random initialization setting. The systematically larger condition numbers  $\frac{\lambda_1}{\lambda_2}$  correspond to more eccentric ellipses in panel A. C. Brushing over pairs with large embedding vs. original distances highlights T-NK cell neighbors whose relative distances are poorly preserved. These cell types are placed close to one another in the spectral initialization of Fig 4. D. Hovering over the fragmented neighborhoods in the bottom right corner of the plot highlights NK cells with more neighbors among monocytes and NK cells, despite their close placement to T cells. This subtype of NK cells bridges the main NK cell cluster and the T cells in Figure 4.





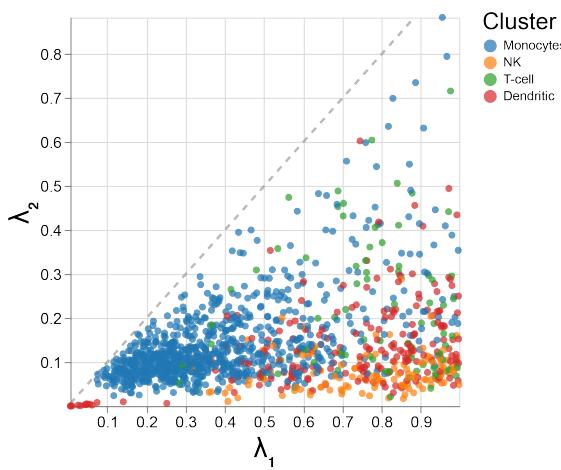
**Fig. 9:** Integrating scDEED dubious embeddings into visualizations made with the `distortions` package. A. The PBMC data with dubious cells flagged by scDEED. B. Hovering over the far left cluster reveals that the scDEED flagged cells have neighbors lying across multiple cell types that are distant in the embedding space. Our visualization functions are designed to accommodate alternative definitions of nonlinear embedding distortion.



**Fig. 10:** A graphical illustration of strategies used to flag fragmented neighborhoods. A. In the bin-based strategy, the original distances are partitioned into evenly-sized intervals  $I_l$ . Within each bin, the interquartile range of embedding distances is computed. Original vs. embedding distance pairs that do not fall within a factor of  $\sigma$  times the IQR of the median embedding distance for that bin are flagged as outliers  $\mathcal{O}_l$ . B. In either the bin or window-based strategies, samples with many neighbor links belonging to  $\cup_l \mathcal{O}_l$  are flagged as being the center of a fragmented neighborhood.

1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242

1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274



1275 **Fig. A1:** A zoomed-in version of Fig 4D. We have restricted to cells with  $\lambda_j^{(i)} < 1$ . A second mode of smaller,  
1276 less eccentric monocytes is visible in this view and contrasts with those that occupy the top right region of  
1277 Fig 4D. We also see a small cluster of dendritic cells with singular values near the origin, corresponding to  
1278 the small cluster placed near T cells in Figure 4A-C.  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296