

# Incorporating cell hierarchy to decipher the functional diversity of single cells

Lingxi Chen<sup>1</sup> and Shuai Cheng Li<sup>1,\*</sup>

<sup>1</sup>Department of Computer Science, City University of Hong Kong, Hong Kong, China

\* Corresponding: shuaicli@cityu.edu.hk

Cells possess functional diversity hierarchically. However, most single-cell analyses renounce the nested structures while detecting and visualizing the functional diversity. Here, we incorporate cell hierarchy to study functional diversity at subpopulation, club (i.e., sub-subpopulation), and cell layers. Accordingly, we implement a package, SEAT, to construct cell hierarchies utilizing structure entropy by minimizing the global uncertainty in cell-cell graphs. With cell hierarchies, SEAT deciphers functional diversity in 36 data sets covering scRNA, scDNA, scATAC, and scRNA-scATAC multiome. First, SEAT finds optimal cell subpopulations with high clustering accuracy. It identifies cell types or fates from omics profiles and boosts accuracy from 0.34 to 1. Second, SEAT detects insightful functional diversity among cell clubs. The hierarchy of breast cancer cells reveals that the specific tumor cell club drives *AREG-EGFT* signaling. We identify a co-accessibility network of dense *cis*-regulatory elements specified by the cell club for GM12878. Third, the cell order from the hierarchy infers periodic pseudo-time of cells, improving accuracy from 0.79 to 0.89. Moreover, we incorporate cell hierarchy layers as prior knowledge to refine nonlinear dimensionality reduction, enabling us to visualize hierarchical cell layouts in low-dimensional space.

## Introduction

Cells in the biological system own functional diversity hierarchically, which signifies cell types or states during development, disease, and evolution, up to the biosystem (1, 2). The heterogeneity of the cell is observed with nested structures (3). In the tumor microenvironment, infiltrated lymphocytes include B cells and T cells. Furthermore, T cells can be classified into helper T cells and cytotoxic T cells (4). Specific expression of the marker CD4 and CD8 will strengthen intra-similarity within helper and cytotoxic T cells, respectively, resulting in nested cell structures. The cellular heterogeneity raised by tumor evolution presents another instance (5, 6). The copy number gain, neutral, and loss classify tumor cells into aneuploid, diploid, and hypodiploid groups, respectively. Fluctuations of copy numbers in focal genome regions further categorize tumor cells into amplification or deletion subtypes. The cell cycle is a rudimentary biological process for cell replications (7). Human cells undergo a cycle G1 - S - G2/M - G1 over a 24-hour period, which signifies that the cycling cells have three flat phase labels (G1, S, and G2/M). In addition, the cycling cells have a hierarchical order (pseudo-time) that records the exact timing in the G1, S, and G2/M phases. The recent maturation of single-cell sequencing technology

offers opportunities to profile large-scale single cells for their transcriptomics (8), genomics (5), epigenomics (9), etc. These technologies have blossomed revolutionary insights into cellular functional diversity under the aegis of assigning cells with similar molecular characteristics to the same group (1, 2). However, most existing clustering tools generate flat cell group (10–14). Moreover, the periodic pseudo-time inference tools neglect the hierarchical order of cycling cells (15–18). Renunciation of the underlying nested structures of cells prevents full-scale detection of cellular functional diversity.

To address the issue, we incorporate *cell hierarchy* to illustrate the nested structure of cellular functional diversity. Cell hierarchy is a tree-like structure with multiple layers that capture cellular heterogeneity. From the root to the tips, the cellular heterogeneity decays. This study focuses on four main layers: global, subpopulation, club, and cell. The global layer is the root that exemplifies the whole cell population, e.g., immune cells. In contrast, the cell groups in the second and third layers resemble *cell subpopulations* and *cell clubs*, respectively. The cell subpopulation is a broad category of cells, such as B cells and T cells (4). Cell clubs within one cell subpopulation catalog the cellular heterogeneity in a finer resolution; that is, the cells share high functional similarity within a single cell club. For example, T cell subpopulation owns CD4 T cell and CD8 T cell clubs (4). The tip layer holds individual cells carrying *cell orders*, which signify the dynamic nuance of cell changes within a cell club, e.g., cellular heterogeneity varies along a periodic time course for cells undergoing a cell cycle process (7).

The actual cell hierarchy is difficult to determine; here, we develop SEAT, Structure Entropy hierArChy deTectiOn, to build a pseudo cell hierarchy leveraging structure entropy (19–22) by diminishing the global uncertainty in cell-cell graphs. SEAT constructs cell hierarchies using a raw or dimensionally reduced single-cell molecular profile as inputs, and computes the global-subpopulation-club-cell layers from the hierarchies. We apply SEAT to 36 data sets that cover single-cell RNA (scRNA), single-cell DNA (scDNA), single-cell assay for transposase-accessible chromatin (scATAC) and scRNA-scATAC multiome. SEAT detects the functional diversity of these single-cell omics data with cell hierarchy from three perspectives: cell subpopulation detection, cell club investigation, and periodic cell cycle pseudo-time inference.

Visualizing the functional diversity of single cells is essen-

tial since the visual inspection is the most direct approach to studying the structure and pattern of cells. Nonlinear dimensional reduction is a trending visualisation method for high-dimensional biological data (23). Nevertheless, state-of-the-art single-cell visualization tools neglect the nested structure of cells by merely capturing at most two levels (global or local) of cell structures (24–26). To tackle the issue, SEAT provides a component to embed the cells into a low-dimensional space by incorporating the multiple layers from the cell hierarchy as prior knowledge. Experiments demonstrate that SEAT consistently visualizes the hierarchical layout of these cells in the two-dimensional space for the above single-cell datasets.

## Result

**Overview of SEAT.** SEAT builds a cell hierarchy annotated with global-subpopulation-club-cell layers computationally from single-cell data (Fig. 1). First, SEAT constructs a pair of dense and sparse cell-cell similarity graphs with a raw or dimensionally reduced single-cell molecular profile as input (Fig. 1 A). Second, we detect cell clubs, determine the order of cells within each cell club, and build the pseudo club hierarchies by minimizing the structure entropy of the sparse graph with agglomerative and divisive heuristics, namely, Agglo(club), Agglo(order), Agglo, Divisive(club), Divisive(order), Divisive (Fig. 1B, Online Methods). Next, we use dynamic programming to find optimal subpopulations from agglomerative and divisive hierarchies, namely, Agglo(sub) and Divisive(sub). We choose the hierarchy carrying the lower subpopulation structure entropy as the final cell hierarchy (Fig. 1C, Online Methods). After that, SEAT outputs the final cell hierarchy carrying with subpopulations, clubs, and orders, namely, SEAT(sub), SEAT(club), and SEAT(order) (Fig. 1A). Furthermore, by incorporating hierarchical cell partition layers, SEAT provides a component, SEAT(viz), to embed cells in a low-dimensional space while preserving their nested structures for improved visualization and interpretation (Fig. 1A).

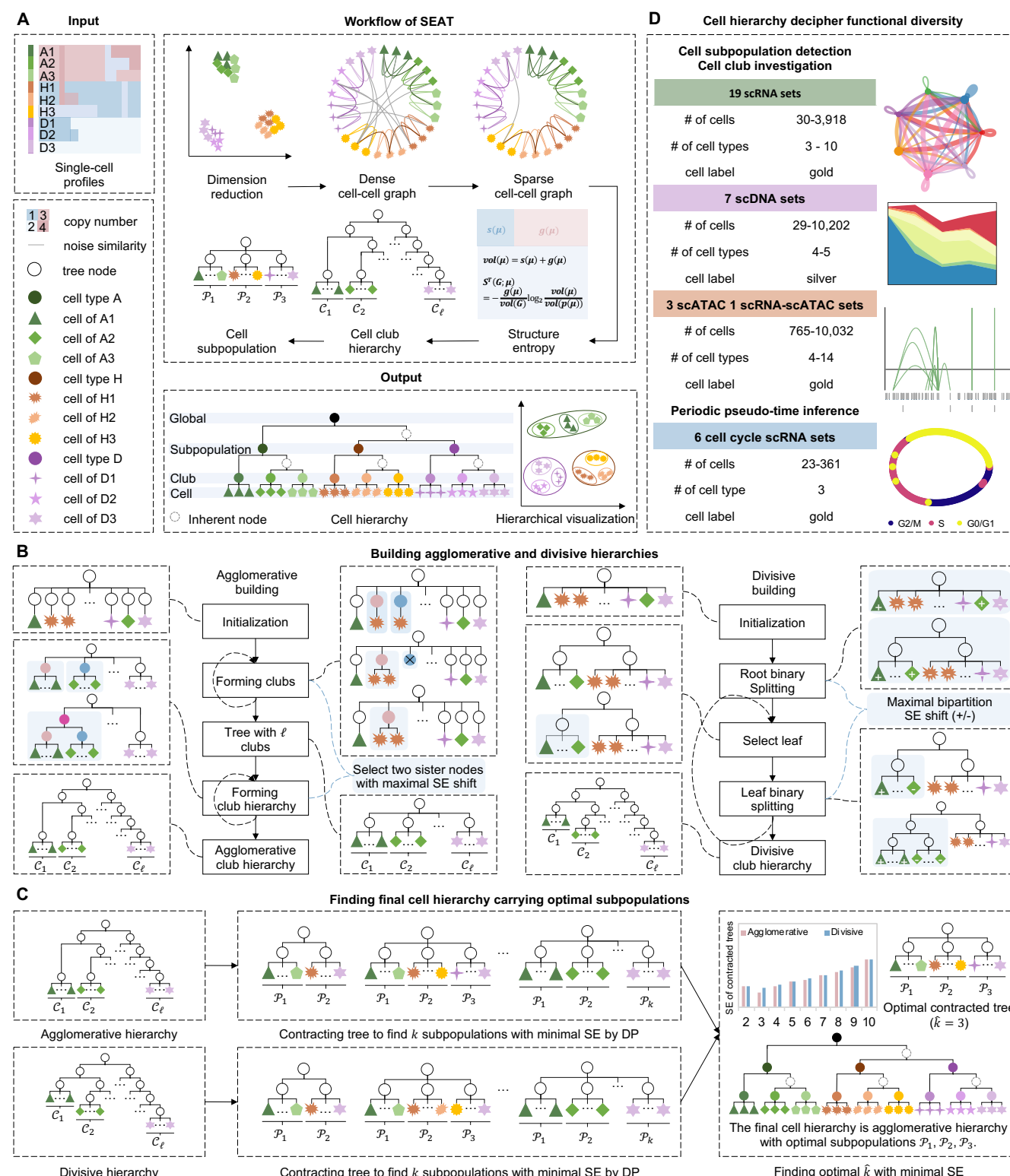
To detect cell subpopulations, some clustering methods require the number of clusters prespecified, while others can determine the number of clusters automatically. The SEAT package supports both. Our package requires no prespecified number of cluster by default, that is, SEAT(sub). If the number of clusters required is as  $k$ , we denote the method as SEAT( $k$ ). When the context is clear, we refer to them as predefined- $k$  and auto- $k$  modes, respectively.

**Cell hierarchy catalogs functional diversity at the subpopulation and club level from scRNA data.** We applied SEAT to nineteen scRNA datasets carrying gold standard cell type labels. The first nine sets are cell line mixtures, including p3cl (27), 3Line-qPCR (28), sc\_10x, sc\_celseq2, sc\_dropseq, sc\_10x\_5cl, sc\_celseq2\_5cl\_p1, sc\_celseq2\_5cl\_p2, and sc\_celseq2\_5cl\_p3 (29). We have four datasets Yan (30), Deng (31), Baise (32), and Goolam (33) which sequence single cells from human or mouse embryos at different stages of development (zygote, 2-cell, early 2-cell, mid 2-cell, late 2-cell, 4-cell, 8-cell, 16-cell, 32-cell, early blast, mid blast, and late blast). The last six datasets are Koh (34), Kumar (35), Trapnell (36), Blakeley (37), Kolodziejczyk (38), and Xin (39), which profile different cell types in single-cell resolution.

2-cell, early 2-cell, mid 2-cell, late 2-cell, 4-cell, 8-cell, 16-cell, 32-cell, early blast, mid blast, and late blast). The last six datasets are Koh (34), Kumar (35), Trapnell (36), Blakeley (37), Kolodziejczyk (38), and Xin (39), which profile different cell types in single-cell resolution.

To access the efficacy of SEAT in cell subpopulations detection, we utilize the adjusted rand index (ARI) and adjusted mutual information (AMI) as clustering accuracy and benchmark SEAT with state-of-the-art clustering tools (spectral clustering (10), K-means (11), hierarchical clustering (12), Louvain (13), and Leiden (14)) with predefined- $k$  and auto- $k$  modes (Online Methods). In predefined- $k$  mode, SEAT( $k$ ) demonstrates comparable or higher clustering accuracy compared to other clustering baselines on most datasets (Fig. 2A). Notably, Louvain( $k$ ) and Leiden( $k$ ) are unable to generate a clustering that exactly matches the number of ground truth labels after 20 different resolution trials for the Goolam and Kolodziejczyk (Fig. 2A). Under the auto- $k$  mode, SEAT(sub) outperforms Louvain and Leiden on all nineteen sets. The clustering accuracies of SEAT(sub) are comparable to or better than the best clustering results with predefined- $k$  clustering tools with the ground truth cluster number provided. This is attributed to the fact that SEAT(sub) finds a cluster number close to the ground truth (Fig. 2 B). Louvain and Leiden have the lowest clustering accuracy because they prefer more clusters. The two-dimensional data embedded by UMAP from raw single-cell expression profiles are inputs of all clustering tools; and the visualizations of them show that the ground truth labels are mixed for the majority of datasets, explaining the low clustering accuracy of both predefined- $k$  and auto- $k$  clustering tools.

SEAT relies on hierarchical structures to study cellular functional diversity. We leverage differential gene expressions to investigate the biological interpretations of these hierarchies. Differentially expressed genes ( $p < 0.05$ ) between cell hierarchy clubs reveal distinct patterns that match ground truth cell subpopulations. Furthermore, visible marker gene patterns reveal the functional diversity among cell clubs within one cell subpopulation. We focus on the top five differentially expressed genes for each data set. As the subpopulation detection accuracy of agglomerative hierarchy is 1 for p3cl dataset, we investigate the functional diversity revealed from agglomerative hierarchy other than the divisive hierarchy. The agglomerative hierarchy revealed three cell subpopulations for p3cl, which correspond to the three ground truth cell types, basal (*KRT81*), luminal (*TFF1*), and fibroblast (*COL1A2* and *VIM*) (Fig. 2D). We observe that each of the basal, luminal, and fibroblast has two major subclasses, controlled by the expression of cell cycle genes (*HIST1H4C*, *CDC20*, *CCNB1*, and *PTTG1*). Cell-cell communication analysis finds a total of 109 significant ( $p < 0.05$ ) ligand-receptor (LR) pair interactions among seven agglomerative hierarchy clubs for breast cancer basal-like epithelial cell line in p3cl, the LR interactions belong to nine signaling pathways AGRN, CD99, CDH, EGF, JAM, LAMININ, MK, NECTIN, and NOTCH (Fig. 2D). In particular, there is a distinct breast cancer cell club (basal-club0) that drives *AREG-EGFR*, an oncogenic signal-



**Fig. 1.** The schematic overview of SEAT. **A** The workflow of SEAT. **B** The algorithm of agglomerative and divisive hierarchy building. **C** The algorithm of finding final cell hierarchy carrying optimal subpopulations. **D** The summary of experimental settings.

ing (40) in breast cancer, to all basal cells, resulting in a high level of *AREG* activated *EGFR* expression (Fig. 2E). The two cell clubs from the luminal subpopulation have six significant ( $p < 0.05$ ) LR interactions involving MK, SEMA3, and CDH signaling pathways. The fibroblast has three significant ( $p < 0.05$ ) LR interactions, including two signaling pathways FN1 and ncWNT. The cell club fibro-club10 release *WNT5B* and then bind *FZD7* from fibro-club9, consistent with the observation that ncWNT is the predominant signaling pathway in skin fibroblasts (41).



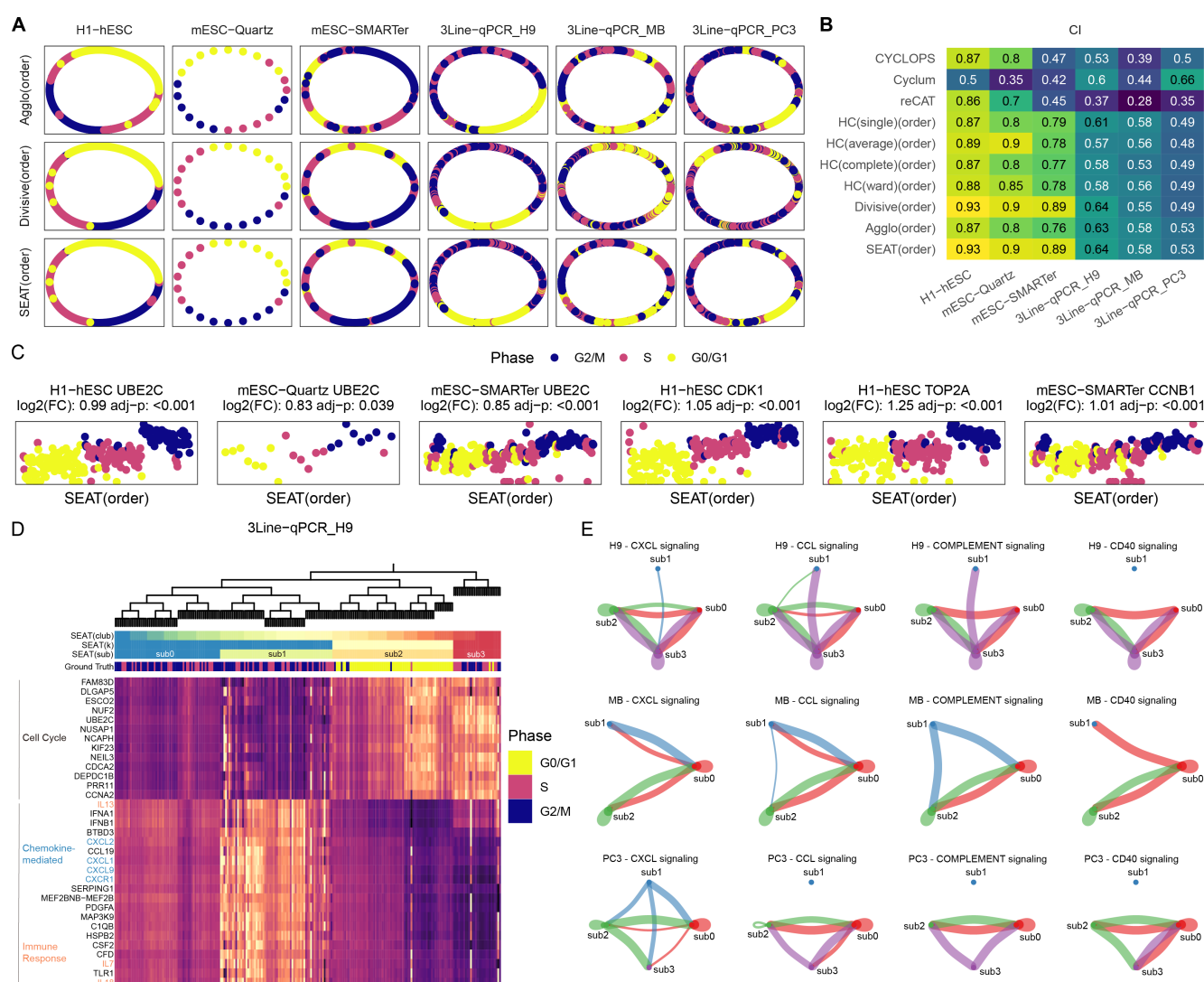


Visualizations of two-dimensional data by UMAP from raw single-cell expression profiles reveal a dense layout. The ground truth cell subpopulations are indistinctly separated in some high clustering accuracy datasets, and the cell clubs are densely arranged in each subpopulation clump. Here, we perform a visualization refinement to check whether SEAT hierarchical visualization eliminates the dense layout of clubs. We use the cell-cell graph constructed by SEAT as input and execute SEAT hierarchical visualization, UMAP, TSNE, and PHATE, independently. In Fig. 2E-F, SEAT hierarchical visualization, UMAP, TSNE, and PHATE separate the ground truth cell type for most datasets. It should be noted that the patterns from SEAT(viz), UMAP, TSNE, and PHATE also correspond to the subpopulation layer annotations, validating SEAT subpopulation finding efficacy. At the cell club level, SEAT(viz) show a clear layout of cell clumps that correspond to the cell hierarchy; each cell club owns a distinct clump, and the distance between clubs belonging to the same subpopulation is within proximity. Although UMAP, TSNE, and PHATE capture the local structures of the clubs, the cell clubs are unclearly segregated.

**Cell hierarchy deciphers periodic cell cycle pseudo-time from single-cell data.** We collect six scRNA cell cycle datasets, H1-hESC (42), mESC-Quartz (43), mESC-SMARTer (44), 3Line-qPCR\_H9, 3Line-qPCR\_MB, and 3Line-qPCR\_PC3 (28) with gold standard G0/G1, S, or G2/M stages and build the cell hierarchies. In predefined-k and auto-k clustering benchmarking, SEAT illustrates higher or comparable clustering accuracy in the six datasets. SEAT predicts the optimal number of clusters closest to ground truth three, while Leiden and Louvain generally predict more clusters than SEAT. Further investigation shows that ground truth labels are mixed or not distinctly separated in two-dimensional data by UMAP for all datasets, explaining the poor performance of 3Line-qPCR data. Likewise, hierarchical visualization plots depict nested layouts corresponding to the cell hierarchies in visualization refinement experiments. If we order the cells in cell cycle progress, cells from the same phase share higher similarity and they should be lined up adjacently. Thus, the cell order obtained from the ideal hierarchy could present a periodic pseudo-time order for cell cycle data. We visualize the cell order periodically with an oval plot. The placements of the cells in the oval represent their pseudo-time in the cell cycle (Fig. 3A). We access the cell ordering accuracy with the change index (CI), which computes how frequently the gold standard cell cycle phase labels switch along the cell order. The benchmark methods are four conventional HC strategies (12) that offer a cell order. We also recruit state-of-the-art tools dedicating to predict the cell cycle pseudo-time, CYCLOPS (15), Cyclum (16), reCAT (17), and CCPE (18). CCPE fails the tasks. SEAT demonstrates the highest ordering accuracy for all datasets, except for 3Line-qPCR\_PC3, where SEAT wins the top two (Fig. 3B). Hence, this suggested that cell hierarchy obtained from SEAT facilitates the cell cycle pseudo-time order inference. SEAT orders cells in H1-hESC, mESC-Quartz, and mESC-

SMARTer alongside the oval that closely matches the G0/G1-S-G2/M cycle (Fig. 3A). Differential expression analysis among ground truth phases reveals distinct cell cycle phase markers. These visible cell cycle marker patterns remain consistent when rearranging with SEAT cell order. The top 20 differential expression genes ( $p < 0.05$ ) for hESC and mESC cells include well-known cell cycle markers *UBE2C*, *TOP2A*, *CDK1*, and *CCNB1*. Their expressions rise progressively with SEAT recovered pseudo-time order and are peaked with significant fold changes at the M phase (Fig. 3C). In H9, MB, and PC3 cell lines, cells in the S and G2/M phases are partially arranged according to the exact time course (Fig. 3A). The differential expression makers of ground truth phases show that there are subpatterns within the S and G2/M phases and similar patterns between the S and G2/M phases, suggesting the cause of poor performance in pseudo-time ordering. Interestingly, after rearranging marker expression with SEAT, we observe distinct marker gene patterns among SEAT discovered cell subpopulations. For the H9 cell line, SEAT detected four cell subpopulations (Fig. 3D), G0/G1 phase corresponds to sub2. Cell cycle S and G2/M phases have three cell subpopulations, sub0, sub1, and sub3. The top 20 differential expression genes ( $p < 0.05$ ) have two groups (Fig. 3D). The genes from the first group enriched GO cell cycle signaling pathways. The genes from the second group enriched in GO chemokine-mediated signaling and immune response pathways with CXC and IL gene family, respectively. We demonstrate the top 20 differential expression genes for MB and PC3. Finally, we verify the cellular interactions among cell subpopulations with cell-cell communication analysis. We find a total of 124, 87, and 77 significant ( $p < 0.05$ ) LR pair interactions among cell subpopulations for H9, MB, and PC3 cell lines, respectively. All datasets exhibit CXCL, CCL, COMPLEMENT, and CD40 signaling interactions among cell subpopulations (Fig. 3E).

**Cell hierarchy detects rare subclones on scDNA data.** With seven scDNA datasets, SEAT catalogs the clonal subpopulations in solid tumor and circulating tumor cells. It identifies the CNV substructures in neuron and gamete cells. Owing to the unique characteristics of CNV profiles, we only adopt SEAT agglomerative hierarchy to investigate the functional diversity of CNV substructures. Navin *et al.* profiled 100 cells from a genetically heterogeneous (polygenetic) triple-negative breast cancer primary lesion T10 (45). Fluorescence-activated cell sorting (FACS) analysis confirmed that T10 carried four main cell subpopulations: diploid (D), hypodiploid (H), aneuploid A (A1), and aneuploid B (A2). Furthermore, Navin *et al.* reported pseudo-diploid cells (P) with varying degrees of chromosome gains and losses from diploids. They are unrelated to the three tumor cell subgroups (H, A1, and A2) (45). Therefore, given whole-genome single-cell CNV profiles as input, we verify whether SEAT and the state-of-the-art clustering tools identify the four major cell groups and the distinct pseudo-diploid cell group (Fig. 4A). In predefined-k mode, SEAT agglomerative hierarchy successfully recognizes five cell subpopulations consistent with the patterns of CNV pro-

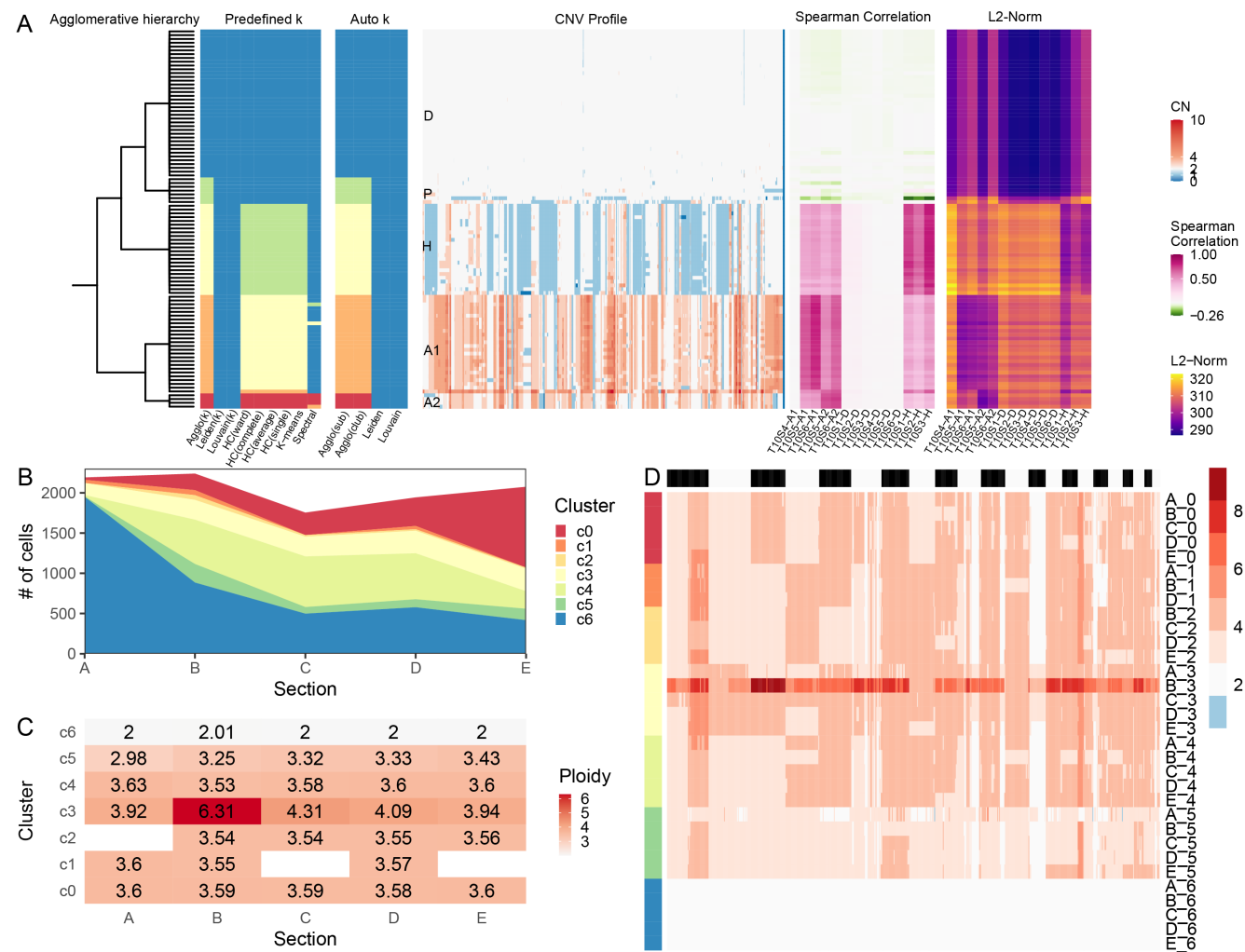


**Fig. 3.** Applying SEAT on six scRNA cell cycle datasets. **A.** The oval visualization of cell pseudo-time. From left to right are H1-hESC, mESC-Quartz, mESC-SMARTer, 3Line-qPCR\_H9, 3Line-qPCR\_MB, and 3Line-qPCR\_PC3. From top and bottom are cell orders obtained from agglomerative hierarchy (Agglo(order)), divisive hierarchy ((Divisive(order))), and SEAT cell hierarchy ((SEAT(order))). **B.** The accuracy of cell pseudo-time order is measured by change index (CI) for hierarchy-building tools. HC(single)(order), HC(average)(order), HC(complete)(order), and HC(ward)(order): the cell order from hierarchical clustering with single, average, complete, and ward linkage. **C.** The normalized expression of M phase marker genes alongside the SEAT cell order. **D.** The top 20 differentially expressed genes in G0/G1, S, and G2/M phases for p3c1, structured with SEAT. SEAT(club): the cell clubs from SEAT cell hierarchy. SEAT(k): the cell subpopulations from SEAT cell hierarchy in predefined-k mode. SEAT(sub): the optimal subpopulations from SEAT cell hierarchy in auto-k mode. **E.** The cell-cell communication among SEAT cell subpopulations for H9, MB, and PC3 cell lines.

files. From top to bottom, the ranks are cancer normal cell group (D), pseudo-diploid cell subgroups (P), subgroups H, and two tumor aneuploid groups, A1 and A2 (Fig. 4A). Leiden(k) and Louvain(k) fail with the same cell-similarity graph as input. Four HC strategies and K-means fail to distinguish the four pseudo-diploid cells as in the Navin *et al.* HC trial (45). Spectral clustering performs poorly by mixing tumor and normal cells. Regarding auto-k clustering algorithms, agglomerative hierarchy identifies five concordant clusters as predefined-k mode. Leiden and Louvain fail at this task. Then, we leverage CNV density signals detected by aCGH from FACS identified D, H, A1, and A2 dissections of T10 (46) as silver-standard to validate the clustering result. We calculate the pairwise Spearman correlation and Euclidean distance (L2-norm) between scaled single-cell

CNV profiles and aCGH CNV signals. As a proof of concept, the three bottom clusters own a higher correlation and a lower distance to aCGH H, A1, and A2 sections, respectively. The cells in the first top cluster detected by SEAT have almost zero correlation and the lowest distance with aCGH D sections, suggesting that they are diploid cells. Pseudo-diploid cells illustrate a low correlation with all aCGH sections, validating their unique CNV profiles. Navin *et al.* sequenced 100 cells from a monogenomic triple-negative breast cancer tumor and its seeded liver metastasis, Navin\_T16 (45). SEAT clusters the 100 samples into four distinct subpopulations. Two are primary and metastasis aneuploid cells, corresponding to the published population structure. Notably, SEAT catalogs diploid cells and pseudo-diploid cells while other tools failed.





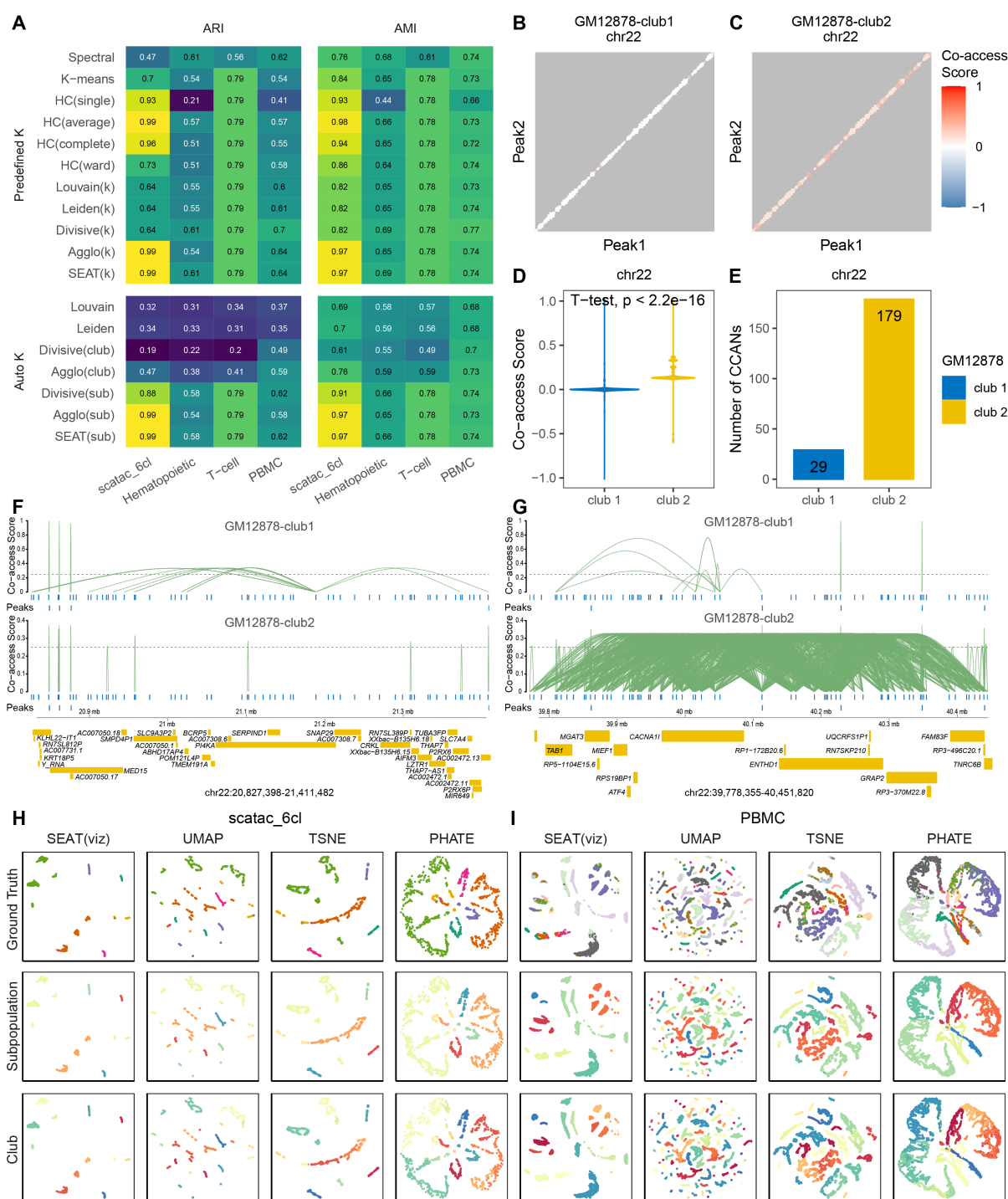
**Fig. 4.** Clustering on scDNA datasets. **A.** The clustering result of Navin\_T10. From left to right is the agglomerative tree yielded by SEAT, clustering results for predefined-k ( $k = 5$ ) and auto-k clustering tools, the whole genome single-cell CNV heatmap of T10, the Spearman correlation, and Euclidean distance (L2-Norm) between scaled copy number profiled by scDNA and copy number density profiled by aCGH, respectively. Spectral: spectral clustering. HC(single), HC(average), HC(complete), and HC(ward): hierarchical clustering with single, average, complete, and ward linkage. Louvain(k) and Leiden(k): Louvain and Leiden in predefined-k mode. Agglo(k): the cell subpopulations from agglomerative hierarchy in predefined-k mode. Agglo(cluster): the cell clubs from agglomerative hierarchy. Agglo(sub): the cell subpopulations from agglomerative hierarchy in auto-k mode. **B.** The stacked area plot illustrates the SEAT subpopulation assignments across 10x\_breast\_S0 tumor sections. Cluster c6 (blue) signifies the diploid cells. **C.** The mean ploidy of SEAT subpopulation assignments across 10x\_breast\_S0 tumor sections. **D.** The whole-genome single-cell CNV heatmap of SEAT subpopulation assignments across 10x\_breast\_S0 tumor sections.

We collect a large-scale 10x scDNA-seq dataset without known subclone labels, 10x\_breast\_S0, where 10,202 cells from five adjacent tumor dissections (A, B, C, D, and E) of triple-negative breast cancer are sequenced. We check whether SEAT seizes the substantial intra-tumor heterogeneity. In Fig. 4B-D, SEAT automatically detects seven subpopulations, and the proportions of the cell subpopulations vary across the five lesions. The blue subpopulation c6 gathers normal cells, with the mean cellular ploidy being diploid for all sections. The number of cells gradually decreases from sections A to E. SEAT identifies six clonal subpopulations (c0-c5), where c3 manifests the highest average ploidy. The distinct amplification events on chr3 and chr4 are mutually exclusive on subclones c0, c1, and c2, indicating an early branching evolution hypothesis consistent with the findings by Wang *et al.*'s (47).

Furthermore, SEAT distinguishes cells with CNV gains and losses in circulating tumor cells in seven patients with lung

cancer (48) and human cortical neurons (49). SEAT also detects the loss of heterogeneity event, validating by successfully classifying chrX-bearing, chrY-bearing, and aneuploid sperm cells (50, 51).

**Cell hierarchy catalogs the accessibility heterogeneity of single-cells.** SEAT catalogs accessibility heterogeneity of single-cells. We utilize three public scATAC-seq data as benchmarking sets with gold-standard cell type labels. scatac\_6cl is a mixture of six cell lines (BJ, GM12878, H1-ESC, HL60, K562, and TF1) (52). Hematopoiesis consists of eight types of human hematopoiesis cells (CLP, CMP, GMP, HSC, LMPP, MEP, MPP, and pDC) (53). T-cell composes of four T-cell subtypes (Jurkat\_T\_cell, Naive\_T\_cell, Memory\_T\_cell, and Th17\_T\_cell) (54). We collect a multiome of scRNA and scATAC dataset, PBMC, for peripheral blood mononuclear cells (PBMCs) with 14 cell types. The order of the cells in agglomerative and divisive hierarchy



**Fig. 5.** Clustering on three scATAC datasets and one scRNA-scATAC multiome dataset. **A.** The adjusted rand index (ARI) and adjusted mutual information (AMI) of predefined-k and auto-k clustering tools. Spectral: spectral clustering. HC(single), HC(average), HC(complete), and HC(ward): hierarchical clustering with single, average, complete, and ward linkage. Louvain(k) and Leiden(k): Louvain and Leiden in predefined-k mode. Divisive(k) and Agglo(k): the cell subpopulations from divisive and agglomerative hierarchy in predefined-k mode. SEAT(k): the cell subpopulations from SEAT cell hierarchy in predefined-k mode. Divisive(club) and Agglo(club): the cell clubs from divisive and agglomerative hierarchy. Divisive(sub) and Agglo(sub): the cell subpopulations from divisive and agglomerative hierarchy in auto-k mode. SEAT(sub): the optimal subpopulations from SEAT cell hierarchy in auto-k mode. **B-D.** The co-accessibility score among peak pairs at chr22 for cells at SEAT club1 and club 2 from scatac\_6cl GM12878 cell type. **E.** The number of cis-co-accessibility networks (CCANs) among pair of peaks at chr22 for cells at SEAT club1 and club 2 from scatac\_6cl GM12878 cell type. **F.** The co-accessibility connections among cis-regulatory elements in chr22:20,827,398-21,441,482. The height of links signifies the degree of the co-accessibility correlation between the pair of peaks. The top panel illustrates cells in scatac\_6cl GM12878-club1, while the bottom shows cells in scatac\_6cl GM12878-club2. **G.** The co-accessibility connections among cis-regulatory elements in chr22:39,778,355-40,451,820. The height of links signifies the degree of the co-accessibility correlation between the pair of peaks. The top panel illustrates cells in scatac\_6cl GM12878-club1, while the bottom shows cells in scatac\_6cl GM12878-club2. **H-I** SEAT hierarchical visualization, UMAP, TSNE, and PHATE plots of scatac\_6cl and PBMC. The cells are colored by subpopulations, clubs, and ground truth. SEAT(viz): the hierarchical visualization from SEAT cell hierarchy.



is consistent with their ground truth cell types. The clustering accuracy of SEAT against its competitors is in Fig. 5A. For the predefined-k mode, SEAT(k) demonstrates the highest clustering accuracy on scatac\_6cl and T-cell sets. For auto-k clustering, SEAT(sub) beats Louvain and Leiden on all four sets. For scatac\_6cl and T-cell, the optimal number of clusters obtained by SEAT matches the ground truth, thus yielding the comparable ARI to predefined-k clustering algorithms. Leiden and Louvain have lower performance due to predicting more clusters than ground truth. We check whether SEAT reveals the functional diversity of single-cell chromatin accessibility. We conduct *cis*-regulatory DNA interaction analysis on chr22 for cells at club1 and club2 predicted by SEAT from the scatac\_6cl GM12878 dataset. Fig. 5B-C depicts the *cis*-regulatory map on chr22 from club1 and club2 cells, respectively. The co-accessibility correlation among peaks from club2 cells is significantly higher ( $p < 0.05$ ) than in club1 cells (Fig. 5D). Meanwhile, we identify 29 and 179 *cis*-co-accessibility networks (CCANs) from GM12878-club1 and GM12878-club2, respectively (Fig. 5E). The genome region where the CCANs are affected shows heterogeneity between GM12878-club1 and GM12878-club2. Fig. 5F shows a GM12878-club1 specified CCANs at chr22:20,827,398-21,441,482. There *cis*-regulatory elements surrounding gene *SNAP29* are co-accessible only in GM12878-club1. Moreover, we found a dense pairwise connection among peaks at chr22:39,778,355-40,451,820 in GM12878-club2 (Fig. 5G), harboring genes *TAB1*, *MGAT3*, *MIEF1*, *CACNA1I*, *ENTHD1*, *GRAP2*, *FAM83F*, *TNRC6B*, etc. Similar to the scRNA visualization refinement experiments, the SEAT hierarchical visualizations reveal a clear pattern of cells corresponding to ground truth, and the nested layouts of subpopulations and clubs are clearly illustrated (Fig. 5H-I). However, UMAP visualizations derived from high-dimensional data mix ground truth cell subpopulations in one clump. Furthermore, UMAP, TSNE, and PHATE visualizations derived from cell-cell graphs fail to place cells from K562 (light green) and TF1 (yellow) within proximity in scatac\_6cl; and they fail to place effector CD8 T cells (magenta) together in PBMC (Fig. 5H-I).

## Discussion

Detecting and visualizing inherent functional diversity is essential in single-cell analysis. Renunciation of the underlying nested structures of cells prevents the capture of full-scale cellular functional diversity. To address this challenge, we incorporate cell hierarchy to investigate the functional diversity of cellular systems at the subpopulation, club, and cell layers, hierarchically. The cell subpopulations and clubs catalog the functional diversity of cells in broad and fine resolution, respectively. In the cell layer, the order of cells further records the slight dynamics among cells locally. Accordingly, we establish SEAT to construct cell hierarchies utilizing structure entropy by diminishing the global uncertainty of cell-cell graphs. In addition, SEAT offers an interface to embed cells into low-dimensional space while preserving the global-

subpopulation-club hierarchical layout in cell hierarchy.

Currently, state-of-the-art clustering tools for cell subpopulation or club investigation renounce the underlying nested structures of cells. Flatten clusterings, such as K-means (11) and spectral clustering (10), do not support the cell hierarchy. Although conventional hierarchical clustering (12), Louvain (13) and Leiden (14) derive cell hierarchy layer by layer via optimizing merging or splitting metrics, computing these metrics merely uses single-layer information. When constructing subsequent layers, they have not incorporated the built-in cell hierarchy in the previous layers. Structure entropy is a metric that encompasses the previously constructed internal cell hierarchy. Experiments validate that SEAT delivers robust cell-type clustering results and forms insightful hierarchical structures of cells.

SEAT is good at finding the optimal subpopulation number with high accuracy. We have collected scRNA, scDNA, and scATAC profiles with the number of cell types ranging from 2 to 14. SEAT consistently predicts the optimal cluster number closest to the gold or silver standards, while Louvain and Leiden predict too many clusters. Especially for scRNA set Kumar, SEAT boosts the accuracy from 0.34 to 1 compared to Louvain and Leiden. Auto-k clustering mode of SEAT is comparable to or better than the best clustering results of predefined-k clustering methods for most datasets.

SEAT specializes in hierarchically deciphering cellular functional diversity at subpopulation and club levels. We observe visible marker gene patterns that match cell clubs within one cell subpopulation. For the p3cl set, the basal, luminal, and fibroblast cell subpopulations have significant cell clubs, determined by the expression of cell cycle genes (*HIST1H4C*, *CDC20*, *CCNB1*, and *PTTG1*). Looking at the seven agglomerative clubs for the basal subpopulation, we find a distinct breast cancer cell club that drives oncogenic *AREG-EGFR* signaling in all basal cells, suggesting a promoting role in tumorigenesis. Cell hierarchy obtained from copy number profiles of 10x\_breast\_S0 demonstrates a mutually exclusive subclones layout, indicating an early branch evolution. Furthermore, we find that there is a cell club specified dense *cis*-regulatory elements co-accessible network at chr22:39,778,355-40,451,820 in GM12878-club2, harboring genes *TAB1*, *MGAT3*, *MIEF1*, *CACNA1I*, *ENTHD1*, *GRAP2*, *FAM83F*, *TNRC6B*, etc.

Inferring the periodic pseudo-time for the cell cycle data is crucial as it reveals the functional diversity of cells undergoing the cell cycle process. Several tools are dedicated to cell cycle pseudo-time inference. CYCLOPS (15) and Cycium (16) utilize deep autoencoders to project expression profiles into cell pseudo-time in the periodic process, which act as black boxes and lack explainability. reCAT (17) employs the Gaussian mixture model to group cells into clusters, and constructs a cluster-cluster graph weighted by the Euclidean distance between the mean expression profile of each cluster, then leverages the traveling salesman path to walk through those clusters with an order. Finding a traveling salesman path is NP-hard, and no polynomial time algorithms are available. CCPE (18) learns a discriminative helix to rep-

resent the periodic process and infer the pseudo-time. However, we fail to run CCPE according to its Github instruction. Moreover, CYCLOPS, Cyclum, reCAT, and CCPE bypass the nested structure of cells when inferring the pseudo-time. In this study, we propose that the cell layer of cell hierarchy encodes the pseudo-time of cells for cell cycle data. We minimize the structure entropy of the kNN cell-cell graph to build the cell hierarchy that carries the nested order between individual cells and their ancestral cell partitions. Then, the order of individual cells is acquired with an in-order traversing of the hierarchy. scRNA data exemplify that SEAT cell orders outperform CYCLOPS, Cyclum, reCAT, and CCPE by accurately predicting the periodic pseudo-time of cells in the cell cycle process. The expressions of M phase marker genes *CDK1*, *TOP2A*, and *CCNB1* rise progressively alongside the SEAT recovered order and are peaked at the M phase with significant fold changes.

Visualizing the hierarchical functional diversity of cells in biological systems is crucial for obtaining insightful biological hypotheses. TSNE (25) preserves the local cell structures. UMAP (24) intends to maintain the global cell structures by minimizing the binary cross entropy. PHATE (26) tackles the general shape and local transition of cells. However, none of them impart the nested structures of cells into the visualization. This study proposes a nonlinear dimension reduction refinement based on UMAP by incorporating a supervised cell hierarchy. We acquire three cell-cell graphs that only store the intra-connections of cells within each global, subpopulation, and club partition. Then, we minimize the weighted binary cross-entropy of the three cell-cell graphs. This approach guarantees the global structure of the cells. Moreover, it ensures that cells within one cell club and cell clubs within one subpopulation are closely placed in the visualization. In contrast, cells from different clubs and subpopulations are kept at a considerable distance. One can adjust the weights of global-subpopulation-cell layers so that the patterns in visualization retain a desired degree of hierarchy. Experiments with scRNA and scATAC data demonstrate that SEAT hierarchical visualization consistently produces a clear layout of cell clumps corresponding to the cell hierarchy.

The structure entropy evaluates the global uncertainty of random walks through a network with a nested structure. The minimum structure entropy interprets a stable nested structure in the network. Li *et al.* has used structure entropy to define tumor subtypes from bulk gene expression data (19) or to detect the hierarchical topologically associating domains from Hi-C data (20). These works utilize greedy merging and combining operations to build a local optimal multi-nary cell hierarchy and cutting hierarchy roughly by keeping the top layers. As we have proven that a binary hierarchy of minimum structure entropy exists for a graph (21), Li *et al.*'s strategy to search for a multi-nary hierarchy is not optimized. Adopted by Louvain and Leiden, modularity is a popular optimization metric to capture community structure in a single-cell network. Agglo(club) is analogous to Louvain's if we switch the merging metric to modularity. Agglo(club) achieves better or comparable performance against Louvain

in most benchmark sets, suggesting the superiority of structure entropy over modularity in measuring the strength of hierarchically partitioning a network into subgroups.

SEAT detects the cell hierarchy, assuming that the entropy codes nested structures of cells. There is no assurance that the resultant cell hierarchy will resemble accurate nested structures of cells. SEAT finds a pseudo cell hierarchy of cells. We show that the pseudo cell hierarchy showcases profound subpopulation detection accuracy and biological insights in single-cell data benchmarking experiments. In future work, we aim to refine the algorithm to find a more accurate and insightful pseudo cell hierarchy.

Recall that the cell hierarchy has multiple layers to present cellular heterogeneity. In this study, we merely utilize four main layers (global, subpopulation, club, and cell) to interpret and visualize functional diversity. In the future, we intend to investigate possible biological insights and visualization layouts derived from more cell hierarchy layers.

Moreover, the order of the cell clubs can be flipped in the cell hierarchy. There is only a partial order among cells bounded by the cell hierarchy. We plan to refine the algorithm to provide a proper non-partial one-dimensional order, which might infer the nuance of pseudo-time or development trajectory among cells outside the periodic cell cycle.

## Data availability

The 25 scRNA, seven scDNA, three scATAC, and one scRNA-scATAC multiome datasets are publicly available.

## Software availability

The source code of SEAT is available at <https://github.com/deepomicslab/SEAT>.

## Competing interests

There is NO competing interest.

## Author contributions statement

LXC conducted the project and wrote the manuscript. SCL supervised the project and revised the manuscript.

## Acknowledgments

Not applicable.

## Funding

This project is supported by CityU/UGC Research Matching Grant Scheme 9229012.

## Online Methods

## Experiment Setting.

**scRNA data.** We collect nineteen scRNA datasets with cell type labels (27–39). For these scRNA datasets, the dimension reduction transformer is UMAP (24). We adopt Seurat (55) for differential expression analysis. Cell-cell communication analysis is conducted with CellChat (41) with default database and parameters. Any ligand-receptor (LR) interaction with less than ten supporting cells is filtered. We also collect six scRNA datasets with gold standard cell cycle labels. Dataset H1-hESC has 247 human embryonic stem cells (hESCs) in G0/G1, S, or G2/M phases identified by fluorescent ubiquitination-based cell cycle indicators (42). The count expression profile and cell cycle labels are obtained with accession code GSE64016. Datasets mESC-Quartz and mESC-SMARTer have 23 and 288 mouse embryonic stem cells (mESCs) sequenced by Quartz-seq and SMARTer, respectively (43, 44). Their G0/G1, S, and G2/M phases are labeled by Hoechst staining. The count expression profile and cell cycle labels are obtained with accession codes GSE42268 and E-MTAB-2805. Datasets 3Line-qPCR\_H9, 3Line-qPCR\_MB, and 3Line-qPCR\_PC3 owns 227 H9 cells, 342 MB cells, and 361 PC3 cells, respectively. The cell cycle stages G0/G1, S, and G2/M are marked by Hoechst staining (28). The raw log2 count expression profiles and cell labels are from the paper’s Data Set S2. The imputation and dimension reduction are conducted by SMURF (56) and UMAP (24). We adopt Seurat (55) for differential expression analysis. Cell-cell communication analysis is conducted with CellChat (41) with default database and parameters. Any ligand-receptor (LR) interaction with less than ten supporting cells is filtered. Gene Ontology (GO) is performed with ShinyGO 0.76 (57).

**scDNA data.** We collect seven scDNA datasets. Navin\_T10 contains 100 cells from a genetically heterogeneous (polygenic) triple-negative breast cancer primary lesion T10, including five cell subpopulations: diploid (D), hypodiploid (H), aneuploid 1 (A1), aneuploid 2 (A2), and pseudo-diploid (P) (45). Navin\_T16 holds 52 cells from genetically homogeneous (monogenetic) breast cancer primary lesion T16P and 48 cells from its liver metastasis T16M, including four cell subpopulations: diploid (D), primary aneuploid (PA), metastasis aneuploid (MA), and pseudo-diploid (P). The Ginkgo CNV profile of T10 and T16 are downloaded from <http://qb.cshl.edu/ginkgo> (58). The silver-standard array comparative genomic hybridization (aCGH) data of T10 and T16 are downloaded with GEO accession code GSE16607 (46). Dataset 10x\_breast\_S0 is a large-scale 10x scDNA-seq set without known cell population labels, where 10,202 cells from five adjacent tumor dissections (A, B, C, D, and E) of triple-negative breast cancer are sequenced. The Bam files are downloaded from 10x official site <https://www.10xgenomics.com/resources/datasets>. We inferred the total CNV profile utilizing Chisel (59). Ni\_CTC sequenced 29 circulating tumor cells (CTCs) across seven lung cancer patients (48). McConnel\_neuron profiles 110 cells from human frontal cortex neurons, with an extensive level of mosaic CNV gains and losses (49). Lu\_sperm

sequenced 99 sperm cells with chrX-bearing, chrY-bearing, and aneuploid groups (50). Wang\_sperm performed single-cell sequencing on 31 sperm cells with CNV gains and losses (51). The Ginkgo CNV profile of these datasets are downloaded from <http://qb.cshl.edu/ginkgo> (58).

**scATAC and scRNA-scATAC multiome data.** We collect three public scATAC-seq data as benchmarking sets with gold standard cell type labels. scatac\_6cl is a mixture of six cell lines (BJ, GM12878, H1-ESC, HL60, K562, and TF1) with 1224 cells (52). Hematopoiesis owns 2210 single-cell chromatin accessibility profiles from eight human hematopoiesis cell populations (CLP, CMP, GMP, HSC, LMPP, MEP, MPP, and pDC) (53). T-cell composes of four T-cell subtypes (Jurkat\_T\_cell, Naive\_T\_cell, Memory\_T\_cell, and Th17\_T\_cell) with a total of 765 cells (54).

We collect a multiome of scRNA and scATAC dataset. PBMC is human peripheral blood mononuclear cells (PBMCs) with 10,032 cells across fourteen cell types.

We downloaded the scOpen (60) processed accessibility profiles and cell labels from <https://github.com/CostaLab/scopen-reproducibility>. UMAP (24) embedded data are used to construct the kNN graphs for each dataset. We adopt Cicero (61) to explore the dynamically accessible element status in different scatac\_6cl GM12878 cell clubs.

**Evaluating cell subpopulation detection.** We access the clustering accuracy of SEAT cell hierarchy, agglomerative hierarchy, and divisive hierarchy with predefined cluster number  $k$ , namely SEAT( $k$ ), Agglo( $k$ ) and Divisive( $k$ ), given by the actual number of ground truth cell types. Competitors are hierarchical clustering (HC) with four linkage strategies (ward, complete, average, and single) (12), K-means (11), and spectral clustering (10). As the leading tool for single-cell clustering, Louvain (13) and Leiden (14) automatically detect how many communities are inside the cell-cell similarity graph. They obtain different numbers of communities at various resolutions. To benchmark Leiden and Louvain in the predefined- $k$  setting, namely Leiden( $k$ ) and Louvain( $k$ ), we heuristically adjusted the resolution 20 times to see if the number of communities was the same as the predefined cluster number  $k$ .

As the predefined  $k$  is undetermined in most real-world scenarios, we evaluate the auto- $k$  clustering efficacy of SEAT against Leiden and Louvain. We also assess the clustering obtained from agglomerative divisive hierarchy clubs, namely Agglo(club) and Divisive(club).

Adjusted Rand index (ARI) (62) and adjusted mutual information (AMI) (63) are adopted as clustering accuracy. They measure the concordance between clustering results and ground truth cell types. A perfect clustering has a value of 1, while random clustering has a value less than or near 0.

**Evaluating cell cycle pseudo-time inference.** SEAT cell hierarchy generates cell order representing the cell cycle pseudo-time for scRNA data. We access the pseudo-time inference accuracy of SEAT given by the actual order of ground truth



cell cycle phases. Benchmark methods are hierarchical clustering (HC) with four linkage strategies (ward, complete, average, and single) (12). An in-order traversal of these hierarchies also generates cell orders. Furthermore, We benchmark our methods with four state-of-the-art tools predicting the cell cycle pseudo-time, CYCLOPS (15), Cyclum (16), reCAT (17), and CCPE (18). CCPE fails the tasks when we follow its Github instruction, so we exclude CCPE for final comparison.

The change index (CI) is used to quantitatively assess the accuracy of cell pseudo-time order against known cell cycle phase labels (17). An ideal cell order changes label  $k-1$  times, where  $k=3$  is the ground truth cell cycle phase number. The change index is defined as  $1 - \frac{c-(k-1)}{n-k}$ , where  $c$  counts the frequency of label alters between two adjacent cells, and  $n$  is the number of cells. A value of 0 suggests the cell order is completely wrong with  $c=n-1$ , while 1 indicates a complete match between cell order and ground truth cell cycle phase with  $c=k-1$ .

**Evaluating hierarchical visualization.** We evaluate the efficacy of SEAT hierarchical visualization with state-of-the-art visualization tools UMAP (24), TSNE (25), and PHATE (26). The dense cell-cell similarity graph  $G$  is used as input, UMAP, TSNE, and PHATE are run with default parameters.

## Reference.

1. Tallulah S Andrews, Vladimir Yu Kiselev, Davis McCarthy, and Martin Hemberg. Tutorial: guidelines for the computational analysis of single-cell rna sequencing data. *Nature protocols*, 16(1):1–9, 2021.
2. Richa Nayak and Yasha Hasija. A hitchhiker’s guide to single-cell transcriptomics and data analysis pipelines. *Genomics*, 2021.
3. Zhijin Wu and Hao Wu. Accounting for cell type hierarchy in evaluating single cell rna-seq clustering. *Genome biology*, 21:1–14, 2020.
4. Yue Gao, Lingxi Chen, Guangyao Cai, Xiaoming Xiong, Yuan Wu, Ding Ma, Shuai Cheng Li, and Qinglei Gao. Heterogeneity of immune microenvironment in ovarian cancer and its clinical significance: a retrospective study. *Oncoimmunology*, 9(1):1760067, 2020.
5. Darlan C Minussi, Michael D Nicholson, Hanghui Ye, Alexander Davis, Kaile Wang, Toby Baker, Maxime Tarabichi, Emi Sei, Haowei Du, Mashiya Rabbani, et al. Breast tumours maintain a reservoir of subclonal diversity during expansion. *Nature*, 592(7853):302–308, 2021.
6. Lingxi Chen, Yuhao Qing, Ruikang Li, Chaohui Li, Hechen Li, Xikang Feng, and Shuai Cheng Li. Somatic variant analysis suite: copy number variation clonal visualization online platform for large-scale single-cell genomics. *Briefings in Bioinformatics*, 23(1):bbab452, 2022.
7. Monika S Kowalczyk, Itay Tirosh, Dirk Heckl, Tata Nageswara Rao, Atray Dixit, Brian J Haas, Rebekka K Schneider, Amy J Wagers, Benjamin L Ebert, and Aviv Regev. Single-cell rna-seq reveals changes in cell cycle and differentiation programs upon aging of hematopoietic stem cells. *Genome research*, 25(12):1860–1872, 2015.
8. Grace XY Zheng, Jessica M Terry, Phillip Belgrader, Paul Ryvkin, Zachary W Bent, Ryan Wilson, Solongo B Ziraldo, Tobias D Wheeler, Geoff P McDermott, Junjie Zhu, et al. Massively parallel digital transcriptional profiling of single cells. *Nature communications*, 8(1):1–12, 2017.
9. Darren A Cusanovich, Riza Daza, Andrew Adey, Hannah A Pliner, Lena Christiansen, Kevin L Gunderson, Frank J Steemers, Cole Trapnell, and Jay Shendure. Multiplex single-cell profiling of chromatin accessibility by combinatorial cellular indexing. *Science*, 348(6237):910–914, 2015.
10. Andrew Y Ng, Michael I Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in neural information processing systems*, pages 849–856, 2002.
11. John A Hartigan and Manchek A Wong. Algorithm as 136: A k-means clustering algorithm. *Journal of the royal statistical society. series c (applied statistics)*, 28(1):100–108, 1979.
12. Stephen C Johnson. Hierarchical clustering schemes. *Psychometrika*, 32(3):241–254, 1967.
13. Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment*, 2008(10):P10008, 2008.
14. Vincent A Traag, Ludo Waltman, and Nees Jan Van Eck. From louvain to leiden: guaranteeing well-connected communities. *Scientific reports*, 9(1):1–12, 2019.

15. Ron C Anafi, Lauren J Francey, John B Hogenesch, and Junhyong Kim. Cyclops reveals human transcriptional rhythms in health and disease. *Proceedings of the National Academy of Sciences*, 114(20):5312–5317, 2017.
16. Shaoheng Liang, Fang Wang, Jincheng Han, and Ken Chen. Latent periodic process inference from single-cell rna-seq data. *Nature communications*, 11(1):1–8, 2020.
17. Zehua Liu, Huazhe Lou, Kaikun Xie, Hao Wang, Ning Chen, Oscar M Aparicio, Michael Q Zhang, Rui Jiang, and Ting Chen. Reconstructing cell cycle pseudo time-series via single-cell transcriptome data. *Nature communications*, 8(1):1–9, 2017.
18. Jiajia Liu, Mengyuan Yang, Weiling Zhao, and Xiaobo Zhou. Ccpe: cell cycle pseudotime estimation for single cell rna-seq data. *Nucleic acids research*, 50(2):704–716, 2022.
19. Angsheng Li, Xianchen Yin, and Yicheng Pan. Three-dimensional gene map of cancer cell types: Structural entropy minimisation principle for defining tumour subtypes. *Scientific reports*, 6(1):1–26, 2016.
20. Angsheng Li, Xianchen Yin, Bingxiang Xu, Danyang Wang, Jimin Han, Yi Wei, Yun Deng, Ying Xiong, and Zhihua Zhang. Decoding topologically associating domains with ultra-low resolution hi-c data by graph structural entropy. *Nature communications*, 9(1):1–12, 2018.
21. Yu Wei Zhang, Meng Bo Wang, and Shuai Cheng Li. Supertad: robust detection of hierarchical topologically associated domains with optimized structural information. *Genome biology*, 22(1):1–20, 2021.
22. Yu Wei Zhang, Lingxi Chen, and Shuai Cheng Li. Detecting tad-like domains from rna-associated interactions. *Nucleic Acids Research*, 2022.
23. Lingxi Chen, Jiao Xu, and Shuai Cheng Li. Deepmf: Deciphering the latent patterns in omics profiles with a deep learning method. *BMC bioinformatics*, 20(23):1–13, 2019.
24. Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
25. Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
26. Kevin R Moon, David van Dijk, Zheng Wang, Scott Gigante, Daniel B Burkhardt, William S Chen, Kristina Yim, Antonia van den Elzen, Matthew J Hirn, Ronald R Coifman, et al. Visualizing structure and transitions in high-dimensional biological data. *Nature biotechnology*, 37(12):1482–1492, 2019.
27. Meichen Dong, Aatish Thennavan, Eugene Urrutia, Yun Li, Charles M Perou, Fei Zou, and Yuchao Jiang. Scdc: bulk gene expression deconvolution by multiple single-cell rna sequencing references. *Briefings in bioinformatics*, 22(1):416–427, 2021.
28. Andrew McDavid, Lucas Dennis, Patrick Danaher, Greg Finak, Michael Krouse, Alice Wang, Philippa Webster, Joseph Beecham, and Raphael Gottardo. Modeling bimodality improves characterization of cell cycle on gene expression in single cells. *PLoS computational biology*, 10(7):e1003696, 2014.
29. Luyi Tian, Xueyi Dong, Saskia Freytag, Kim-Anh Lê Cao, Shian Su, Abolfazl Jalal-Abadi, Daniela Amann-Zalcenstein, Tom S Weber, Azadeh Seidi, Jafar S Jabbari, et al. Benchmarking single cell rna-sequencing analysis pipelines using mixture control experiments. *Nature methods*, 16(6):479–487, 2019.
30. Liying Yan, Mingyu Yang, Hongshan Guo, Lu Yang, Jun Wu, Rong Li, Ping Liu, Ying Lian, Xiaoying Zheng, Jie Yan, et al. Single-cell rna-seq profiling of human preimplantation embryos and embryonic stem cells. *Nature structural & molecular biology*, 20(9):1131–1139, 2013.
31. Qiaolin Deng, Daniel Ramsköld, Björn Reinis, and Rickard Sandberg. Single-cell rna-seq reveals dynamic, random monoallelic gene expression in mammalian cells. *Science*, 343(6167):193–196, 2014.
32. Fernando H Biase, Xiaoyi Cao, and Sheng Zhong. Cell fate inclination within 2-cell and 4-cell mouse embryos revealed by single-cell rna sequencing. *Genome research*, 24(11):1787–1796, 2014.
33. Mubeen Goolam, Antonio Scialdone, Sarah JL Graham, Iain C Macaulay, Agnieszka Jedrusik, Anna Hupalowska, Thierry Voet, John C Marioni, and Magdalena Zernicka-Goetz. Heterogeneity in oct4 and sox2 targets biases cell fate in 4-cell mouse embryos. *Cell*, 165(1):61–74, 2016.
34. Pang Wei Koh, Rahul Sinha, Amira A Barkal, Rachel M Morganti, Angela Chen, Irving L Weissman, Lay Teng Ang, Anshul Kundaje, and Kyle M Loh. An atlas of transcriptional, chromatin accessibility, and surface marker changes in human mesoderm development. *Scientific data*, 3(1):1–15, 2016.
35. Roshan M Kumar, Patrick Cahan, Alex K Shalek, Rahul Satija, A Jay DaleyKeyser, Hu Li, Jin Zhang, Keith Pardee, David Gennert, John J Trombetta, et al. Deconstructing transcriptional heterogeneity in pluripotent stem cells. *Nature*, 516(7529):56–61, 2014.
36. Cole Trapnell, Davide Cacchiarelli, Jonna Grimsby, Prapti Pokharel, Shuang Li, Michael Morse, Niall J Lennon, Kenneth J Livak, Tarjei S Mikkelsen, and John L Rinn. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nature biotechnology*, 32(4):381–386, 2014.
37. Paul Blakeley, Norah ME Fogarty, Ignacio Del Valle, Sissy E Wamaita, Tim Xiaoming Hu, Kay Elder, Philip Snell, Leila Christie, Paul Robson, and Kathy K Niakan. Defining the three cell lineages of the human blastocyst by single-cell rna-seq. *Development*, 142(18):3151–3165, 2015.
38. Aleksandra A Kolodziejczyk, Jong Kyoung Kim, Jason CH Tsang, Tomislav Ilicic, Johan Henriksson, Kedar N Natarajan, Alex C Tuck, Xuefei Gao, Marc Bühler, Pentao Liu, et al. Single cell rna-sequencing of pluripotent states unlocks modular transcriptional variation. *Cell stem cell*, 17(4):471–485, 2015.
39. Yurong Xin, Jinrang Kim, Haruka Okamoto, Min Ni, Yi Wei, Christina Adler, Andrew J Murphy, George D Yancopoulos, Calvin Lin, and Jesper Gromada. Rna sequencing of single human islet cells reveals type 2 diabetes genes. *Cell metabolism*, 24(4):608–615, 2016.
40. Christina S Kappler, Stephen T Guest, Jonathan C Irish, Elizabeth Garrett-Mayer, Zachary Kratche, Robert C Wilson, and Stephen P Ethier. Oncogenic signaling in amphiregulin and egfr-expressing pten-null human breast cancer. *Molecular oncology*, 9(2):527–543, 2015.



41. Suoqin Jin, Christian F Guerrero-Juarez, Lihua Zhang, Ivan Chang, Raul Ramos, Chen-Hsiang Kuan, Peggy Myung, Maksim V Plikus, and Qing Nie. Inference and analysis of cell-cell communication using cellchat. *Nature communications*, 12(1):1–20, 2021.
42. Ning Leng, Li-Fang Chu, Chris Barry, Yuan Li, Jee Choi, Xiaomao Li, Peng Jiang, Ron M Stewart, James A Thomson, and Christina Kendzierski. Osco identifies oscillatory genes in unsynchronized single-cell rna-seq experiments. *Nature methods*, 12(10):947, 2015.
43. Yohei Sasagawa, Itoshi Nikaido, Tetsutaro Hayashi, Hiroki Danno, Kenichiro D Uno, Takeshi Imai, and Hiroki R Ueda. Quartz-seq: a highly reproducible and sensitive single-cell rna sequencing method, reveals non-genetic gene-expression heterogeneity. *Genome biology*, 14(4):1–17, 2013.
44. Florian Buettner, Kedar N Natarajan, F Paolo Casale, Valentina Proserpio, Antonio Scialdone, Fabian J Theis, Sarah A Teichmann, John C Marioni, and Oliver Stegle. Computational analysis of cell-to-cell heterogeneity in single-cell rna-sequencing data reveals hidden subpopulations of cells. *Nature biotechnology*, 33(2):155–160, 2015.
45. Nicholas Navin, Jude Kendall, Jennifer Troge, Peter Andrews, Linda Rodgers, Jeanne McIndoo, Kerry Cook, Asya Stepansky, Dan Levy, Diane Esposito, et al. Tumour evolution inferred by single-cell sequencing. *Nature*, 472(7341):90, 2011.
46. Nicholas Navin, Alexander Krasnitz, Linda Rodgers, Kerry Cook, Jennifer Meth, Jude Kendall, Michael Riggs, Yvonne Eberling, Jennifer Troge, Vladimir Grubor, et al. Inferring tumor progression from genomic heterogeneity. *Genome research*, 20(1):68–80, 2010.
47. Rujin Wang, Dan-Yu Lin, and Yuchao Jiang. Scope: a normalization and copy-number estimation method for single-cell dna sequencing. *Cell systems*, 10(5):445–452, 2020.
48. Xiaohui Ni, Minglei Zhuo, Zhe Su, Jianchun Duan, Yan Gao, Zhijie Wang, Chenghang Zong, Hua Bai, Alec R Chapman, Jun Zhao, et al. Reproducible copy number variation patterns among single circulating tumor cells of lung cancer patients. *Proceedings of the National Academy of Sciences*, 110(52):21083–21088, 2013.
49. Michael J McConnell, Michael R Lindberg, Kristen J Brennand, Julia C Piper, Thierry Voet, Chris Cowing-Zitron, Svetlana Shumilina, Roger S Lasken, Joris R Vermeesch, Ira M Hall, et al. Mosaic copy number variation in human neurons. *Science*, 342(6158):632–637, 2013.
50. Sijia Lu, Chenghang Zong, Wei Fan, Mingyu Yang, Jinsen Li, Alec R Chapman, Ping Zhu, Xuesong Hu, Liya Xu, Liying Yan, et al. Probing meiotic recombination and aneuploidy of single sperm cells by whole-genome sequencing. *science*, 338(6114):1627–1630, 2012.
51. Jianbin Wang, H Christina Fan, Barry Behr, and Stephen R Quake. Genome-wide single-cell analysis of recombination activity and de novo mutation rates in human sperm. *Cell*, 150(2):402–412, 2012.
52. Jason D Buenrostro, Beijing Wu, Ulrike M Litzenburger, Dave Ruff, Michael L Gonzales, Michael P Snyder, Howard Y Chang, and William J Greenleaf. Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature*, 523(7561):486–490, 2015.
53. Jason D Buenrostro, M Ryan Corces, Caleb A Lareau, Beijing Wu, Alicia N Schep, Martin J Aryee, Ravindra Majeti, Howard Y Chang, and William J Greenleaf. Integrated single-cell analysis maps the continuous regulatory landscape of human hematopoietic differentiation. *Cell*, 173(6):1535–1548, 2018.
54. Ansuman T Satpathy, Naresha Saligrama, Jason D Buenrostro, Yuning Wei, Beijing Wu, Adam J Rubin, Jeffrey M Granja, Caleb A Lareau, Rui Li, Yanyan Qi, et al. Transcript-indexed atac-seq for precision immune profiling. *Nature medicine*, 24(5):580–590, 2018.
55. Yuhao Hao, Stephanie Hao, Erica Andersen-Nissen, William M Mauck III, Shiwei Zheng, Andrew Butler, Maddie J Lee, Aaron J Wilk, Charlotte Darby, Michael Zager, et al. Integrated analysis of multimodal single-cell data. *Cell*, 2021.
56. Bingchen Wang, Juhua Pu, Lingxi Chen, and Shuaicheng Li. Smurf: embedding single-cell rna-seq data with matrix factorization preserving selfconsistency. *bioRxiv*, 2022.
57. Steven Xijin Ge, Dongmin Jung, and Runan Yao. Shinygo: a graphical gene-set enrichment tool for animals and plants. *Bioinformatics*, 36(8):2628–2629, 2020.
58. Tyler Garvin, Robert Aboukhalil, Jude Kendall, Timour Baslan, Gurinder S Atwal, James Hicks, Michael Wigler, and Michael C Schatz. Interactive analysis and assessment of single-cell copy-number variations. *Nature methods*, 12(11):1058, 2015.
59. Simone Zaccaria and Benjamin J Raphael. Characterizing allele-and haplotype-specific copy numbers in single cells with chisel. *Nature biotechnology*, 39(2):207–214, 2021.
60. Zhijian Li, Christoph Kuppe, Susanne Ziegler, Mingbo Cheng, Nazanin Kabgani, Sylvia Menzel, Martin Zenke, Rafael Kramann, and Ivan G Costa. Chromatin-accessibility estimation from single-cell atac-seq data with scopen. *Nature communications*, 12(1):1–14, 2021.
61. Hannah A Pliner, Jonathan S Packer, José L McFaline-Figueroa, Darren A Cusanovich, Riza M Daza, Delasa Aghamirzaie, Sanjay Srivatsan, Xiaojie Qiu, Dana Jackson, Anna Minkina, et al. Cicero predicts cis-regulatory dna interactions from single-cell chromatin accessibility data. *Molecular cell*, 71(5):858–871, 2018.
62. William M Rand. Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical association*, 66(336):846–850, 1971.
63. Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons, 2012.