

1 **cellstruct: Metrics scores to quantify the biological preservation between two**
2 **embeddings**

3 Jui Wan LOH¹, John F. OUYANG¹

4 ¹Centre for Computational Biology; Programme in Cardiovascular and Metabolic
5 Disorders, Duke-NUS Medical School, Singapore

6

7 ***Address for correspondence and reprint requests:**

8 Dr. John F. Ouyang, PhD

9 Centre for Computational Biology; Programme in Cardiovascular and Metabolic

10 Disorders, Duke-NUS Medical School, 8 College Rd, Singapore 169857

11 Email: john.ouyang@duke-nus.edu.sg

12

13 **KEYWORDS:** single-cell dimension reductions, global/local preservation, t-SNE,
14 UMAP, pairwise distances

15

16

17

18

19

20

21

22

23

24

25

1 **Abstract**

2 Single-cell transcriptomics (scRNA-seq) is extensively applied in uncovering biological
3 heterogeneity. There are different dimensionality reduction techniques, but it is unclear
4 which method works best in preserving biological information when creating a two-
5 dimensional embedding. Therefore, we implemented cellstruct, which calculates three
6 metrics scores to quantify the global or local biological similarity between a two-
7 dimensional and its corresponding higher-dimensional PCA embeddings at either
8 single-cell or cluster level. These scores pinpoint cell populations with low biological
9 information preservation, in addition to visualizing the cell-cell or cluster-cluster
10 relationships in the PCA embedding. Two study cases illustrate the usefulness of
11 cellstruct in exploratory data analysis.

12

13 **Background**

14 Single-cell transcriptomics (scRNA-seq) is increasingly used to interrogate the
15 biological heterogeneity and disease progression in different complicated biological
16 systems. Various dimensionality reduction methods [1-3] are employed to reduce the
17 high-dimensional features i.e. highly variable genes to two dimensions. t-SNE [4, 5]
18 and UMAP [6, 7] are commonly used in single-cell analysis packages (e.g. Seurat [8]
19 and Scanpy [9]) for removing technical noise while maintaining the biological signals
20 both globally and locally. These two-dimensional embeddings evaluate the
21 performance of multi-modal/batch integration, define cluster relationship including cell
22 type annotation, present gene expression changes between different conditions, and
23 infer trajectories driving developmental or disease progression [10-15], inherently
24 assuming the preservation of biological information from the underlying gene
25 expression space. However, dimension reduction often results in information loss

1 because it is difficult to represent all the complex biological variation in 2D. Moreover,
2 dimension reduction involves highly non-linear mathematical operations, introducing
3 different degree of transformation to different parts of data. Therefore, the assessment
4 of local and global structures preservation in the reduced embeddings from the
5 untransformed data is critical for the selection of most accurate representation of the
6 underlying variance for making rigorous biological inferences [1, 16-19], minimizing
7 the chances of data misinterpretations caused by distortions introduced in dimension
8 reduction [18, 20-22].

9
10 We present cellstruct to evaluate a reduced embedding's ability to retain global or local
11 relationships as compared to the reference embedding, at both single-cell and cluster
12 level. Cell populations with low metric scores indicate poor biological information
13 retention, guiding users to subset certain cell populations for closer inspection, or to
14 tune the dimension reduction hyperparameters for the generation of new embeddings
15 with better structure preservation. Thus, cellstruct is indispensable in scRNA-seq
16 exploratory analysis, by assessing the fidelity of different two-dimensional embeddings
17 and by revealing the underlying biological distances/relationships between cells or
18 clusters.

19

20 **Results and Discussion**

21 Cellstruct implements three metrics for assessing the preservation of global or local
22 relationships between a reduced (e.g. UMAP) and a reference (e.g. PCA) embeddings
23 at both cluster and single-cell level. It also provides heatmaps and dimension
24 reductions for visualizing the pairwise cell/cluster distances in different embeddings
25 (Figure 1A). We discovered that the metric assessing local relationships in the

1 preservation of the k-nearest neighbors does not contribute to better interpretation of
2 datasets, particularly in refining ambiguous/mixed cell types (detailed analysis in
3 supplementary text; Figure S1). Thus, we focus on the preservation of global
4 relationships at both single-cell and cluster level. To achieve this, we devised a global
5 single-cell (GS) score quantifying the correlation between the global position of a
6 single cell in the reduced and reference embeddings. Here, the global position of a
7 single cell is given by its distance against a fixed number of randomly selected
8 waypoint cells. Similarly, we devised a global cluster (GC) score by correlating the
9 cluster-cluster distances.

10

11 To illustrate these metrics, we applied cellstruct to a human liver cell atlas [23-27]
12 (Figures 1A-C) where we calculated the GS and GC scores for different embeddings
13 i.e. UMAP, t-SNE and Force-Directed Layout (FDL) (Figure S2). Note that we varied
14 the hyperparameters to generate a tuned UMAP embedding, which will be used for
15 the remaining analysis (further details in supplementary text; Figure S3). The cluster-
16 cluster distance heatmap, which is used to calculate the GC score, revealed that
17 hepatocytes are highly dissimilar from other cell types in the reference embedding
18 (refDR), but this is not accurately reflected in reduced embeddings, particularly in t-
19 SNE (Figure 1D). The lymphoid cells (B/T/NK cells, purple text in Figure 1D) were
20 clustered closely with mast cells in all embeddings, resulting in high GC scores in
21 general. Intuitively, one would expect the myeloid cells (macrophages/DCs, blue text)
22 to be transcriptomically similar to the lymphoid cells, as correctly reflected in the refDR.
23 However, the UMAP and FDL embeddings placed the myeloid closer to
24 fibroblast/endothelial cells instead, indicated by the shorter distances in the cluster-
25 cluster distance heatmap. Overall, all reduced embeddings recapitulated the

1 relationships within the “islands” of endothelial and fibroblasts respectively, but t-SNE
2 failed to reproduce it in the myeloid and lymphoid groups. The distance heatmap for t-
3 SNE also less resembled the refDR counterpart, resulting in the lowest mean GC
4 score amongst the different reduced embeddings (Figures 1D, S2).

5
6 At the single-cell level, we observed that the endothelial, lymphoid, and mast cells
7 showed high GS scores, preserving the reference DR’s structure in the UMAP
8 embedding (Figure 1E), and hepatocytes have the highest variance in GS score
9 among all cell types (Figure S4), suggesting substantial heterogeneity in the refDR
10 and/or UMAP embeddings. To investigate this heterogeneity, we sampled 284
11 hepatocytes (1%) to visualize the normalized pairwise distances between these cells
12 and 1,000 randomly selected waypoint cells. Three distinct groups were observed to
13 have very varied GS scores for all reduced embeddings. This variability in the GS
14 score is alleviated using the FDL embedding, supported by the higher GS scores,
15 particularly in Group3, due to the elongated projection of hepatocytes in FDL, which
16 provided a better separation between Group1 and Group3 cells (Figure S5).
17 Surprisingly, based on the cell-waypoint distance heatmap, Group2 and Group3 cells
18 are transcriptomically more similar to other cell types, than the remaining hepatocytes
19 (Figures 1F, S5).

20
21 We also interrogated the cholangiocytes due to their low GS scores. We sampled 428
22 cholangiocytes (~10%) and identified 81 “outliers” (Group1 and Group 2 cells) with
23 distinct cell-waypoint distance patterns (Figure S6). Three groups of cells were
24 identified: Group1 cells are very likely mislabeled as hepatocytes, since they clustered
25 together with the hepatocytes, Group2 cells represent a rare subpopulation, while

1 Group3 cells are the main cholangiocyte population. Also, we noticed that FDL
2 provided the worst embedding for cholangiocytes, particularly for Group3 cells,
3 suggesting that different embeddings might be needed when interpreting different cell
4 types (Figure S2). Overall, the FDL embedding is best at preserving the underlying
5 global cell-cell and cluster-cluster relationships (further analysis in supplementary text;
6 Figure S3).

7

8 We next applied cellstruct to peripheral blood single-cell data from COVID-19 patients
9 and healthy controls [28]. The T/NK cells showed relatively high GS scores, while the
10 remaining cell types have relatively lower GS scores, with COVID-CD14 Monocytes
11 showing the most variability in GS score (Figures 2A, S8A). Thus, we decided to
12 investigate the COVID-CD14 Monocytes further, sampled 1,657 cells (20%) for which
13 we plotted the cell-waypoint distance heatmaps (Figures 2B-C, S8B). A small
14 population of CD14 Monocytes did not cluster with the remaining cells of its type, and
15 this small population comprises of only COVID-CD14 Monocytes, suggesting that
16 these cells (black boxed) are different from other COVID-CD14 Monocytes (Figure
17 S8B), driving the differences within these 20% sampled cells. Four groups of cells
18 were identified, and they were well delineated in the tuned UMAP embedding of
19 COVID-CD14 Monocytes only (Figures 2C-D). We revealed that these groups are
20 significantly associated with the patient severity (Floor.NonVent, ICU.NonVent, and
21 ICU.Vent) ($p < 2.2 \times 10^{-16}$, Table S1), and Groups A-C shared similar gene set
22 enrichment with more severe patients (GroupA-ICU.Vent: neutrophil degranulation,
23 GroupC-ICU.Nonvent: interferon [IFN] signaling, GroupB: both biological processes)
24 (Figures 2D-G). In addition, the monocytes analysis from Wilk et al. corroborated with
25 our findings that very few cells within patients C2(1), C3(1), and C7(0) are found in

1 Groups B and C (Table S2), as these patients have an absence of predicted IFN and
2 IFN regulatory factor activities (Figure S9).

3

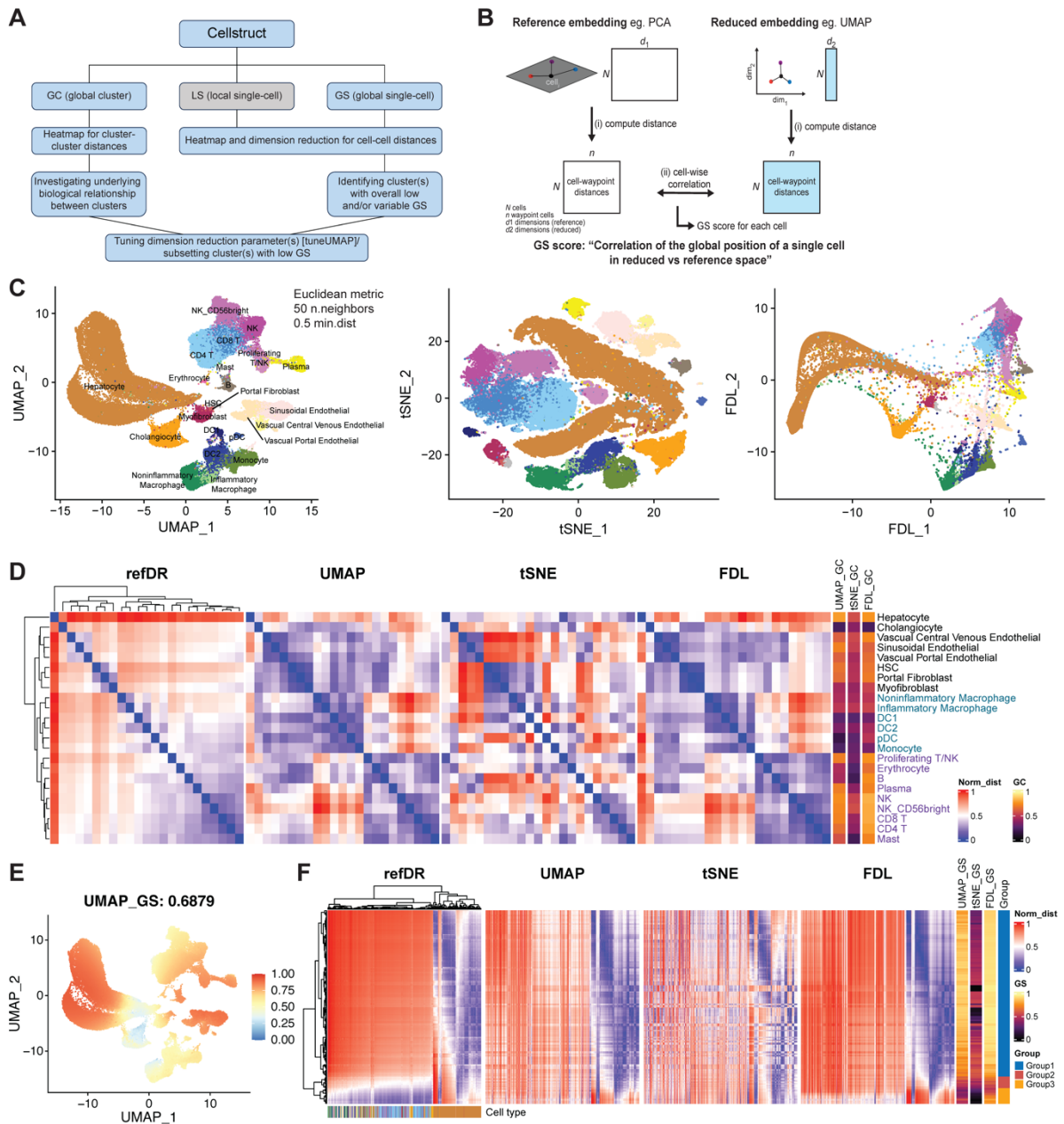
4 Finally, we compared cellstruct with scDEED using their 20 simulated datasets
5 (detailed comparison of simulated dataset 1 in supplementary text; Figure S10).
6 scDEED assesses the reliability of reduced embeddings and classifies the cells into
7 trustworthy and dubious [29]. We employed the same statistical approach to classify
8 trustworthy cells using our GS score. With the respective set of trustworthy cells
9 separately determined by cellstruct and scDEED, we measured the preservation of
10 neighboring information using K-nearest neighbors (KNN) and K-nearest clusters
11 (KNC) metrics. Cellstruct showed significantly higher KNN values in both t-SNE and
12 UMAP embeddings (mean difference ≥ 0.15 , $p < 1 \times 10^{-3}$), while scDEED exhibited a
13 higher KNC value in t-SNE (mean difference: 0.17, $p = 0.01$) (Figure S11), suggesting
14 that the trustworthy cells identified by cellstruct are more robustly preserved in the cell-
15 cell relationships than the scDEED counterpart. Moreover, the computational time for
16 cellstruct is one-third shorter than scDEED (cellstruct: ≤ 20 seconds, scDEED: ≤ 60
17 seconds) for both t-SNE and UMAP embeddings. Similar to scDEED, our tool is
18 applicable to different embedding methods, and only the embeddings (and cell
19 annotation for GC score) are required to run cellstruct.

20

21 **Conclusions**

22 In this work, we demonstrated the utility of cellstruct for exploratory data analysis given
23 different reduced embeddings. We provided several metrics and visualizations, which
24 quantify the biological information preservation in a reduced embedding, facilitating
25 the process in making observation-based biological interpretation. Users can now

1 evaluate cell-cell or cluster-cluster similarity in the underlying high-dimensional space,
2 and they are cautioned with cell populations comprising low/variable metric scores,
3 avoiding potential misinterpretations from the highly non-linear dimension reduction
4 procedure. More importantly, we hope cellstruct can bring awareness to the single-cell
5 analysis community that while 2D embeddings are useful for visualization and
6 interpretation, such embeddings often transform the underlying data to different extent
7 for different cell populations.
8



1

2 **Figure 1. The cellstruct package and application on human liver cell atlas.**

3 (A) Schematic of cellstruct package, which includes three metrics scores and various

4 visualizations to interrogate cell-cell or cluster-cluster relationships in single-cell data.

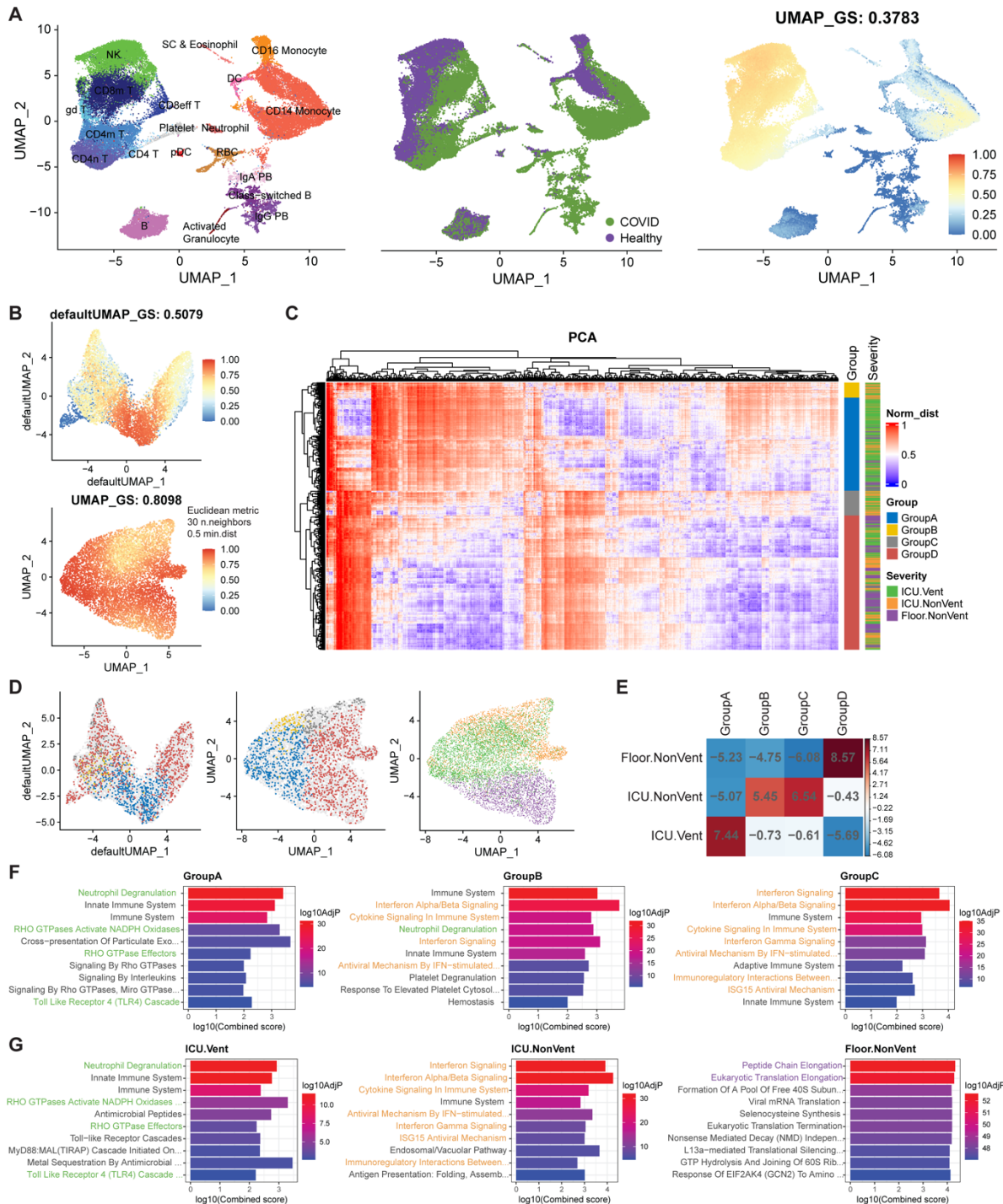
5 (B) Schematic showing the calculation of global single-cell (GS) metric. (C) Cell type

6 annotation of the liver cell atlas in UMAP, t-SNE, and FDL embeddings. (D) Heatmaps

7 illustrating the normalized cluster-cluster distances in reference (i.e. refDR) and

8 reduced (UMAP, t-SNE, and FDL) embeddings. GC scores of each cluster were

1 depicted as single-column heatmaps for each reduced embedding, and cluster labels
2 were colored by major grouping (blue: myeloid and purple: lymphoid and mast cell).
3 (E) The distribution of GS scores on UMAP projection, with the mean score indicated
4 in the title. (F) Heatmaps showing the cell-cell distances between 284 randomly
5 sampled hepatocytes and 1,000 waypoint cells, in different embeddings. Again, GS
6 scores were colored in the single-column heatmaps for each embedding. Three
7 groups of hepatocytes were identified and visualized in each reduced embedding in
8 Figure S5.
9



1

2 **Figure 2. Application of cellstruct to the peripheral immune atlas of healthy and**
 3 **COVID patients.**

4 (A) UMAP embeddings showing the cell type annotation, COVID/healthy samples, and
 5 GS scores (mean score in the title) of the immune atlas. (B-G) The downstream
 6 analysis was focused on the 8,285 COVID-CD14 Monocytes subset, due to their high

1 variability in GS score. (B) GS scores of COVID-CD14 Monocytes shown on default
2 and tuned UMAP embeddings. (C) Heatmap showing cell-cell distances between
3 1,657 randomly sampled COVID-CD14 Monocytes (same cells as Figure S8B) and
4 1,000 waypoint cells in reference PCA embedding, divided into four groups of cells
5 and annotated patient severity. (D) These four groups were delineated on default and
6 tuned UMAP embeddings, and patient severity was only shown on the tuned UMAP,
7 which does not show a separation of GroupD cells. (E) A contingency table of the four
8 monocyte groups and patient severity, colored by the residuals of Chi-squared test.
9 (F-G) Enriched pathways in different monocyte groups (F) and patient severity (G),
10 taken from up-regulated genes between the group/severity of interest against
11 remaining cells. Here, the randomly sampled COVID-CD14 Monocytes from Figure
12 2C were used, and GroupD monocytes were omitted due to small number of up-
13 regulated genes. Groups A-C are mainly enriched for neutrophil degranulation and
14 interferon signaling processes, which are respectively detected in ICU.Vent and
15 ICU.NonVent cells. Floor.Nonvent cells are enriched for translation process. Pathways
16 specific to patient severity were colored green, orange and purple.

17

18 **Methods**

19 **Cellstruct**

20 Cellstruct implements three metrics scores to quantify the preservation of global or
21 local relationships between a reduced and a reference embedding at both global and
22 local level. Each cell is assigned with GS (global single-cell) and optionally LS (local
23 single-cell) scores that measure the correlation of its global position between both
24 embeddings (Figure 1B) and the distances of its nearest neighbors within each
25 embedding respectively (Figure S1A). Similarly, each cluster (i.e. cell type) is assigned

1 with a GC (global cluster) score to describe the preservation of cluster-cluster
2 relationship.

3

4 Global single-cell (GS) metric score

5 GS score is the Pearson (default) cell-wise correlation of the cell-waypoint distances
6 calculated in the reference and reduced embeddings, given by this formula for cell i :

7

$$8 \quad GS_i = cor(dist_{ref}(cell_i, cell_{waypoint}), dist_{reduced}(cell_i, cell_{waypoint}))$$

9

10 where 10,000 cells are randomly selected as “waypoint” cells that serve to describe
11 the global position of each single cell, and pairwise distances between a single cell
12 and these waypoints are computed in both reference and reduced embeddings to give
13 the cell-waypoint distances. Note that the same set of waypoints is used across all
14 single cells. Thus, the GS score indicates how well the global location of each cell is
15 preserved from the reference to reduced embeddings.

16

17 Local single-cell (LS) metric score

18 LS score is the ratio of the mean reciprocal-squared-distance of the 30 nearest
19 neighbors (NN) in the reduced embedding to the 30NN in the reference embedding,
20 given by this formula for cell i :

21

$$22 \quad LS_i = \frac{mean\left(\frac{1}{dist_{ref}(cell_i, cell_{30NN, reduced})^2}\right)}{mean\left(\frac{1}{dist_{ref}(cell_i, cell_{30NN, ref})^2}\right)}$$

23

1 where the two sets of 30NN of the single cell (i.e. target) are determined using RANN
2 R package (v.2.6.1), in the reduced and reference embeddings respectively. Thus, the
3 LS score measures how far the reduced embedding NN are in the reference space,
4 and the denominator serves as a normalization factor, so that LS varies from 0 to 1.

5

6 Global cluster (GC) metric score

7 GC score is the Pearson (default) correlation between the cluster-cluster distances
8 calculated in the reference and reduced embeddings. Similar to the GS score, GC
9 score evaluates how well the global location of a cluster is preserved between
10 reference and reducing embeddings. Here, the distances are computed between the
11 centroid of each cluster. The centroids are determined by averaging each dimension
12 across all the cells in a cluster.

13

14 Overall, cellstruct takes in a Seurat object and requires users to specify the two input
15 embeddings (i.e. the reduced and reference embeddings) and optionally cell
16 annotations (e.g. cell type) for the GC score calculation. If cell annotations were not
17 provided, Seurat clusterID (i.e. `seurat_clusters`) will be used instead, giving rise to less
18 biological interpretable scores. For example, no biological inference could be made
19 with the observation of short distance between the arbitrarily labeled clusters 1 and 2,
20 as compared to the more meaningful labels e.g. CD4 and CD8 T cells. By default,
21 cellstruct will calculate both the GC and GS scores and add these scores into the
22 metadata of the returned Seurat object. In addition, dimension reduction plots
23 illustrating the distribution of GC and GS scores as well as the cluster annotation will
24 be generated for each reduced embedding (Figure S2).

25

1 Heatmaps and dimension reduction plots are provided as functions to illustrate the
2 normalized pairwise distances in different embeddings, for better visualization of the
3 metrics scores. For each cell (or each row in the cell-waypoint distance heatmap), the
4 cell-cell distances are normalized by the 99 percentile of the distances between the
5 single cell of interest and the remaining cells (or waypoint cells) in each embedding.
6 For the purpose of illustration, 1,000 randomly selected cells are used in plotting
7 heatmap. As for the GC metric visualization (i.e. the cluster-cluster distance heatmap),
8 cluster-cluster distances are normalized by the maximum distance in each embedding.

9

10 **Datasets**

11 Two single-cell datasets (human liver cell atlas from Azimuth
12 [<https://zenodo.org/record/7770308>] and peripheral immune atlas from Fred Hutch
13 [<https://atlas.fredhutch.org/fredhutch/covid/dataset/wilk>]) were used to demonstrate
14 our tool. These datasets were preprocessed by the authors, and the Seurat objects
15 were downloaded for our study.

16

17 Human liver cell atlas consists of 79,492 cells, which are composed of 23 cell types
18 from the broader groups of hepatocytes, cholangiocytes, fibroblasts, immune cells
19 (myeloid, T/NK, B, plasma cells, and erythrocytes), mast, endothelial, and stem cells
20 (Figure 1C). This reference atlas was collated from several liver studies [23-27],
21 involving 29 donors across a range of ages, for better liver cells annotation and
22 understanding of the liver-related diseases. We generated UMAP, t-SNE, and FDL
23 embeddings using the default hyperparameters, and as discussed in the
24 supplementary text, we also selected another UMAP embedding with different
25 hyperparameters (Euclidean metric, 50 n.neighbors, and 0.5 min.dist), which has the

1 highest mean GS score for the exploratory analysis. For the investigation on
2 cholangiocytes, we identified 78 “outliers” from all 3,634 cells via the clustering pattern
3 detected in the cell-waypoint reference distance heatmap (data not shown), and these
4 78 “outliers”, along with the 350 randomly chosen non-repeating cholangiocytes, were
5 used to plot the cell-waypoint distance heatmaps in Figure S6A.

6
7 This peripheral immune atlas was generated by Wilk et al. group to study the
8 pathophysiology of COVID-19. It consists of 44,172 cells, comprising of T/NK, myeloid,
9 B, red blood cells, plasmablasts, platelets, and granulocytes. The peripheral blood
10 mononuclear cells (PBMCs) were collected from seven hospitalized COVID-19
11 patients (four of them developed acute respiratory distress syndrome) and six controls
12 [28] (Figure 2A left, middle). For this study case, we focused on the COVID-19 CD14
13 Monocytes, as explained in the Results and Discussion section. 8,285 COVID-CD14
14 Monocytes were extracted, and UMAP embedding was tuned due to non-uniform
15 distribution of GS scores across the same cell population and the detection of a group
16 of partly separated cells with low GS scores on default UMAP (tuned UMAP
17 parameters: Euclidean metric, 30 n.neighbors, and 0.5 min.dist) (Figure 2B).
18 Pearson’s Chi-squared test was used to evaluate the association between the four
19 detected groups of COVID-CD14 Monocytes and patient severity within 1,657 cells
20 (Table S1, Figure 2E).

21

22 **Generation of UMAP, t-SNE, and FDL embeddings**

23 Both UMAP and t-SNE embeddings were implemented with the functions RunTSNE()
24 and RunUMAP() in the Seurat package. The hyperparameters that we tuned for UMAP
25 embedding are metric, n.neighbors, and min.dist in the RunUMAP() function. The

1 default hyperparameters of RunUMAP() and RunTSNE() are respectively cosine
2 metric; 30 n.neighbors; 0.3 min.dist and 30 perplexity value. As for FDL, it was
3 generated using scanpy.tl.draw_graph function in the Scanpy package. All the
4 arguments were kept as default, unless indicated as otherwise, when we ran these
5 functions.

6

7 **Stability analysis of cellstruct**

8 We performed a stability analysis on the number of waypoint cells for GS score
9 calculation, by varying the number of waypoints from 1K to 10K, with 1K increment on
10 UMAP, t-SNE, and FDL embeddings of human liver cell atlas. The GS scores are
11 independent of the number of waypoint cells from 1K to 10K cells, regardless of the
12 dimension reduction methods (Figure S7A). Hence, we use 10K cells for GS score
13 computation and 1K cells for heatmap illustration. In addition, we inspected the stability
14 of tuneUMAP function, by downsampling the number of cells in human liver cell atlas
15 to 5K, 8K, 10K, 20K, ..., 60K, and 70K cells respectively, studying the performance of
16 GS score across different dataset size. It was shown that GS scores are relatively
17 consistent, with at least 20K cells being sampled (Figure S7B).

18

19 **Comparison with scDEED**

20 scDEED (v0.1.0) is implemented as an R package. It provides a reliability score, which
21 is the Pearson correlation of the target's distances to its closest 50% neighboring cells
22 in both reference and reduced embeddings, and a classification of dubious and
23 trustworthy cells, by comparing to a null distribution of reliability score [29]. To compare
24 cellstruct with scDEED, the same dubious/trustworthy cell classification, implemented
25 by scDEED, was performed by comparing our GS score to a null distribution, which is

1 the GS score assigned to a permuted object generated by scDEED. Using the default
2 t-SNE and UMAP embeddings from the 20 simulated datasets generated by the
3 scDEED authors, the preservation of neighboring information was evaluated using the
4 same two metrics (K-nearest neighbors, KNN and K-nearest clusters, KNC) employed
5 by them. Trustworthy cells, which were separately identified by cellstruct and scDEED
6 in the t-SNE and UMAP, in the simulated datasets were retained (Table S3). Default t-
7 SNE and UMAP embeddings were regenerated for these trustworthy cells, and
8 dubious/trustworthy cells were re-classified for the evaluation of KNN and KNC
9 metrics, assessing the biological preservation and also the robustness of each tool in
10 identifying trustworthy cells (i.e. the first round of trustworthy cells). Paired t-test was
11 used to statistically evaluate the differences in mean KNN and KNC values (Figure
12 S11).

13

14 **Gene set enrichment analysis**

15 Differential expression analysis was performed among the four monocyte groups,
16 identified by cellstruct, and the three patient severity groups respectively from the 20%
17 COVID-CD14 Monocytes subset using FindAllMarker function in Seurat (v4.3.0). The
18 differential genes were tabulated in Tables S4 (monocyte groups) and S5 (patient
19 severity) respectively. Since Group D has only seven up-regulated genes (adjusted
20 $p < 0.05$), gene set enrichment analysis was performed in Groups A-C and the three
21 patient severity groups using enrichR (v3.2) [30-32] and Reactome 2022 database
22 [33]. The top 10 significant biological processes were shown on the bar plots in Figures
23 2F-G.

24

25 **Availability of data and materials**

1 The cellstruct R package can be installed from [https://github.com/the-ouyang-](https://github.com/the-ouyang-lab/cellstruct)
2 [lab/cellstruct](https://github.com/the-ouyang-lab/cellstruct). Data objects and codes for reproducing the figures and analysis can be
3 found at <https://github.com/the-ouyang-lab/cellstruct-reproducibility>.

4

5 **Acknowledgements**

6 Not applicable.

7

8 **Funding**

9 Both JWL and JFO are supported by the Singapore National Medical Research
10 Council (NMRC) under OF-YIRG funding (MOH-OFYIRG21nov-0004).

11

12 **Authors' contributions**

13 JWL and JFO wrote and edited the manuscript. JWL implemented the tool. JFO
14 supervised the work.

15

16 **Ethics approval and consent to participate**

17 Not applicable.

18

19 **Consent for publication**

20 Not applicable.

21

22 **Competing interests**

23 The authors declare that they have no competing interests.

24

25 **Supplementary Information**

1 Additional file 1: Supplementary text (detailed discussions of LS metric and tuning
2 UMAP hyperparameters, as well as comparative analysis of reduced embeddings and
3 cellstruct vs scDEED)

4 Additional file 2: Supplementary figures (**Figure S1**. Evaluation of local cell-cell
5 relationships via local single-cell (LS) metric in human liver cell atlas. **Figure S2**. An
6 example of cellstruct output using the human liver cell atlas. **Figure S3**. Four single
7 cells were selected as an illustration for the comparison between different reduced
8 embeddings. **Figure S4**. Distribution of GS scores for UMAP embedding in human
9 liver cell atlas. **Figure S5**. Dimension reductions showing three groups of hepatocytes.
10 **Figure S6**. Investigating cholangiocytes in the human liver cell atlas using cellstruct.
11 **Figure S7**. Stability analysis of GS metric using human liver cell atlas. **Figure S8**.
12 Applying cellstruct to the COVID peripheral immune atlas. **Figure S9**. Corroboration
13 of COVID-CD14 Monocyte groups, identified by cellstruct, with the author's original
14 analysis. **Figure S10**. Comparative analysis between scDEED and cellstruct using
15 Simulated Dataset 1. **Figure S11**. Scatterplot of KNC and KNN values of t-SNE and
16 UMAP embeddings in 20 simulated datasets.)

17 Additional file 3: Table S1. Number of cells for each group detected in different patient
18 severity level.

19 Additional file 4: Table S2. Number of cells for each sample in each COVID-CD14
20 Monocyte group.

21 Additional file 5: Table S3. Number of trustworthy, dubious, and neither cells classified
22 by scDEED and cellstruct respectively.

23 Additional file 6: Table S4. Differential genes expression analysis among four
24 monocyte groups detected in 20% of COVID-CD14 Monocytes.

1 Additional file 7: Table S5. Differential genes expression analysis among three patient
2 severity groups detected in 20% of COVID-CD14 Monocytes.

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

1 **References**

- 2 1. Sun S, Zhu J, Ma Y, Zhou X: **Accuracy, robustness and scalability of dimensionality**
3 **reduction methods for single-cell RNA-seq analysis.** *Genome Biol* 2019, **20**:269.
- 4 2. Xiang R, Wang W, Yang L, Wang S, Xu C, Chen X: **A Comparison for Dimensionality**
5 **Reduction Methods of Single-Cell RNA-seq Data.** *Front Genet* 2021, **12**:646936.
- 6 3. Ding J, Regev A: **Deep generative model embedding of single-cell RNA-Seq profiles on**
7 **hyperspheres and hyperbolic spaces.** *Nat Commun* 2021, **12**:2554.
- 8 4. Kobak D, Berens P: **The art of using t-SNE for single-cell transcriptomics.** *Nat Commun*
9 2019, **10**:5416.
- 10 5. van der Maaten L, Hinton G: **Visualizing Data using tSNE.** *J Mach Learn Res* 2008,
11 **9**:2579-2605.
- 12 6. Becht E, McInnes L, Healy J, Dutertre CA, Kwok IWH, Ng LG, Ginhoux F, Newell EW:
13 **Dimensionality reduction for visualizing single-cell data using UMAP.** *Nat Biotechnol*
14 2018.
- 15 7. McInnes L, Healy J, Melville J: **UMAP: Uniform Manifold Approximation and**
16 **Projection for Dimension Reduction.** *arXiv* 2018.
- 17 8. Hao Y, Hao S, Andersen-Nissen E, Mauck WM, 3rd, Zheng S, Butler A, Lee MJ, Wilk AJ,
18 Darby C, Zager M, et al: **Integrated analysis of multimodal single-cell data.** *Cell* 2021,
19 **184**:3573-3587 e3529.
- 20 9. Wolf FA, Angerer P, Theis FJ: **SCANPY: large-scale single-cell gene expression data**
21 **analysis.** *Genome Biol* 2018, **19**:15.
- 22 10. Dou J, Liang S, Mohanty V, Miao Q, Huang Y, Liang Q, Cheng X, Kim S, Choi J, Li Y, et al:
23 **Bi-order multimodal integration of single-cell data.** *Genome Biol* 2022, **23**:112.

- 1 11. Hie B, Bryson B, Berger B: **Efficient integration of heterogeneous single-cell**
2 **transcriptomes using Scanorama.** *Nat Biotechnol* 2019, **37**:685-691.
- 3 12. Peyvandipour A, Shafi A, Saberian N, Draghici S: **Identification of cell types from single**
4 **cell data using stable clustering.** *Sci Rep* 2020, **10**:12349.
- 5 13. Szabo PA, Levitin HM, Miron M, Snyder ME, Senda T, Yuan J, Cheng YL, Bush EC, Dogra
6 P, Thapa P, Farber DL, Sims PA: **Single-cell transcriptomics of human T cells reveals**
7 **tissue and activation signatures in health and disease.** *Nat Commun* 2019, **10**:4706.
- 8 14. Saelens W, Cannoodt R, Todorov H, Saeys Y: **A comparison of single-cell trajectory**
9 **inference methods.** *Nat Biotechnol* 2019, **37**:547-554.
- 10 15. Cao J, Spielmann M, Qiu X, Huang X, Ibrahim DM, Hill AJ, Zhang F, Mundlos S,
11 Christiansen L, Steemers FJ, Trapnell C, Shendure J: **The single-cell transcriptional**
12 **landscape of mammalian organogenesis.** *Nature* 2019, **566**:496-502.
- 13 16. Wang X, Yang L, Wang YC, Xu ZR, Feng Y, Zhang J, Wang Y, Xu CR: **Comparative analysis**
14 **of cell lineage differentiation during hepatogenesis in humans and mice at the single-**
15 **cell transcriptome level.** *Cell Res* 2020, **30**:1109-1126.
- 16 17. Liu X, Ouyang JF, Rossello FJ, Tan JP, Davidson KC, Valdes DS, Schroder J, Sun YBY, Chen
17 J, Knaupp AS, et al: **Reprogramming roadmap reveals route to human induced**
18 **trophoblast stem cells.** *Nature* 2020, **586**:101-107.
- 19 18. Chari T, Pachter L: **The specious art of single-cell genomics.** *PLoS Comput Biol* 2023,
20 **19**:e1011288.
- 21 19. Heiser CN, Lau KS: **A Quantitative Framework for Evaluating Single-Cell Data**
22 **Structure Preservation by Dimensionality Reduction Techniques.** *Cell Rep* 2020,
23 **31**:107576.

- 1 20. Kobak D, Linderman GC: **Initialization is critical for preserving global data structure in**
2 **both t-SNE and UMAP.** *Nat Biotechnol* 2021, **39**:156-157.
- 3 21. Alquicira-Hernandez J, Powell JE, Phan TG: **No evidence that plasmablasts**
4 **transdifferentiate into developing neutrophils in severe COVID-19 disease.** *Clin Transl*
5 *Immunology* 2021, **10**:e1308.
- 6 22. Cooley SM, Hamilton T, Aragoes SD, Ray JCJ, Deeds EJ: **A novel metric reveals**
7 **previously unrecognized distortion in dimensionality reduction.** *bioRxiv* 2019.
- 8 23. Aizarani N, Saviano A, Sagar, Mailly L, Durand S, Herman JS, Pessaux P, Baumert TF,
9 Grun D: **A human liver cell atlas reveals heterogeneity and epithelial progenitors.**
10 *Nature* 2019, **572**:199-204.
- 11 24. Ramachandran P, Dobie R, Wilson-Kanamori JR, Dora EF, Henderson BEP, Luu NT,
12 Portman JR, Matchett KP, Brice M, Marwick JA, et al: **Resolving the fibrotic niche of**
13 **human liver cirrhosis at single-cell level.** *Nature* 2019, **575**:512-518.
- 14 25. Zhang M, Yang H, Wan L, Wang Z, Wang H, Ge C, Liu Y, Hao Y, Zhang D, Shi G, et al:
15 **Single-cell transcriptomic architecture and intercellular crosstalk of human**
16 **intrahepatic cholangiocarcinoma.** *J Hepatol* 2020, **73**:1118-1130.
- 17 26. MacParland SA, Liu JC, Ma XZ, Innes BT, Bartczak AM, Gage BK, Manuel J, Khuu N,
18 Echeverri J, Linares I, et al: **Single cell RNA sequencing of human liver reveals distinct**
19 **intrahepatic macrophage populations.** *Nat Commun* 2018, **9**:4383.
- 20 27. Payen VL, Lavergne A, Alevra Sarika N, Colonval M, Karim L, Deckers M, Najimi M,
21 Coppieters W, Charloteaux B, Sokal EM, El Taghdouini A: **Single-cell RNA sequencing**
22 **of human liver reveals hepatic stellate cell heterogeneity.** *JHEP Rep* 2021, **3**:100278.

- 1 28. Wilk AJ, Rustagi A, Zhao NQ, Roque J, Martinez-Colon GJ, McKechnie JL, Ivison GT,
2 Ranganath T, Vergara R, Hollis T, et al: **A single-cell atlas of the peripheral immune**
3 **response in patients with severe COVID-19.** *Nat Med* 2020, **26**:1070-1076.
- 4 29. Xia L, Lee C, Li JJ: **scDEED: a statistical method for detecting dubious 2D single-cell**
5 **embeddings and optimizing t-SNE and UMAP hyperparameters.** *bioRxiv* 2023.
- 6 30. Chen EY, Tan CM, Kou Y, Duan Q, Wang Z, Meirelles GV, Clark NR, Ma'ayan A: **Enrichr:**
7 **interactive and collaborative HTML5 gene list enrichment analysis tool.** *BMC*
8 *Bioinformatics* 2013, **14**:128.
- 9 31. Kuleshov MV, Jones MR, Rouillard AD, Fernandez NF, Duan Q, Wang Z, Koplev S, Jenkins
10 SL, Jagodnik KM, Lachmann A, et al: **Enrichr: a comprehensive gene set enrichment**
11 **analysis web server 2016 update.** *Nucleic Acids Res* 2016, **44**:W90-97.
- 12 32. Xie Z, Bailey A, Kuleshov MV, Clarke DJB, Evangelista JE, Jenkins SL, Lachmann A,
13 Wojciechowicz ML, Kropiwnicki E, Jagodnik KM, Jeon M, Ma'ayan A: **Gene Set**
14 **Knowledge Discovery with Enrichr.** *Curr Protoc* 2021, **1**:e90.
- 15 33. Gillespie M, Jassal B, Stephan R, Milacic M, Rothfels K, Senff-Ribeiro A, Griss J, Sevilla
16 C, Matthews L, Gong C, et al: **The reactome pathway knowledgebase 2022.** *Nucleic*
17 *Acids Res* 2022, **50**:D687-D692.

18