

METHOD

Open Access



# CellWalker integrates single-cell and bulk data to resolve regulatory elements across cell types in complex tissues

Pawel F. Przytycki<sup>1</sup> and Katherine S. Pollard<sup>1,2,3\*</sup> 

\* Correspondence: [katherine.pollard@gladstone.ucsf.edu](mailto:katherine.pollard@gladstone.ucsf.edu)

<sup>1</sup>Gladstone Institutes, San Francisco, CA, USA

<sup>2</sup>Chan-Zuckerberg Biohub, San Francisco, CA, USA

Full list of author information is available at the end of the article

## Abstract

Single-cell and bulk genomics assays have complementary strengths and weaknesses, and alone neither strategy can fully capture regulatory elements across the diversity of cells in complex tissues. We present CellWalker, a method that integrates single-cell open chromatin (scATAC-seq) data with gene expression (RNA-seq) and other data types using a network model that simultaneously improves cell labeling in noisy scATAC-seq and annotates cell type-specific regulatory elements in bulk data. We demonstrate CellWalker's robustness to sparse annotations and noise using simulations and combined RNA-seq and ATAC-seq in individual cells. We then apply CellWalker to the developing brain. We identify cells transitioning between transcriptional states, resolve regulatory elements to cell types, and observe that autism and other neurological traits can be mapped to specific cell types through their regulatory elements.

## Background

Gene regulatory elements are critical determinants of tissue and cell type-specific gene expression [1, 2]. Annotation of putative enhancers, promoters, and insulators has rapidly improved through large-scale projects such as ENCODE [3], PsychENCODE [4], B2B [5], and Roadmap Epigenomics [6]. However, both predictions and validations of regulatory elements have been made largely in cell lines or bulk tissues lacking anatomical and cellular specificity [7]. Bulk measurements miss regulatory elements specific to one cell type, especially minority ones [8]. This lack of specificity limits our ability to determine how genes are differentially regulated across cell types and to discover the molecular and cellular mechanisms through which regulatory variants affect phenotypes.

Single-cell genomics is an exciting avenue to overcoming limitations of bulk tissue studies [8, 9]. However, these technologies struggle with low-resolution measurements featuring high rates of dropout and few reads per cell [8, 9]. Many methods have been developed to address these problems in single-cell expression data (scRNA-seq) [8, 9]. However, these strategies generally fail on scATAC-seq data because there are fewer reads per cell, and the portion of the genome being sequenced is typically much larger



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

than the transcriptome [10]. Consequently, scATAC-seq has much lower coverage and worse signal-to-noise than scRNA-seq.

Several scATAC-seq analysis methods have been developed to increase the number of informative reads used per cell. These include Cicero [11], which aggregates reads from peaks that are co-accessible with gene promoters to emulate gene-focused scRNA-seq data, and SnapATAC [12], which computes cell similarity based on genome-wide binning of reads. Other methods search for informative reads based on known or predicted regulatory regions [13, 14]. However, these approaches often miss rare but known cell types [10]. Other methods attempt to detect cell types in scATAC-seq data by either mapping the data into the same low-dimensional space as scRNA-seq data or by labeling cells in scATAC-seq to known cell-type expression profiles [15, 16]. While these provide a promising avenue towards adding labels to clusters of cells observed in scATAC-seq data, they do not help to increase the resolution of cell type detection.

We present CellWalker, a generalizable network model that improves the resolution of cell populations in scATAC-seq data, determines cell label similarity, and generates cell type-specific labels for bulk data by integrating information from scRNA-seq and a variety of bulk data. These labels can be generated concurrently from the same tissue, but could also be from cell lines, sorted cells, or related tissues. Our method goes beyond co-embedding or directly labeling cells with this prior knowledge about cell types, instead propagating cell-type signatures over a network of cells and cell types so that they are weighted with evidence of cell types in scATAC-seq. Diffusion through this network allows labeling information to indirectly influence cells with similar genome-wide open chromatin profiles even if they could not be initially labeled. A major benefit of our model is that it allows us to compute the level of influence of each label and cell on every other label and cell, thus providing an avenue for additional inferences. These include deconvoluting bulk measurements and assessing their relevance to specific cell populations, as well as quantifying similarity between known cell types in the tissue where scATAC-seq was performed.

The developing human brain presents a complex landscape of cell types, each with unique regulatory programs [17–19]. Using CellWalker, we mapped cell types derived from scRNA-seq data to a large set of scATAC-seq data. The derived influence matrix made it possible to examine changes in regulation across neuronal development and map enhancers to specific cell types. Using this cell type-specific atlas of putative regulatory elements (pREs), we found that autism spectrum disorder (ASD) genes are enriched for pREs specifically active in inhibitory interneurons, while developmental delay genes are enriched for pREs specifically active in radial glia. The ability to map psychiatric traits to cell types is a crucial step towards understanding the mechanisms through which disease develops and responds to treatment. As more large-scale single-cell studies are released, generalizable methods such as CellWalker will be fundamental towards integrating them with existing bulk data to increase our understanding of cell type-specific regulatory programs.

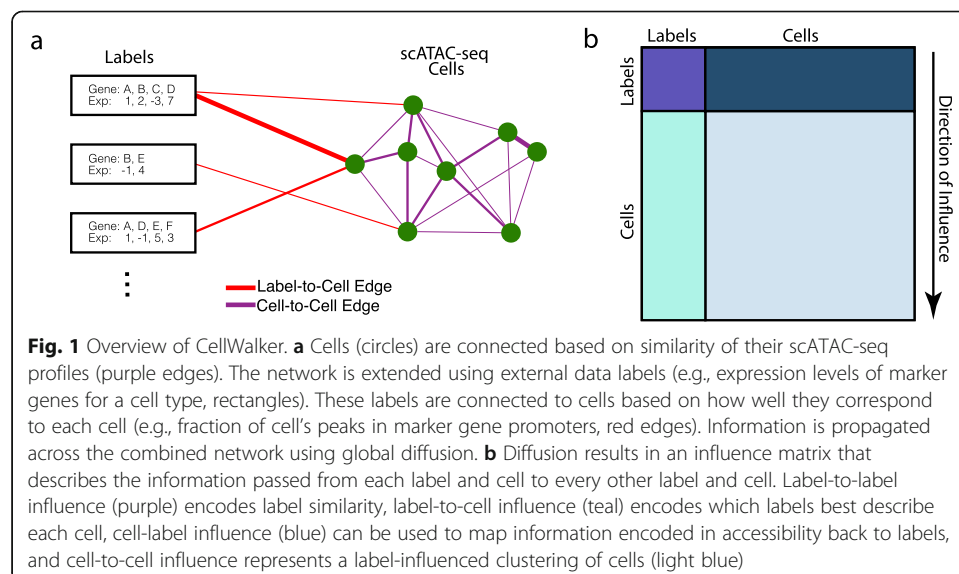
## Results

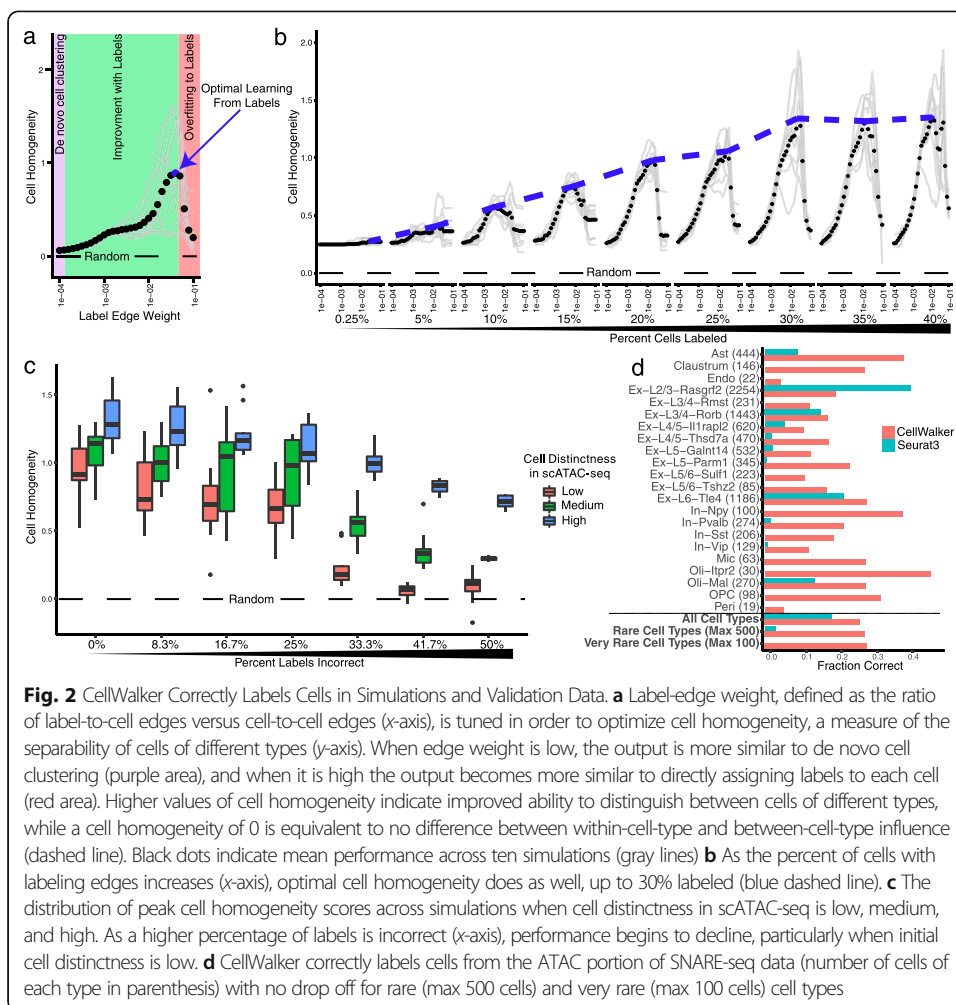
### Overview of method

CellWalker resolves cell types and differentially accessible regions in scATAC-seq data by integrating information from scRNA-seq and bulk data (see Additional File 1: Fig.

S1 for full pipeline). This integration relies on building a combined network featuring nodes representing cells in scATAC data and nodes for external labeling data, e.g., cell types derived from scRNA-seq data (Fig. 1a). Briefly, cells from scATAC-seq are nodes in the network, and edges between them encode information about cell similarity. A second set of nodes represents labeling datasets connected to cell nodes by edges that encode the similarity between a label and a cell.

Using a graph diffusion implemented via a random walk with restarts, CellWalker computes a global influence matrix that relates every cell and label to every other cell and label based on information flow between them in the network (Fig. 1b). In this matrix, each column describes where walks that begin at a given node end. Different portions of this matrix can be used to map information between and within domains: cell-to-cell for clustering cells, label-to-label for exploring label similarity, label-to-cell for cell type labeling, and cell-to-label for distributing bulk signatures to labels. The user can set a single label edge weight parameter ( $s$ ) defining the ratio of label-to-cell edges versus cell-to-cell edges (Fig. 2a). This parameter represents a trade-off between cell similarity information in the scATAC-seq data versus in the external labels. When  $s$  is low, the output is similar to de novo cell clustering using only scATAC-seq, and when it is high, the output converges towards directly assigning labels to each cell. Thus,  $s$  can be chosen to reflect user preference across these strategies, or it can be tuned as CellWalker is run to optimize a criterion that assesses the quality of the resulting cell clusters. We developed a measure called cell homogeneity for this purpose. It can be computed directly from the influence matrix as the median ratio of information between cells within the same cell type to information between cells of different cell types. A higher cell homogeneity score indicates a greater ability to differentiate between different cell types.





**Method validation and evaluation**

To assess the ability of CellWalker to distribute labeling information across cells, we first tested it on a compendium of simulated datasets. Using cell homogeneity to quantify performance, we found that as few as 10% of cells being labeled is sufficient for CellWalker to improve cell labeling, and that there is no further improvement after ~ 30% of cells are labeled (Fig. 2b). As expected, cell homogeneity degrades as more cells are initially mislabeled, and improves when cells of different types are more distinct from each other in the scATAC-seq data (Fig. 2c and Additional File 1: Fig. S2a). Furthermore, CellWalker performs well with noisy data, even when up to 50% of reads are dropped or random reads are added (Additional File 1: Fig. S2b). Finally, we observed that CellWalker is able to distribute labels to novel cell populations (Additional File 1: Fig. S2c). These results establish the network diffusion strategy implemented in CellWalker as a robust approach to integrate scATAC-seq with scRNA-seq or other labeling data.

In all cases, we observed that tuning the label edge weight parameter was very important for assigning accurate labels to cells. In particular, a setting of *s* near 0, where labels are assigned to de novo clusters, was never the optimal setting of *s*. Similarly,

very high settings of  $s$ , where labels are directly assigned to cells without considering cell-to-cell similarity, were also not among the most accurate. This indicates that blending these strategies is beneficial for accurately labeling cells in scATAC-seq data.

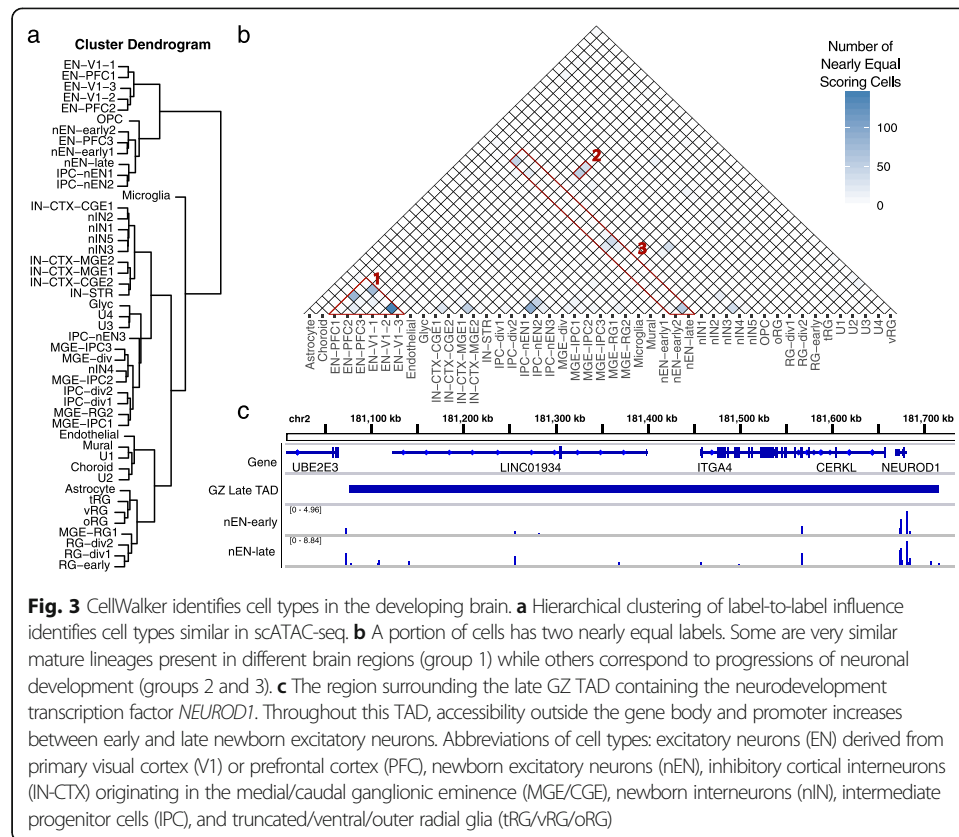
Importantly, CellWalker is efficient enough in terms of both compute time and memory usage to be practical for analysis of current single-cell data set sizes. Running the method beginning to end on 10,000 cells requires only 8 min on a single core on a personal computer (Additional File 1: Fig. S3). On a high-performance cluster, we estimate that 100,000 cells could be analyzed in ~ 80 h of total clock time.

Next, we tested CellWalker on adult mouse cortex SNARE-seq data which includes both scRNA-seq and scATAC-seq reads for each cell [20]. We analyzed the scATAC-seq portion of the data with CellWalker. For cell type labeling, we integrated the scATAC-seq data with differentially expressed marker genes previously derived from clustering the scRNA-seq portion of the SNARE-seq data. Performance was evaluated using the held-out scRNA-seq label for each cell that was identified in the original publication. We tuned the edge weight parameter  $s$  to optimize cell homogeneity (as in our simulations) and observed that this closely mirrors optimization of the fraction of cells labeled correctly, validating cell homogeneity as a measure of how well cell types are resolved (Additional File 1: Fig. 4a and b). We compared CellWalker to label transfer, as implemented in Seurat [15], and found that CellWalker labels more cells correctly (Fig. 2d). This advantage is greater when considering only rare cell types and very rare cell types (Fig. 2d, bottom). We saw similar results on a second set of SNARE-seq data derived from developing mouse cortex, as well as for a 10x Single Cell Multiome ATAC + Gene Exp chip of healthy human brain tissue (Additional File 1: Fig. 4c and d) [20, 21]. Taken together, these analyses of multiome data indicate that CellWalker's integration of label data provides a substantial advantage towards resolving cell types in scATAC-seq data, particularly for analysis datasets with poorer data quality and for identifying rarer cell types.

### Identification of cell types in the developing brain

Given the ability of CellWalker to identify rare cell types in brain SNARE-seq data, we next applied it to a scATAC-seq study of the human telencephalon with multiple biological replicates spanning mid-gestation [19]. Previous work generated a cell type atlas in similar samples based on extensive analysis of scRNA-seq data [17]. Using this atlas as external labeling data, we used CellWalker to compute a full influence matrix across all labels and 30,000 scATAC-seq cells. First, using the label-to-label portion of the influence matrix, we hierarchically clustered all labels and observed high agreement with the scRNA-seq clustering from the previously published results from Nowakowski et al. [17] (Fig. 3a). In terms of broad cell types, we observed that all radial glia and all interneurons group together. Other more local similarities were reflected as well, such as early newborn excitatory neurons (nEN-early1) being more similar to prefrontal cortex excitatory neurons (EN-PFC3) than to other newborn excitatory neuron types.

Next, we scored each cell using label-to-cell influence. This produces a “fuzzy” labeling of cells, representing the fact that a scATAC-seq cell may be strongly connected through the network to multiple cell types. For most cells in the scATAC-seq data (27,



270 out of 30,000), CellWalker assigned a label without much ambiguity. Thus, most transcriptional states observed in scRNA-seq are associated with a distinct open chromatin signature in scATAC-seq. While previous work assigning cell types based on clustering of this scATAC-seq data only allowed for identification of broad cell types found in scRNA-seq (see Ziffra et al. [19]), CellWalker was able to identify cells of all specific types, notably distinguishing between different subtypes of radial glia and separating the two major types of medial ganglionic eminence (MGE)-derived inhibitory interneurons.

In a few cases, we observed cells with multiple nearly equally scoring labels, indicating intermediate membership in multiple cell types and revealing transcriptional states that correspond to highly similar open chromatin profiles (Fig. 3b). Some of these relationships, such as visual cortex (V1) and prefrontal cortex (PFC) excitatory neurons, represent similar types of maturing neurons that are present in two brain regions (Fig. 3b, group 1). Others correspond to progressions of neuronal development. For example, the newborn interneuron and caudal ganglionic eminence (CGE) cortical interneuron cell types have shared network influence on a large group of scATAC-seq cells (Fig. 3b, group 2). Some transcriptional states are difficult to distinguish from each other, because cells receive influence from multiple different scRNA-seq cell types. Specifically, we found groups of cells that scored highly as two or more of the following types: intermediate progenitor cells, early newborn excitatory neurons, late newborn excitatory neurons, and maturing excitatory neurons (Fig. 3b, group 3).

To explore whether these indeterminate cell types represent limitations of scATAC-seq data, failures of the CellWalker model, or cases where transcription changes without large changes in open chromatin, we took a closer look at early and late newborn excitatory neurons (nEN-early and nEN-late respectively). These are fairly abundant, identifiable cell types in scRNA-seq [17]. However, 92% of nEN-early ATAC peaks are also found in nEN-late cells. We assigned each scATAC-seq cell an excitatory neuron progression score based on the difference between the influence of the early and late newborn excitatory neuron labels, such that a higher excitatory neuron progression score indicates a later newborn excitatory neuron. Using this score, we observed that while there is a small distinct set of early newborn excitatory neurons, the majority of newborn excitatory neurons fall evenly between the two types with many scores near zero (Additional File 1: Fig. S5a). This indicates that there is a continuous gradient of changes in chromatin accessibility rather than large-scale difference between transitioning cell types. However, this observed difference between the dynamics of gene expression versus open chromatin during developmental transitions may not hold up with higher coverage scATAC-seq data, which could elucidate distinct chromatin profiles.

#### Cell type-specific annotation of loci

We next sought to determine if cell type annotations could be used to characterize the biology of loci based on chromatin accessibility at distal regulatory elements. Because most distal regulation occurs within Topologically Associated Domains (TADs) [22], we asked if the transition from early to late excitatory neurons could be attributed to differences in TAD accessibility between cell types. It is generally believed that Hi-C contact maps derived from bulk data represent the average of a mixture of cells [23]. We correlated the distal accessibility (defined as outside a gene body or promoter) of TADs derived from the germinal zone (GZ) of the mid-gestation developing human cerebral cortex [22] with excitatory newborn progression score and found that the distribution of correlations is significantly bimodal (empirical  $p$  value = 0.021, Additional File 1: Fig. 5b and 5c). This means that the accessibility of GZ TADs distinctly either correlates or anti-correlates with cell state progression from early to late excitatory neuron. As a control, we find that the median distance of peaks to genes and the number of peaks per TAD do not correlate with excitatory neuron progression (Additional File 1: Fig. 5d and 5e). We therefore classified GZ TADs as early or late depending on their correlation with excitatory neuron progression. As a validation of the classification of these TADs, we find that the expression of genes in early TADs negatively correlates with excitatory neuron progression score, while the expression of genes in late TADs correlates positively (median correlations of  $-0.62$  and  $0.22$  respectively). Thus, subtle changes in chromatin accessibility between early and late newborn excitatory neurons may be associated with cell type-specific TAD activity. The ability to separate TADs by cell type enables a greater understanding of gene regulation in complex tissues such as the human brain. A similar strategy could be applied to other annotations of loci, such as linkage disequilibrium (LD) blocks or expression quantitative trait loci (eQTLs).

Several key genes involved in neuronal development lie in early or late TADs, indicating their expression may be distally regulated. Notably, the neurogenic differentiation gene *NEUROD1* lies in a late TAD with higher levels of accessibility late than early

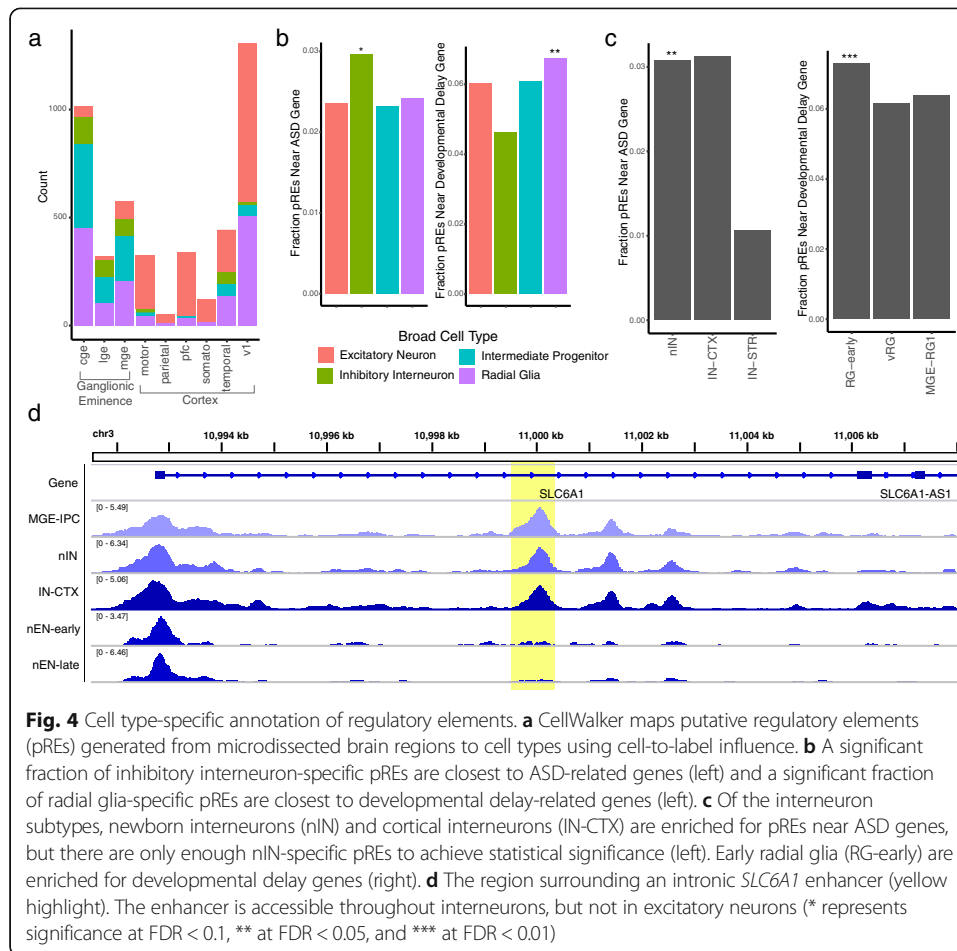
throughout the TAD but similar accessibility in the gene body and promoter (Fig. 3c). Correspondingly, *NEUROD1* has two-fold higher mean transcripts in late than early newborn excitatory neurons (mean 73 TPM early vs 131 late in scRNA-seq data [17]). This indicates that the gene expression differences of *NEUROD1* are potentially driven by distal enhancers. Conversely, *TENM4*, which is involved in establishing neuronal connectivity during development [24], lies in an early TAD (Additional File 1: Fig. S5f) and is less expressed in late newborn excitatory neurons (mean 350 TPM early vs 249 late in scRNA-seq data [17]). Deciphering the cell type-specific regulation of these genes is an important step towards understanding how differences in genotype lead to their misexpression and linked diseases.

### Cell type-specific annotation of regulatory elements

It is generally believed that many enhancers involved in brain development function in a cell type-specific manner [19]. CellWalker provides a way to explore this idea. We mapped pREs derived from bulk ATAC-seq on microdissected tissue across the mid-gestation human telencephalon [18] to cell types based on cell-to-label influence (Fig. 4a). As expected, we found that the many pREs specific to the ganglionic eminence map to intermediate progenitor cell types, while pREs in other regions primarily map to types of excitatory neurons [17]. As further validation, we also observed that pREs from cell types labeled to specific regions such as the MGE map to regions sampled from the ganglionic eminence (Additional File 1: Fig. S6a). These findings demonstrate that cell types resolved in scATAC-seq data with CellWalker can be used to annotate regulatory elements discovered in bulk ATAC-seq. This strategy combines the benefits of high coverage in bulk data with the cell-type information in scATAC-seq.

We next examined whether cell type-specific pREs are associated with disease genes. First, we considered sets of genes near significant variants detected in a collection of Genome-Wide Association Studies (GWAS) [25]. Testing for associations with pREs active in four broad cell types, among these gene sets, we found significant enrichment for pREs near genes associated with a collection of neurological diseases as well as many measures for developmental delay (significant at FDR < 0.1, Additional File 2: Table S1). We therefore decided to take a closer look at curated lists of genes linked to autism spectrum disorders (ASD) and developmental delay [26]. We found that pREs specific to radial glia are significantly associated with genes linked to developmental delay (Fig. 4b, right), among which pREs specific to early radial glia are significant (Fig. 4c, right). pREs specific to inhibitory interneurons are significantly associated with genes linked to ASD, in agreement with previous studies [27–29] (Fig. 4b, left). Among these, enrichment was significant for newborn interneurons (Fig. 4c, left). Recently, an enhancer of *SLC6A1* which is accessible in cells in the ganglionic eminence was linked to ASD [18]. We found that this enhancer maps to both intermediate progenitor cells located in the ganglionic eminence as well as to newborn interneurons and is accessible in cells predicted to belong to these cell types (Fig. 4d). As validation, we found that this peak is accessible in intermediate progenitors and interneurons but not excitatory neurons in ATAC-seq data on FACS sorted cells [30] (Additional File 1: Fig. S6c). Interestingly, while the enhancer is most strongly linked to early stages of interneuron development, expression of *SLC6A1* increases throughout the course of interneuron





development (Additional File 1: Fig. S6b). It is possible therefore that de novo mutations observed in this enhancer in ASD individuals contribute to changes in the initiation of *SLC6A1* expression, which could influence the timing of interneuron development. This strategy for determining cell type-specific effects can be applied to other loci to better understand their potential roles in disease and cell differentiation.

## Discussion

The development of high-throughput sequencing technologies has enabled an explosion of data generation, necessitating techniques to integrate these data into knowledge and testable hypotheses. Using CellWalker, we were able to uncover cell type-specific signals based on a combination of bulk and single-cell data. This was made possible with external data about known cell types which helped overcome the low signal-to-noise ratio present in scATAC-seq data. This strategy is broadly applicable, as there are already vast amounts of bulk RNA-seq, bulk epigenomics, scRNA-seq, and other cell atlas data for tissues, organoids, and cell lines related to samples where scATAC-seq is being performed. Here, we applied CellWalker to neurodevelopment. In another study, we used CellWalker to map transcriptional disease states from mouse heart data to scATAC-seq data in matched tissues and uncovered cell type-specific enhancers activated by stress [31]. CellWalker is computationally efficient enough to enable even

larger scale integrations of such data, with 100,000 cells being feasible to analyze on a high-performance cluster (Additional File 1: Fig. S3).

CellWalker extends naturally to incorporate multiple cell atlases simultaneously. This presents a myriad of possible opportunities such as measuring the influence between disease and non-disease labels for similar cells. For example, an atlas of the developing brain could be used together with cell types derived from post-mortem brains of individuals with ASD to directly measure the relationships between those sources of labeling data. An alternative possibility is to use CellWalker to transfer labels across species. One question the simultaneous use of multiple atlases raises is how label edge weight parameters vary across data, and how these can be interpreted for scoring the relevance of different data sets to labeling. Alternative integrations of data are also possible by, for example, using each label multiple times, but weighing edges differently based on bulk measurements of histone marks. Much of the power of CellWalker lies in its generalizable network model.

The cellular complexity of the developing human brain presented an ideal testbed for CellWalker. We were able to detect rare cell types and tease out distal regulatory programs by integrating bulk data with scATAC-seq data. However, as more single-cell data is generated with higher read depths and greater cell coverage, it may turn out that our data simply did not have the power to uncover the true underlying regulatory landscape. Rare intermediate cell types that correspond better with identified transcriptional states may exist. Furthermore, new approaches that simultaneously measure multiple epigenetic and transcriptional attributes in the same cell will soon begin to enable the detection of cell-specific links between regulation and expression [21]. However, while these technologies continue to be developed and improved, exploiting existing troves of bulk data provides a powerful avenue towards understanding cell type-specific regulation.

## Conclusions

In summary, we have shown that CellWalker is a powerful tool for combining bulk data with single-cell analysis. We were able to define cell types in a complex tissue and annotate genomic loci with the cell types in which they are active. This in turn enabled us to dig into the cell type-specific dynamics of chromatin accessibility in the developing human brain and dissect genetic causes of neurological disease.

## Methods

### CellWalker initialization

CellWalker takes as input scATAC-seq data and labeling information, either directly in the form of marker genes, or by processing scRNA-seq data to generate labels (for example using Seurat). scATAC-seq data can optionally be converted into a cell-by-gene matrix using software such as SnapATAC, Cicero, or ArchR [32]. Overall, initializing CellWalker requires either a cell-by-peak or cell-by-gene matrix and labels with associated marker genes, optionally with the log-fold change in expression for each marker gene in each label.

### CellWalker network construction

The network constructed by CellWalker consists of two types of edges: cell-to-cell and label-to-cell. Cell-to-cell edges have an edge weight corresponding to the similarity

between pairs of cells in scATAC-seq data. This can be computed flexibly using measures such as the Jaccard similarity of cell peak profiles or binned cell accessibility (as in SnapATAC), or as the distance in PCA space (as in Seurat). Cell-to-label edges have an edge weight corresponding to the similarity between the given labeling feature and each cell. If the data is in the form of a cell-by-gene matrix, this is directly used to compute weights as the fraction of each cell's gene score that falls on a label's marker genes, optionally scaled by log-fold change in expression. If, for example, distal peaks should be included in the edge weights, the gene-by-peak matrix outputted by Cicero can be used to compute cell-to-label peaks. Alternatively, if the data is in the form of a peak-by-peak matrix, the fraction of a cell's peaks in the promoters or promoters and gene bodies of marker genes can be used. With this general approach, it is possible to add a large variety of external data to the model. Although these edges may be sparsely connected to cells, the edges between cells distribute information. CellWalker includes a single parameter, label edge weight, which determines the ratio of the weight of label-to-cell edges relative to the weight of cell-to-cell edges.

#### CellWalker diffusion

To diffuse the information from all data sources across the network, we implemented a random walk with restarts. A unit amount of information is initialized at each node. Then, at each time step, a fixed portion restarts and the remainder propagates across each edge connected to the node, proportionally to edge weights. Even cells poorly annotated with external data will receive information about those annotations via cells that are similar. This algorithm is equivalent to an insulated heat diffusion graph kernel. To implement diffusion, we first compute a  $q$ -by- $q$  walk matrix  $W$  encoding the fraction of information that must move to each neighboring node in each time step, where  $q$  is equal to the total number of nodes in the graph. This is 0 if the nodes have no edge between them and the fraction of total weight of edges for each node otherwise. In matrix notation, the computation is  $W = D^{-1}A$ , where  $D$  is a diagonal matrix of the sums of edge weights for each node and  $A$  is the adjacency matrix representing the graph. Given this formulation of the walk matrix  $W$  and a non-zero restart probability  $\alpha$ , the walk always converges to a stationary distribution. Due to this property, there is a closed form solution for the  $q$ -by- $q$  influence matrix  $F$ , which defines the amount of information that reaches each node from each other node and is computed as  $F = \alpha(I - (1 - \alpha)W)^{-1}$ . Prior work has examined how different settings of alpha distribute information to neighboring nodes and found that a restart probability between 0.4 and 0.6 encodes graph structure well with only minor variance in information in that range [33]. Based on this, we set our restart probability to 0.5.

#### CellWalker influence matrix

The output of the diffusion process is an influence matrix  $F$  that is a square matrix with dimensions equal to the number of cells plus the number of labels. In this matrix, each column represents the amount of influence that cell or label has on each other cell and label. Thus, the matrix includes three portions used for downstream analysis: label-to-cell influence, label-to-label influence, and cell-to-label influence.

### Optimizing label edge weight

The label edge weight parameter  $s$  can be internally optimized based on cell homogeneity. Cell homogeneity is computed directly from the influence matrix  $F$  as the log of the median ratio of information between cells within the same label to information between cells with different labels. Cell homogeneity is fast to compute, taking less than a second per label with 10,000 cells, and thus serves as a fast way to set the label edge weight parameter.

### CellWalker applications: cell labeling, label clustering, and bulk data mapping

Label-to-cell influence is used for *cell labeling*. Each cell receives a score from each label allowing for a “fuzzy” allocation of labels. A label can be assigned to each cell based on maximum influence, or multiple labels can be considered for each cell to determine if there is ambiguity in cell labeling. Cell labels can then be used for further downstream analysis for peak calling, transcription factor binding, and other analyses provided by software such as SnapATAC and cisTopic. Label-to-label influence is used for *label clustering*. Via indirect diffusion through cells, there is a score for how much each label influences each other label. This determines how similar labels are to each other, and these influence scores can directly be used as distances for hierarchical clustering. Finally, cell-to-label influence is used for *bulk data mapping*. Bulk data can be mapped to labels by computing the fraction of each cell’s reads that overlap the bulk data and then summing the influence of each cell on each label. For example, an enhancer can be mapped to labels by using the fraction of reads from each cell that overlaps that enhancer and taking the sum of multiplying those fractions by the influence of each cell on each label. Like with cell labeling, this generates a fuzzy allocation of labels to the bulk data.

### Simulations

We generated artificial cells that emulate high-quality scATAC-seq processed by the SnapATAC [12] pipeline as follows. For each of  $n$  cells, we sampled the number of total reads for that cell from the distribution of reads per cell we observed in real data (median 5500 reads per cell, based on SnapATAC processed scATAC-seq data from Zifra et al. [19]). We then distributed those reads across  $p$  bins proportionally to the distribution of reads per bin observed in real data. This resulted in a  $p$ -by- $n$  count matrix of  $n$  cells and  $p$  bins. We split the pool of generated cells into two cell types and gave the cells low, medium, or high within-type distinctness by splitting bins evenly across cell types and adding a fixed percent (1, 5, or 10 respectively) of additional reads across those bins to each cell. In order to label cells, we generated two label nodes and created edges from these nodes to cells with a weight of 1 depending on the simulation scenario. Cell-to-cell edges were given a weight of the Jaccard similarity of each cell’s bins. For each simulation, we ran CellWalker on ten different assignments of cell-label edges for each of a range of label edge weights between  $10^{-4}$  and  $10^{-1}$  for 400 cells of each cell type. We evaluated the ability of CellWalker to separate the two cell types using cell homogeneity which we computed directly from the influence matrix  $F$  as the log of median ratio of information between cells within the same cell type to information between cells of different cell types. To test the importance of label-to-cell edges, we tested labeling between a single cell in each cell type up to

labeling 40% of cells using medium cell distinctness. We tested the importance of cell mis-labeling by labeling 15% of cells correctly and adjusting the number of mislabeled cells between a single cell and 15% of cells (not necessarily mutually exclusively). In tests for robustness to noisy reads, we randomly added or removed a fixed percentage of all reads in each cell. Finally, to test if CellWalker is able to distribute labels to cell populations even without any initial labeling, we generated an additional set of 400 cells with no labeling edges. Rather than give these cells medium, low, or high within group distinctness, they were made more similar to one of the previous cell types by being generated by randomly sampling reads proportionally to bins from either of the other two cell types, with the proportion of bins from the cell type adjusted between 10 and 50%. We additionally used simulated data to determine how CellWalker's runtime and memory usage scales with the number of non-zero bins in the cell-by-bin matrix, and found that both relationships are linear (Additional File 1: Fig. S3).

### Data processing and analysis

We downloaded the cell-by-peak matrix for the scATAC-seq portion and the cell-by-gene matrix for the scRNA-seq portion of the SNARE-seq data for the adult and developing mouse cerebral cortex (GEO accession number GSE126074) [20]. We additionally downloaded the cell type labels assigned to each cell, as well as marker genes for each cell type, which includes the log-fold change of expression for each marker in the given cell type compared to other cells. We ran CellWalker on this data by computing the Jaccard similarity between binarized peak accessibility vectors for cells for cell-to-cell edges and the fraction of each cell's peaks that are in marker's gene body or promoter (2 kb upstream of TSS) for a given cell type scaled by the log-fold change in expression of each marker for label-to-cell edges. We tested label edge weights between  $10^{-2}$  and  $10^4$  and computed both the cell homogeneity and the fraction of exact label matches at each weight (Additional File 1: Fig. S4). We found that the two follow nearly identical trends implying that cell homogeneity is a good proxy for correct labeling. For comparison, we ran Seurat3 [15] on the cell-by-peak and cell-by-gene matrices and assigned labels using default parameters for anchor transfer between the two datasets. We downloaded 10x Single Cell Multiome ATAC + Gene Exp chip data for human healthy brain tissue from 10x Genomics [34], which has scRNA-seq and scATAC-seq in the same cell, but no prior known labeling. We used Seurat to cluster the scRNA-seq and used those cluster assignments as labels. We then analyzed the data as described for SNARE-seq data.

Multi-sample mid-gestation human telencephalon scATAC-seq data from PsychENCODE ([synapse.org](https://synapse.org) id syn21392931) was previously processed using SnapATAC to generate a large cell-by-bin matrix [19], and a previously derived set of marker genes was used for labeling [17]. This data was generated using 10x scATAC-seq (as described in Ziffra et al. [19]) on primary samples of the human forebrain at midgestation from six individuals. These samples were derived from dorsolateral prefrontal cortex (PFC), primary visual cortex (V1), primary motor cortex (M1), primary somatosensory cortex, dorso-lateral parietal cortex, temporal cortex, insular cortex, and the medial ganglionic eminence (MGE). As before, cell-to-cell edge weights were computed using Jaccard similarity and label-to-cell edge weights were computed as the sum of normalized SnapATAC generated gene accessibility score for each marker scaled by that

marker gene's log-fold change in expression. We tested label edge weights between  $10^{-2}$  and  $10^4$  and selected a weight of one as optimal. We hierarchically clustered labeling nodes using the Euclidean distance between label-to-label vectors and “hclust” with default parameters in R [35]. To compute cell label scores from label-to-cell influence, we computed  $z$ -scores for each column and then rescaled to a maximum score of one. We considered a cell to have two nearly identical labeling scores if the top two highest scoring labels were within the bottom 10% of all differences between the two highest scoring labels for all cells.

To compute an excitatory neuron progression score for each cell, we took the difference between the nEN-early2 and nEN-late score for each cell. For our analysis, we only considered the subset of cells with nEN-early2 or nEN-late as their top label scores. 1367 germinal zone (GZ) TADS were previously generated [36] based on HiC data taken from Won et al. [22] For each TAD, we computed the fraction of each cell's distal (non-promoter, not in gene body) peaks that were accessible in that TAD. Then, for each TAD, we correlated excitatory neuron progression scores for each cell with the distal accessibility of each cell. Negative correlations imply early active TADs and positive correlations imply late active TADs. To determine if this distribution of correlations is significantly bimodal, we permuted TAD locations 1000 times using the randomizeRegions function in regioneR [37] (restricted to exclude blacklisted locations) and compared the median absolute values of correlations. Of the full set of TADs, we classified 290 as early TADs with correlation less than  $-0.5$  and 247 as late TADs with correlation  $> 0.5$ . Mean transcripts (TPM) were computed from the same scRNA data that cell type marker genes were derived from [17].

Regionally microdissected developmental brain pREs were downloaded from Markenscoff-Papadimitriou et al. (GEO accession number GSE149268) [18]. These pREs are regions predicted to be regulatory using machine learning trained on bulk epigenetic data and vista enhancers. To score each label for each pRE, we calculated the sum of cell-to-label influence across all cells which had a peak in the given pRE. We split cell types into four primary labels following the groupings in Nowakowski et al. [17] GWAS data was downloaded from the NHGRI-EBI GWAS catalog [25]. A list of genes was generated for each disease or trait using all mapped or reported genes for each significant variant as annotated in the catalog. We omitted diseases and traits that had fewer than 100 associated genes. Disease gene sets were downloaded from Werling et al. [26]. To estimate disease gene set enrichment, we computed a one-sided empirical  $p$  value by comparing the fraction of enhancers that were closest to disease genes to the fraction closest to an equally sized random sample of brain expressed genes, resampling 10,000 times. FDRs were computed using the Benjamini-Hochberg procedure. Cell tracks were generated using SnapATAC [12]. Peaks from FACS sorted cells were taken from Song et al. [30].

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-021-02279-1>.

**Additional file 1: Fig. S1.** Flowchart of CellWalker Pipeline. **Fig. S2.** Additional Simulation Results. **Fig. S3.** Runtime Analysis. **Fig. S4.** CellWalker Performance on SNARE-seq Data. **Fig. S5.** nEN Progression. **Fig. S6.** Cell Type-Specific Regulatory Elements.

**Additional file 2: Table S1.** Enrichment for cell type-specific pREs near genes with significant variants detected in a collection of GWAS studies.

**Additional file 3.** Review history.

### Acknowledgements

GZ TADs were generated and made available by Sean Whalen [36]. Ryan Ziffra and Tom Nowakowski provided early access to processed scATAC-seq data and helpful feedback on the project [19]. Thank you to Kathleen Keough for suggestions to improve manuscript clarity and John Rubenstein for help interpreting our results in the context of neurodevelopmental gene regulation.

### Review history

The review history is available as Additional file 3.

### Peer review information

Barbara Cheifet and Yixin Yao were the primary editors of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

### Authors' contributions

P.F.P. and K.S.P. designed the study, P.F.P. performed the analysis, and P.F.P. and K.S.P. wrote the manuscript. All authors read and approved the final manuscript.

### Funding

This work was supported by research grants to K.S.P. from: NIH/NIMH awards U01-MH116438, R01-MH109907, and R01-MH123178. Funding bodies played no role in the design of the study and collection, analysis, and interpretation of data, or in writing the manuscript.

### Availability of data and materials

No datasets were generated during the current study. All analyzed data is publicly available or available by request from the corresponding publications. SNARE-seq data for the adult and developing mouse cerebral cortex is available from GEO accession number GSE126074 [20]. 10x Single Cell Multiome ATAC + Gene Exp chip data for human healthy brain tissue is available from 10x Genomics [34]. Multi-sample mid-gestation human telencephalon scATAC-seq data is available from [synapse.org](https://synapse.org) id syn21392931. GZ TADS were previously generated [36] based on HiC data taken from Won et al. [22] and are available from those authors upon request. Regionally microdissected developmental brain pREs are available from GEO accession number GSE149268 [18]. GWAS data was downloaded from the NHGRI-EBI GWAS catalog [25]. Disease gene sets were downloaded from Werling et al. [26]. CellWalker code and simulated data is available under the GNU GPL-2.0 License at <https://github.com/PollardLab/CellWalker> (DOI: <https://doi.org/10.5281/zenodo.4456095>) along with a readme demonstrating how the method can be applied to sample data [38].

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

Not applicable.

### Competing interests

The authors declare that they have no competing interests.

### Author details

<sup>1</sup>Gladstone Institutes, San Francisco, CA, USA. <sup>2</sup>Chan-Zuckerberg Biohub, San Francisco, CA, USA. <sup>3</sup>Institute for Computational Health Sciences, Institute for Human Genetics, and Department of Epidemiology and Biostatistics, University of California, San Francisco, CA, USA.

Received: 3 August 2020 Accepted: 25 January 2021

Published online: 14 February 2021

### References

1. Visel A, Minovitsky S, Dubchak I, Pennacchio LA. VISTA Enhancer Browser—a database of tissue-specific human enhancers. *Nucleic Acids Res.* 2007;35(suppl\_1):D88–92. <https://doi.org/10.1093/nar/gkl822>.
2. Bulger M, Groudine M. Enhancers: the abundance and function of regulatory sequences beyond promoters. *Dev Biol.* 2010;339(2):250–7. <https://doi.org/10.1016/j.ydbio.2009.11.035>.
3. Dunham I, Kundaje A, Aldred SF, et al. An integrated encyclopedia of DNA elements in the human genome. *Nature.* 2012;489(7414):57–74. <https://doi.org/10.1038/nature11247>.
4. Wang D, Liu S, Warrell J, et al. Comprehensive functional genomic resource and integrative model for the human brain. *Science.* 2018;362(6420). doi:<https://doi.org/10.1126/science.aat8464>
5. Hoang TT, Goldmuntz E, Roberts AE, et al. The Congenital Heart Disease Genetic Network Study: cohort description. *PLoS ONE.* 2018;13(1). doi:<https://doi.org/10.1371/journal.pone.0191319>
6. Kundaje A, Meuleman W, Ernst J, et al. Integrative analysis of 111 reference human epigenomes. *Nature.* 2015;518(7539):317–30. <https://doi.org/10.1038/nature14248>.
7. Inoue F, Ahituv N. Decoding enhancers using massively parallel reporter assays. *Genomics.* 2015;106(3):159–64. <https://doi.org/10.1016/j.ygeno.2015.06.005>.

8. Hwang B, Lee JH, Bang D. Single-cell RNA sequencing technologies and bioinformatics pipelines. *Exp Mol Med*. 2018; 50(8):1–14. <https://doi.org/10.1038/s12276-018-0071-8>.
9. Haque A, Engel J, Teichmann SA, Lönnberg T. A practical guide to single-cell RNA-sequencing for biomedical research and clinical applications. *Genome Med*. 2017;9(1):75. <https://doi.org/10.1186/s13073-017-0467-4>.
10. Chen H, Lareau C, Andreani T, et al. Assessment of computational methods for the analysis of single-cell ATAC-seq data. *Genome Biol*. 2019;20(1):241. <https://doi.org/10.1186/s13059-019-1854-5>.
11. Pliner HA, Packer JS, McFaline-Figueroa JL, et al. Cicero predicts cis-regulatory DNA interactions from single-cell chromatin accessibility data. *Mol Cell*. 2018;71(5):858–871.e8. doi:<https://doi.org/10.1016/j.molcel.2018.06.044>
12. Fang R, Preissl S, Hou X, et al. Fast and accurate clustering of single cell epigenomes reveals cis-regulatory elements in rare cell types. *Bioinformatics*. 2019. <https://doi.org/10.1101/615179>.
13. Schep AN, Wu B, Buenostro JD, Greenleaf WJ. chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data. *Nat Methods*. 2017;14(10):975–8. <https://doi.org/10.1038/nmeth.4401>.
14. Bravo González-Blas C, Minnoye L, Pappasokrati D, et al. cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq data. *Nat Methods*. 2019;16(5):397–400. <https://doi.org/10.1038/s41592-019-0367-1>.
15. Stuart T, Butler A, Hoffman P, et al. Comprehensive integration of single-cell data. *Cell*. 2019;177(7):1888–902.e21. <https://doi.org/10.1016/j.cell.2019.05.031>.
16. Pliner HA, Shendure J, Trapnell C. Supervised classification enables rapid annotation of cell atlases. *Nat Methods*. 2019; 16(10):983–6. <https://doi.org/10.1038/s41592-019-0535-3>.
17. Nowakowski TJ, Bhaduri A, Pollen AA, et al. Spatiotemporal gene expression trajectories reveal developmental hierarchies of the human cortex. *Science*. 2017;358(6368):1318–23. <https://doi.org/10.1126/science.aap8809>.
18. Markenscoff-Papadimitriou E, Whalen S, Przytycki P, et al. A chromatin accessibility atlas of the developing human telencephalon. *Cell*. 2020;0(0). doi:<https://doi.org/10.1016/j.cell.2020.06.002>
19. Ziffra RS, Kim CN, Wilfert A, et al. Single cell epigenomic atlas of the developing human brain and organoids. *bioRxiv*. Published online January 8, 2020:2019.12.30.891549. doi:<https://doi.org/10.1101/2019.12.30.891549>
20. Chen S, Lake BB, Zhang K. High-throughput sequencing of the transcriptome and chromatin accessibility in the same cell. *Nat Biotechnol*. 2019;37(12):1452–7. <https://doi.org/10.1038/s41587-019-0290-0>.
21. Zhu C, Preissl S, Ren B. Single-cell multimodal omics: the power of many. *Nat Methods*. 2020;17(1):11–4. <https://doi.org/10.1038/s41592-019-0691-5>.
22. Won H, de la Torre-Ubieta L, Stein JL, et al. Chromosome conformation elucidates regulatory relationships in developing human brain. *Nature*. 2016;538(7626):523–7. <https://doi.org/10.1038/nature19847>.
23. Finn EH, Pegoraro G, Brandão HB, et al. Extensive heterogeneity and intrinsic variation in spatial genome organization. *Cell*. 2019;176(6):1502–15.e10. <https://doi.org/10.1016/j.cell.2019.01.020>.
24. Hor H, Francescato L, Bartesaghi L, et al. Missense mutations in TENM4, a regulator of axon guidance and central myelination, cause essential tremor. *Hum Mol Genet*. 2015;24(20):5677–86. <https://doi.org/10.1093/hmg/ddv281>.
25. Buniello A, MacArthur JAL, Cerezo M, et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res*. 2019;47(D1):D1005–12. <https://doi.org/10.1093/nar/gky1120>.
26. Werling DM, Brand H, An J-Y, et al. An analytical framework for whole-genome sequence association studies and its implications for autism spectrum disorder. *Nat Genet*. 2018;50(5):727–36. <https://doi.org/10.1038/s41588-018-0107-y>.
27. Rapanelli M, Frick LR, Pittenger C. The role of interneurons in autism and Tourette syndrome. *Trends Neurosci*. 2017; 40(7):397–407. <https://doi.org/10.1016/j.tins.2017.05.004>.
28. Lunden JW, Durens M, Phillips AW, Nestor MW. Cortical interneuron function in autism spectrum condition. *Pediatr Res*. 2019;85(2):146–54. <https://doi.org/10.1038/s41390-018-0214-6>.
29. Jin X, Simmons SK, Guo A, et al. In vivo Perturb-Seq reveals neuronal and glial abnormalities associated with autism risk genes. *Science*. 2020;370(6520). doi:<https://doi.org/10.1126/science.aaz6063>
30. Song M, Pebworth M-P, Yang X, et al. 3D epigenomic characterization reveals insights into gene regulation and lineage specification during corticogenesis. *bioRxiv*. Published online February 25, 2020:2020.02.24.963652. doi:<https://doi.org/10.1101/2020.02.24.963652>
31. Alexanian M, Przytycki PF, Micheletti R, et al. A transcriptional switch governing fibroblast plasticity underlies reversibility of chronic heart disease. *bioRxiv*. Published online July 22, 2020:2020.07.21.214874. doi:<https://doi.org/10.1101/2020.07.21.214874>
32. ArchR: An integrative and scalable software package for single-cell chromatin accessibility analysis | *bioRxiv*. Accessed December 16, 2020. <https://www.biorxiv.org/content/10.1101/2020.04.28.066498v1.full>
33. Leiserson MDM, Vandin F, Wu H-T, et al. Pan-cancer network analysis identifies combinations of rare somatic mutations across pathways and protein complexes. *Nat Genet*. 2015;47(2):106–14. <https://doi.org/10.1038/ng.3168>.
34. Maheshwari S, Chatterjee S, Sapida J, et al. Massively parallel simultaneous profiling of the transcriptomic and epigenomic landscape at single cell resolution.1.
35. R Core Team. R: A language and environment for statistical computing. Vienna: R Foundation for Statistical Computing; 2013. ISBN 3-900051-07-0, <http://www.R-project.org>.
36. Ryu H, Inoue F, Whalen S, et al. Massively parallel dissection of human accelerated regions in human and chimpanzee neural progenitors. *bioRxiv*. Published online January 29, 2018:256313. doi:<https://doi.org/10.1101/256313>
37. Gel B, Díez-Villanueva A, Serra E, Buschbeck M, Peinado MA, Malinverni R. regioneR: an R/Bioconductor package for the association analysis of genomic regions based on permutation tests. *Bioinformatics*. 2016;32(2):289–91. <https://doi.org/10.1093/bioinformatics/btv562>.
38. Przytycki P. CellWalker. GitHub; 2021. doi:<https://doi.org/10.5281/zenodo.4456095>.

## Publisher's Note

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