

# Correlation-based Linking of Epigenomic Regions to Target Genes in Bulk and Single Cell Epigenomic Data

by

Benjamin Thomas James

B.S. Computer Science and Mathematics  
University of Tulsa (2019)

Submitted to the Department of Electrical Engineering and Computer Science,  
School of Engineering,

in partial fulfillment of the requirements for the degree of

Master of Science

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2023

© Benjamin Thomas James 2023. License: CC BY-NC 4.0.

The author hereby grants to MIT a nonexclusive, worldwide, irrevocable, royalty-free license to exercise any and all rights under copyright, including to reproduce, preserve, distribute and publicly display copies of the thesis, or release the thesis under an open-access license.

Authored by Benjamin James  
Department of Electrical Engineering and Computer Science  
May 19, 2023

Certified by Manolis Kellis  
Professor of Electrical Engineering and Computer Science  
Thesis Supervisor

Accepted by Leslie A. Kolodziejski  
Professor of Electrical Engineering and Computer Science  
Chair, Departmental Committee on Graduate Students



# Correlation-based Linking of Epigenomic Regions to Target Genes in Bulk and Single Cell Epigenomic Data

by

Benjamin Thomas James

Submitted to the Department of Electrical Engineering and Computer Science,  
School of Engineering,  
on May 19, 2023, in partial fulfillment of the  
requirements for the degree of  
Master of Science

## Abstract

Single-cell sequencing has catalyzed a significant shift in biological modeling and hypothesis generation. In this study, we present advanced algorithms that utilize multiple modalities for peak-gene linking via correlation-based mechanisms. This approach provides a robust framework for bias correction and noise reduction at single-cell resolution.

Our research introduces a novel algorithm that employs peak-gene linking to integrate epigenomic data, thereby translating it into transcriptomic data. This integration facilitates comprehensive secondary analyses on meta-cells and sub-cell types. Crucially, our methodology enables swift computation of modules from any single-cell assay, promoting exploration of intricate biological and disease mechanisms that might remain unmodeled with a pseudo-bulk approach.

By combining snRNA-seq and snATAC-seq data, our method substantially outperforms equivalent tasks of gene expression estimation, showing promising results even with unpaired real data. Furthermore, the modeling of genomic peak modules using our algorithms uncovers additional signal potentially overlooked when examining single peaks.

We envision correlation-based linking as a key aspect of future single-cell multiomic technology, as it allows for correlation between assays at the single-UMI level. As such, improved modeling of data at the single-cell level will enhance our understanding of complex gene regulatory networks and disease mechanisms.

Thesis Supervisor: Manolis Kellis

Title: Professor of Electrical Engineering and Computer Science,



# Contents

<b>Abstract</b>	<b>3</b>
<b>1 Introduction</b>	<b>9</b>
<b>2 Correlation-based enhancer-gene linking on bulk samples</b>	<b>11</b>
2.1 Initial correlation-based enhancer-gene linking on bulk samples . . . .	12
<b>3 Integration of snATAC-seq and snRNA-seq data through correlation-based peak-gene linking</b>	<b>15</b>
3.1 Integration of snATAC-seq and snRNA-seq data as a peak-gene linking task . . . . .	15
3.2 An iterative strategy for integrating epigenomic and transcriptomic data	17
3.3 More accurate estimation of gene expression from joint snRNA and snATAC peak-gene linking . . . . .	20
3.4 Results on snRNA-seq and snATAC-seq aged human brain data . . .	28
<b>4 Peak modules from snATAC-seq data</b>	<b>33</b>
4.1 Peak modules . . . . .	33
<b>5 Discussion</b>	<b>37</b>



# List of Figures

2-1	Statistics of initial peak-gene links from EpiMap [3]. <b>a</b> Histogram of the median number of genes per enhancer per bin across biosamples. <b>b</b> Histogram of the median number of enhancers per gene per bin across biosamples. <b>c</b> Total number of sample groups in which a unique link is active, out of 3.3M unique enhancer-gene links. <b>d</b> Histogram plot of parent unique elements (enhancers in yellow, links in green) for which the element is active in a given number of biosamples. <b>e</b> Histogram showing the median distance between the enhancer and TSS of a link across biosamples. <b>f</b> Mean (blue) or median (red) enhancer-gene link distance for all enhancers active in a given number of biosamples. <b>a-f</b> : All means represented by blue dashed lines and text and medians by red dashed lines and text. <b>g-h</b> : Comparison (prediction F1 score and AUPRC) of gene-enhancer link predictions (blue) with distance, activity by distance, correlation by distance, correlation and activity by distance, and Roadmap Epigenomicslinks on functional gene set-based links ( <b>g</b> ) and physical, genetic, and perturbation-based links across four cell lines ( <b>h</b> ). . . . .	14
3-1	Integration framework cartoon. Note the iterative process that uses peak-gene linking to improve correlations, and that each “integration” is done at the major cell type level instead of overall in order to model intra cell type variance instead of inter cell type variance . . . . .	18
3-2	Overview of a first-pass iteration of the algorithm, including the binning and adjustments to model sparsity. After the first iteration, the computed meta-cells are included. . . . .	21
3-3	Plots of marker genes for each cell subtype, computed for each integration method (the one described), scGLUE, and Seurat/Signac . . .	30
3-4	Pearson correlation coefficient across marker genes per subtype . . . .	31
3-5	AUROC/AUPRC linking plots at the major cell type level (Ex=Excitatory neurons, In=Inhibitory neurons, Oligo=Oligodendrocytes, Microglia=Microglia and immune cells) using Nott et al PLAC-seq data. This method performs on average better than the currently used peak-gene linking methods for single cell. Notice that the scale on the AUROC plots — all the AUC values are near 0.5 due to the difficulty of this task. Additionally, note that depending on the peak set, distance from peak to gene varies in accuracy. . . . .	32

4-1 Cell type specific ATAC-seq modules capture distinct biological programs via Gene Ontology in aged human brains . . . . . 35



# Chapter 1

## Introduction

Human disease is incredibly complex, as many mechanisms lead to both onset as well as disease progression. There are many factors contributing to disease, including both genetic and environmental effects. However, many of these effects are either unknown or not very well characterized at the molecular level. Since the sequencing of the human genome, fields such as genomics, epigenomics, and transcriptomics have emerged to study these effects at different levels. Along with genetic testing, technologies include RNA sequencing, which measures the level of RNA transcripts which can code for proteins in a given tissue, and chromatin accessibility, which measures the openness of the DNA in a given tissue or cell type. Massive studies have demonstrated numerous regions of the genome that contribute to disease, mostly at the single nucleotide polymorphism (SNP) level. However, the vast majority of these disease critical regions are in regions of the genome that do not code for proteins. Projects such as Roadmap Epigenomics [18], ENCODE [7], IHEC [20], and EpiMap [3] have found that these disease critical SNPs are located in control regions such as enhancers and promoters that modify gene expression levels instead of directly altering function. Previous large-scale projects have been hindered by cell type heterogeneity within tissues of interest. For example, the brain is composed of several different cell types,

including many types of neurons, glia, and vasculature. Tissue level sequencing is not sufficient to uncover molecular level contributions to disease. However, recent technology allows for sequencing of RNA transcripts or DNA accessibility at the single-cell or single-nucleus resolution using single cell/nucleus RNA sequencing (scRNA-seq) and single nucleus assay for transposable accessible chromatin sequencing (snATAC-seq). Recently, methods to measure chromatin accessibility and gene expression from the same cell (deemed single-cell multiomics) has been commercially available and is starting to become widely adopted. When many of these experimental breakthroughs in technology occur, the computational methods and analysis using them has to also catch up to fully take advantage. Methods to link non-coding regulatory elements to target genes has been studied at the bulk level [8] [10], but not as well at the single-cell level, let alone the single cell multiomics level. First, we have created a methodology to link enhancers to target genes across 833 bulk samples in EpiMap based on epigenome-transcriptome correlations. We then adapted this technology to work at the single-cell level. We discovered this linking allows for better estimation of gene expression, which allows us to better link enhancers to the correct target genes. Next, we adapt this technology to find co-accessible modules of chromatin regulation at the single-cell level, capturing more regulation than what can be found from bulk or pseudo-bulk based chromatin modules, and finally, we can estimate 3D chromatin architecture from this single-cell accessibility additionally from single-cell chromatin accessibility. We apply these methods to single-nucleus ATAC sequencing and single-nucleus multiomic sequencing.

## Chapter 2

# Correlation-based enhancer-gene linking on bulk samples

Finding target genes of disease-relevant regulatory elements is an important step in understanding disease circuitry. Most disease-relevant SNPs are distal to their target genes, and of these are not the closest gene to the regulatory element. For example, the FTO locus identified via GWAS is not connected to the closest gene, FTO, but is instead linked to distal IRX3 and IRX5 genes [5] Previously, correlation based linking was used to determine regulatory regions connected to genes using correlations between histone marks and gene expression[8]. We started from this methodology and updated it for larger samples.

## 2.1 Initial correlation-based enhancer-gene linking on bulk samples

We established an estimated 3.3 million correlations between specific tissue enhancers and genes. We did this by integrating data from epigenomic-transcriptional correlation and genomic proximity. On average, each gene was connected to 13 enhancers and each enhancer was linked to approximately 1.5 genes. The median separation between them was around 42,359 base pairs. The specificity of these links was roughly six times greater than that of the enhancers. The links specific to certain samples covered larger distances compared to constitutive links. Compared to previous methodologies, our links demonstrated superior performance. This conclusion was based on gene-set enrichment metrics and carefully curated benchmark datasets. Furthermore, we substantially increased the number of biosamples with predicted links, rising from 127 in Roadmap Epigenomics to 833. This data was released in EpiMap [3].

In our study, we grappled with a significant issue that emerged in our Roadmap-inspired linking strategy, and was further exacerbated in the EpiMap dataset: the heterogeneity of the samples in EpiMap. Particularly, the abundance of blood and immune samples appeared to inflate counts, thus potentially skewing correlations to mirror the sampling distribution instead of the correlation we aimed to model in our tissue of interest. To address this, our work concentrated on eliminating this inflated correlation due to sample composition.

We employed zero-phase component analysis (ZCA) whitening [2], a technique that standardized our dataset and enhanced the discernability of less common cell types. This process, essentially, de-correlated the dataset and ensured that the data's variance was well modeled, which significantly improved our predictive capabilities.

To further refine our data analysis, we turned to the autocorrelation derived from

the singular value decomposition (SVD) of the data to establish sample-to-sample correlations. Armed with these correlations, we shifted our focus from raw Pearson correlations to weighted correlations to re-introduce the critical factor of cell type specificity. This method proved to be particularly effective, as these weighted correlations within the ZCA-whitened space resulted in improved capture of gene regulatory networks.

For better modeling of extremely high and low correlations in the data, we applied Fisher-transformed correlations, specifically the hyperbolic arctangent [9]. This technique is designed to normalize and stabilize the variance of correlation coefficients, thus making them more comparable and better suited for hypothesis testing.

Our normalization approach was inspired by the Activity-by-Contact (ABC) model [10], with necessary modifications to accommodate the predicted link probability and the activity of the sample as the "activity". This normalization process was crucial to the successful interpretation of our correlation models:

$$\frac{Activity * Contact * Probability}{\sum_{|Peak-Gene| < 1\text{Mbp}} Activity * Contact * Probability} \quad (2.1)$$

In particular, we used the Hi-C power of -0.87 and scale parameter of 5.41 (parameters to the ABC model) to estimate contact, with activity defined by pseudo-bulk accessibility, and the probability as mentioned above.

One crucial realization stemming from our work with the EpiMap dataset was the importance of a well-structured negative training set and an appropriate normalization scheme in correlation modeling. We found these to be key components in obtaining accurate and reliable results. As such, future work must pay careful attention to these aspects to ensure the validity and applicability of correlation models.

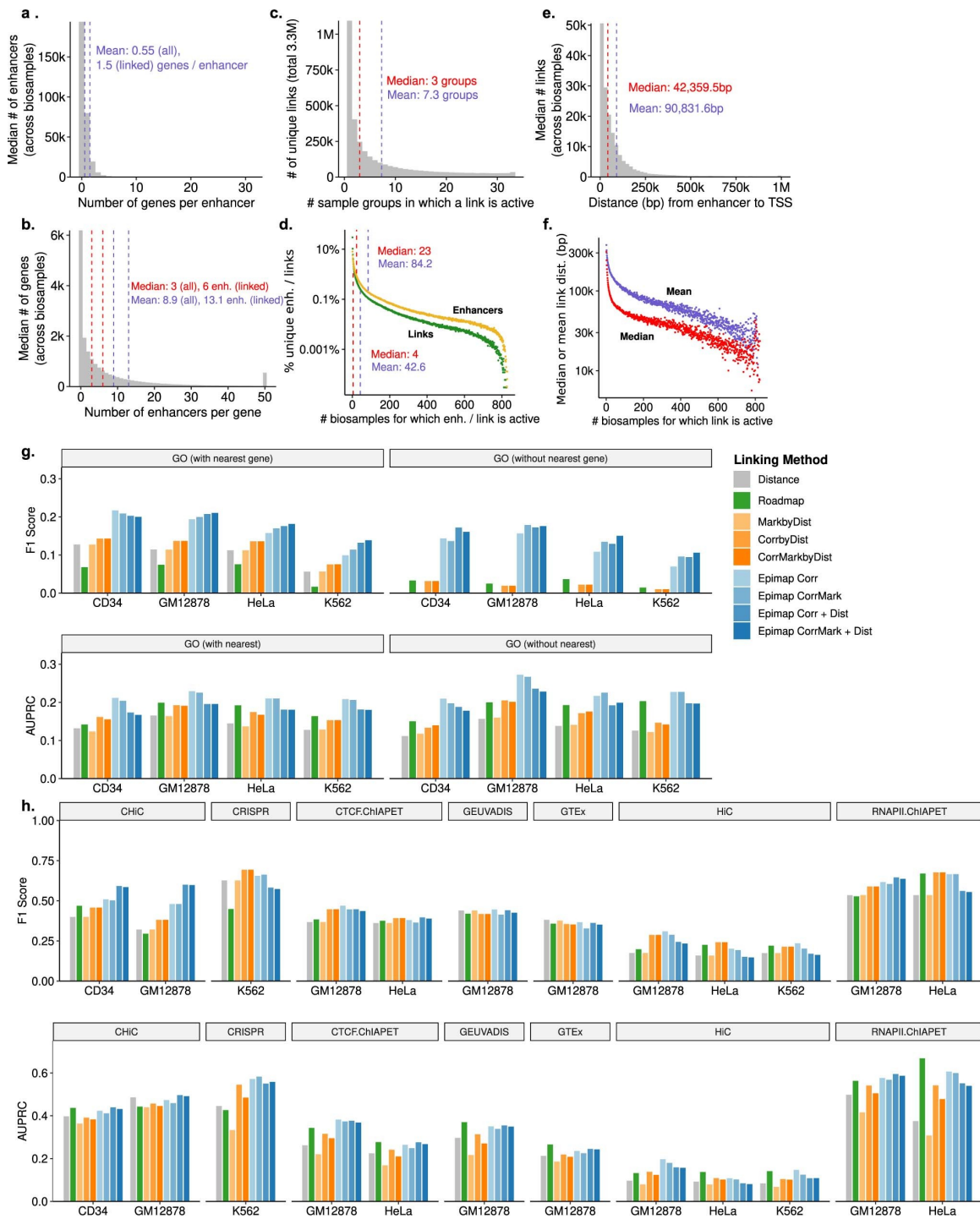


Figure 2-1: Statistics of initial peak-gene links from EpiMap [3]. **a** Histogram of the median number of genes per enhancer per bin across biosamples. **b** Histogram of the median number of enhancers per gene per bin across biosamples. **c** Total number of sample groups in which a unique link is active, out of 3.3M unique enhancer-gene links. **d** Histogram plot of percent unique elements (enhancers in yellow, links in green) for which the element is active in a given number of biosamples. **e** Histogram showing the median distance between the enhancer and TSS of a link across biosamples. **f** Mean (blue) or median (red) enhancer-gene link distance for all enhancers active in a given number of biosamples. **a-f**: All means represented by blue dashed lines and text and medians by red dashed lines and text. **g-h**: Comparison (prediction F1 score and AUPRC) of gene-enhancer link predictions (blue) with distance, activity by distance, correlation and activity by distance, and Roadmap Epigenomlinks with functional gene set-based links (**g**) and physical, genetic, and perturbation-based links across four cell lines (**h**).

# Chapter 3

## Integration of snATAC-seq and snRNA-seq data through correlation-based peak-gene linking

### 3.1 Integration of snATAC-seq and snRNA-seq data as a peak-gene linking task

Advances in single-cell sequencing technologies have revolutionized our ability to investigate cellular heterogeneity and regulatory processes at unprecedented resolution. Among these techniques, snATAC-seq and snRNA-seq have emerged as powerful tools for studying chromatin accessibility and gene expression patterns, respectively. Integrating these two datasets offers a promising avenue to gain a deeper understanding of cellular function and regulation.

Due to the nature of the central dogma of molecular biology, DNA has only two copies, and there are two strands in normal diploid mammalian cells. Therefore, for

a single cell, when the Tn5 cuts the DNA, the quantitative readout is not as reliable as counts of RNA molecules or proteins. Since most applications of this technology have no common axis, it can be challenging to find shared molecular function between snATAC-seq and snRNA-seq. While snRNA-seq provides insights into gene expression patterns, snATAC-seq captures open chromatin. This often overlaps with regulatory elements such as enhancers, promoters, transcription factor binding sites, and gene bodies. Therefore, estimation of gene expression from open chromatin not only allows for a task that helps integration, but this estimation helps understand regulatory mechanisms in an understandable fashion.

Because of the differences mentioned between DNA and RNA, cell type identification in snRNA-seq is much easier due to the distribution of RNA molecules in the nucleus or cytosol. Cell type identification is often the most important first step for understanding complex biological systems. So, by integrating epigenomic and transcriptomic data, researchers can achieve improved cell type identification.

Additionally, the combination of chromatin and expression data allows for the modeling of regulatory elements and their influence on gene expression. Often times, snRNA-seq studies find differentially expressed genes along common gene pathways and programs. Enrichment of the binding sites of the promoters of these genes often hints at regulatory mechanisms. However, just like with snRNA-seq, snATAC-seq integrative analysis allows for researchers to identify potential regulatory elements associated with differentially expressed genes. This integration helps in identifying candidate transcription factors and cis-regulatory elements that may drive cell type-specific gene expression programs. This is often the cases in cell state/fate analysis such as hematopoiesis. Then by combining information on chromatin accessibility and gene expression, researchers can infer direct regulatory interactions between transcription factors and target genes. This integration facilitates the precise reconstruction of regulatory networks and the identification of key regulatory modules controlling



specific cellular processes.

## 3.2 An iterative strategy for integrating epigenomic and transcriptomic data

We have formulated a unique strategy for integrating single-nucleus assay for transposase-accessible chromatin (snATAC) and single-nucleus RNA sequencing (snRNA) data. This approach cyclically refines the snATAC gene score matrix, the snATAC-snRNA integration, and the peak-gene connections. Initiation of this integration requires an initial estimated gene score matrix for the snATAC cells, but no preliminary information about distal regulatory links or snRNA-snATAC integration is required. The process involves the following key steps:

- Integration of snRNA data (gene expression matrix) and snATAC data (estimated gene score matrix)
- Assignment of sub-cell types to snATAC cells derived from the joint snATAC-snRNA space
- Refreshing peak-gene connections based on snRNA-snATAC integration results
- Recalculation of the gene score matrix for snATAC data using the updated peak-gene connections
- Reiteration of steps 1 through 4 to re-integrate snRNA and snATAC data.

Each of these steps is further elucidated in the following sections (integration, cell type assignment, peak-gene linking). The algorithm is deemed to have converged when the Pearson correlation between snRNA gene expression and the revised snATAC gene

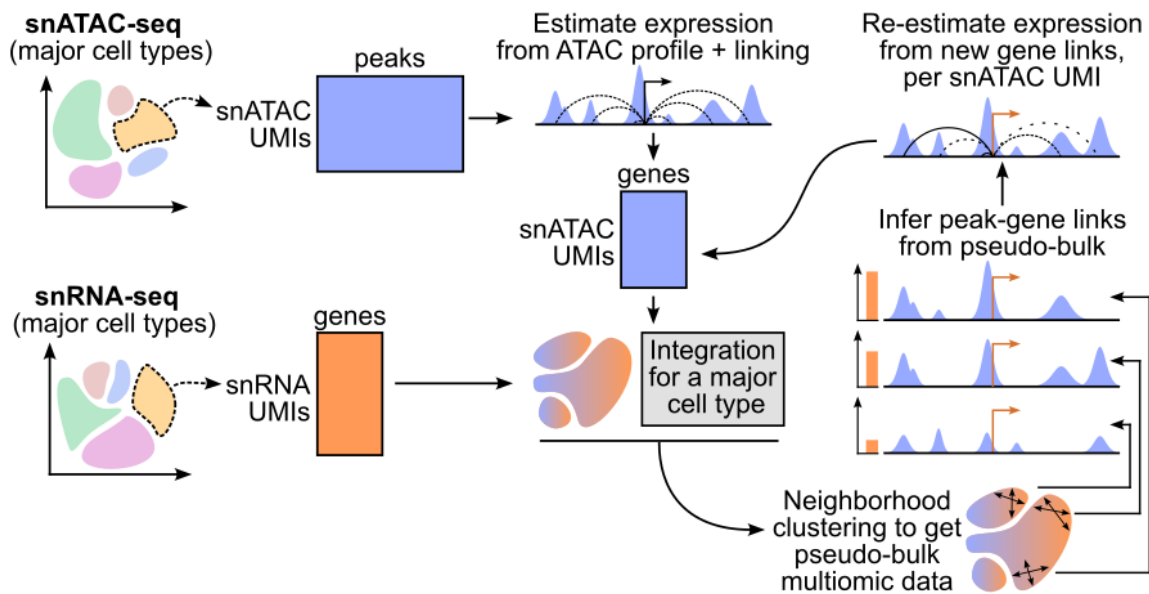


Figure 3-1: Integration framework cartoon. Note the iterative process that uses peak-gene linking to improve correlations, and that each “integration” is done at the major cell type level instead of overall in order to model intra cell type variance instead of inter cell type variance

score plateaus. For instance, in the case of the prefrontal cortex (PFC) data, we observed minimal to no improvement beyond two iterations of updating links and recalculating the gene score matrix. However, given that each iteration is performed manually, the user is granted the discretion to decide when integration sufficiency is reached.

Since major cell types can be discerned in both snATAC and snRNA, we separate and integrate each major cell type independently. This allows for more accurate examination of the variability within cell types, with neurons as a notable example.

**The method is as follows:** Our preliminary computation of the single-nucleus assay for transposase-accessible chromatin (snATAC) gene score matrix mirrors the methodology used by ArchR [11], which infers gene-level expression from chromatin accessibility. To elaborate, we generate a set of 200 base pair (bp) tiles spanning the genome and count the number of reads located within 5 kilobases (kb) of each gene

body. This count is then weighted using a score that decreases exponentially based on distance, specifically using the formula  $e^{-\frac{|distance\_to\_genebody|}{5000}} + e^{-1}$ . This method allows us to emphasize proximal peaks while minimizing the influence of distal peaks.

We initially restrict the contributing distal peaks for each gene to terminate at the next gene’s body. In addition, we apply an extra weight to the gene score to counterbalance gene length. This inverse weight means that shorter genes receive a score of 5 while the longest genes receive a score of 1. Lastly, the unique molecular identifier (UMI) counts for each snATAC cell are adjusted so that the total count across all genes equals 10,000, aligning with single-nucleus RNA sequencing normalization strategies.

Once the snRNA-based peak-to-gene link information becomes available in subsequent integration rounds, we incorporate these links into the snATAC gene score matrix estimation. In particular, we adjust the contribution of distal peaks based on the peak-to-gene linking probability scores computed by the logistic regression model. We don’t limit the calculation of the gene score to nearby gene boundaries when using peak-gene links but rather use all gene-linked peaks (peak-gene links computed for peaks within 1Mb of the gene TSS).

After obtaining gene scores for both ATAC and RNA modalities, we carry out cell type integration separately. For each cell type, we scale the ATAC and RNA gene matrices to 10,000 total counts per cell, followed by  $\log(1 + x)$  scaling. Given the substantial corrections required between ATAC and RNA, simple batch correction is insufficient to fully counteract the effects. To rectify the ATAC-RNA differences, we sequentially run ComBat [13], Harmony [14], and BB-kNN [17] through ScanPy [23] [22]. Subsequently, we apply a diffusion map as implemented in Scanpy [6] on the nearest neighbor matrix to produce a latent space, which is then used to calculate a UMAP embedding [15].

Finally, using an ATAC-RNA integration, cell types are assigned by counting the

number of sub cell types that are nearest neighbors from the nearest neighbor matrix used to build the integrated UMAP. Leiden clustering [21] is also applied on this neighbor matrix, with two graphs from resolution parameters of 1 and 2 to build two levels of meta-cells from which peak-gene linking is estimated. However, for the final integration UMAP, we instead use predicted cell type labels to create meta-cells to predict links. This is useful because meta-cells [1] allow coarse resolution alignment and also fix many sparsity issues inherent with this high dimensional data.

### **3.3 More accurate estimation of gene expression from joint snRNA and snATAC peak-gene linking**

We develop a method for computing peak-gene links from either snATAC alone (first iteration of the integration) or from integrated snATAC-snRNA (subsequent iterations). At a high level, we take the following steps for the first iteration.

- Estimate and rescale the gene accessibility (gene score) matrix
- Perform TF-IDF on and rescale the peak matrix
- Concatenate the peak and gene accessibility matrices
- Perform PCA on this concatenated matrix
- Re-scale the PCA embeddings with different weights and concatenate these rescaled embeddings
- Take dot product between peak and gene across embeddings to estimate their correlation

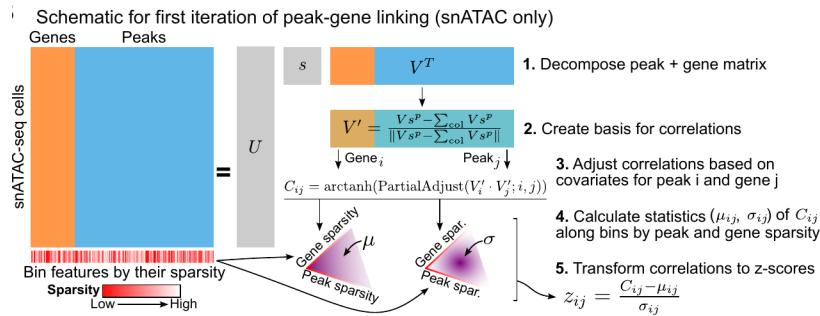


Figure 3-2: Overview of a first-pass iteration of the algorithm, including the binning and adjustments to model sparsity. After the first iteration, the computed meta-cells are included.

- Compute the partial correlation adjusted for quality covariates, and transform by  $\text{arctanh}$
- Bin each peak-gene pair into a bin based on the peak sparsity and the gene sparsity
- Compute per-bin statistics (mean, sd) of the correlation (over all (peak, gene) pairs in a bin)
- Adjust the scores for individual pairs by z-scoring with the (mean, sd) statistics from their bin
- Predict link probability from these scores using a logistic regression model
- Multiply link probability by peak accessibility and normalize this as in ABC (accounting for gene density and distance) to obtain the final linking scores

We describe these steps and the modifications for subsequent iterations in more detail below: Peak-gene correlation is based on combined co-accessibility (i.e. within ATAC-seq) and peak-gene correlation (i.e. across assays).

**Estimation and re-scaling of gene score matrix** We calculated snATAC gene accessibility matrix (or gene score matrix) as part of the integration as outlined above,

either using the ArchR method alone (in initialization) or peak-gene links reweighting the relative contribution of distal peaks (where the links used are from the previous round of integration). We first calculate peak-gene accessibility PC embeddings per UMI for co-accessibility and peak-gene accessibility-gene expression PC embeddings per metacell for peak-gene expression linking. Only peaks within 1Mb of the gene TSS are considered for peak-gene linking.

**Perform TF-IDF on and rescale the peak matrix** To accurately scale peaks and gene accessibility, we ensure that each of these modalities has a per-UMI sum of 10,000. To do so, gene accessibility scores are scaled to 10,000 per UMI. However, due to the sparsity inherent in an snATAC peak matrix, we utilize term frequency inverse document frequency (TF-IDF) to initially scale the peak matrix. Next, we scale the TF-IDF matrix to 10,000 per-UMI counts to ensure compatibility with gene accessibility.

**Dimensionality reduction of peak matrix** Finally, we concatenate the 10,000 per-UMI counts and  $\log(1 + x)$  scale the concatenated matrix before we call truncated PCA with 100 principal components. After we perform an initial integration with snRNA, we include that data in order to more accurately capture correlations with transcripts.

**Re-scale the PCA embeddings with different weights and concatenate these rescaled embeddings** Similar to co-accessibility embeddings, we calculate peak-gene embeddings in a similar manner. However, instead of calculating the embeddings per-UMI, we utilize a metacell like framework to combine ATAC and RNA modalities in a high enough resolution that preserves accuracy. Here we use sub cell types per individual as meta cells. Similar to co-accessibility estimation, we use TF-IDF for the peak matrix followed by scaling to 10000 counts per meta cell as well as 10,000 counts per UMI for both gene accessibility and gene expression. In total, when combined, each meta-cell has 30,000 counts. Then, we utilize  $\log(1+x)$  scaling

followed by truncated PCA with 100 components to form a similar embedding. This framework utilizes varying powers of the eigenvalues of the PC embeddings to generate five distinct embeddings scaling from zero-phase component analysis (ZCA) [2] to principal component analysis (PCA): From a peak matrix, initially correlations can be done over PCA space from a decomposition  $U \cdot s \cdot V^T$ :

$$\rho_{ij}^* = \text{cor}(V_i * s, V_j * s) \quad (3.1)$$

From a PCA decomposition, ZCA can be computed as  $UV^T$ . Thus,  $VV^T$  gives the covariance matrix as  $U$  is unitary.

When scaling  $V$  by subtracting the mean and  $L_2$  norm-ing, the covariance of  $V'$  is the autocorrelation matrix of  $V$ .

$$V'_i = \frac{V_i - \mu_{V_i}}{\sqrt{V_i - \mu_{V_i}}} \quad (3.2)$$

Reweighting is performed following the equation below, with eigenvalues  $s$ , eigenvectors  $V$  from PCA and power  $p$  as one of (0, 0.25, 0.5, 0.75, 1), where the resulting embeddings are adjusted by the eigenvalues, centered, and normalized:

When done for both co-accessibility and peak-gene embeddings, we use 10 total embeddings; when gene expression meta-cells are included, we use 15 total embeddings. In these embeddings, we subtract their per-UMI means and Euclidean-norm the resulting matrix so that the dot product in that space approximates Pearson correlation.

**Take dot product between peak and gene across embeddings to estimate their correlation** From these embeddings, we calculate initial scores between peaks and genes by taking the dot product between two indices in the embeddings. Often, single cell data is noisy is affected by sample quality or condition such as post-mortem interval or TSS enrichment. Correlations per-sample can be driven by these metrics.

Partial correlations allow for adjustment due to these factors. In order to remove the effect of TSS enrichment and number of fragments on these correlations, we compute partial correlations over the peak-gene embeddings instead of raw correlations. Partial correlations allow for adjustment due to these factors but the classical partial correlation algorithm requires a matrix inverse of the autocorrelation matrix, which is infeasible for peak-peak or peak-gene correlations.

However, the classical partial correlation algorithm requires a matrix inverse of the autocorrelation matrix, which is infeasible for peak-peak or peak-gene correlations.

So, given a full correlation matrix with covariate correlations with features  $i$  as  $\mathbf{b}_i$  and  $j$  as  $\mathbf{b}_j$ , and covariate autocorrelation matrix as  $\mathbf{A}$ .

$$\Sigma = \begin{pmatrix} 1 & \rho_{ij} & \mathbf{b}_i \\ \rho_{ij} & 1 & \mathbf{b}_j \\ \mathbf{b}_i^T & \mathbf{b}_j^T & \mathbf{A} \end{pmatrix} \quad (3.3)$$

$$= \begin{pmatrix} \mathbf{R} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{A} \end{pmatrix} \quad (3.4)$$

Using the matrix inverse of  $\Sigma$ , partial correlations can be computed:

$$\rho_{ij|\mathbf{V}\{\mathbf{A}\}} = \frac{-\Sigma_{ij}^{-1}}{\sqrt{\Sigma_{ii}^{-1}} \sqrt{\Sigma_{jj}^{-1}}} \quad (3.5)$$

Since any autocorrelation matrix  $\Sigma$  is guaranteed to be full rank, the inverse of  $\Sigma$ , the precision matrix, exists and must be:



$$\Sigma^{-1} = \begin{pmatrix} \mathbf{R} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{A} \end{pmatrix}^{-1} \quad (3.6)$$

$$= \left( \begin{array}{c|c} (\mathbf{R} - \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B})^{-1} & \dots \\ \dots & \dots \end{array} \right) \quad (3.7)$$

$$= \left( \begin{array}{c|c} \begin{pmatrix} 1 - \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_i & \rho_{ij} - \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_j \\ \rho_{ij} - \mathbf{b}_j^T \mathbf{A}^{-1} \mathbf{b}_i & 1 - \mathbf{b}_j^T \mathbf{A}^{-1} \mathbf{b}_j \end{pmatrix}^{-1} & \dots \\ \dots & \dots \end{array} \right) \quad (3.8)$$

$$= \left( \begin{array}{c|c} \frac{1}{(1 - \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_i)(1 - \mathbf{b}_j^T \mathbf{A}^{-1} \mathbf{b}_j) - (\rho_{ij} - \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_j)(\rho_{ij} - \mathbf{b}_j^T \mathbf{A}^{-1} \mathbf{b}_i)} & \begin{pmatrix} 1 - \mathbf{b}_j^T \mathbf{A}^{-1} \mathbf{b}_j & \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_j - \rho_{ij} \\ \mathbf{b}_j^T \mathbf{A}^{-1} \mathbf{b}_i - \rho_{ij} & 1 - \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_i \end{pmatrix} \\ \dots & \dots \end{array} \right) \begin{array}{c} \dots \\ \dots \end{array} \quad (3.9)$$

$$(3.10)$$

Then, using this to determine the partial correlations:

$$\rho_{ij\mathbf{V}\{\mathbf{A}\}} = \frac{-\Sigma_{ij}^{-1}}{\sqrt{\Sigma_{ii}^{-1}} \sqrt{\Sigma_{jj}^{-1}}} \quad (3.11)$$

$$= \frac{\frac{\rho_{ij} - \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_j}{|\mathbf{R} - \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B}|}}{\sqrt{\frac{1 - \mathbf{b}_j^T \mathbf{A}^{-1} \mathbf{b}_j}{|\mathbf{R} - \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B}|}} \sqrt{\frac{1 - \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_i}{|\mathbf{R} - \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B}|}}} \quad (3.12)$$

$$= \frac{S (\rho_{ij} - \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_j)}{\sqrt{S (1 - \mathbf{b}_j^T \mathbf{A}^{-1} \mathbf{b}_j)} \sqrt{S (1 - \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_i)}} \quad (\text{Using } S = \text{sgn } |\mathbf{R} - \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B}|) \quad (3.13)$$

Since  $\mathbf{A}^{-1}$  and  $\mathbf{b}$  can be computed beforehand, the pairwise partial correlations can be easily calculated from this formula. The sign  $S$  can be trivially calculated using

intermediate products.

In the rest of this document, partial correlations can be substituted in formulae instead of correlations for ease of reading.

In sum, given a correlation  $\rho_{ij}$ , covariate correlations as  $\mathbf{b}_i$ , and covariate auto-correlation matrix as  $\mathbf{A}$ . We reduce this to the following equation, which we use to compute the partial correlation:

$$\rho_{ij\mathbf{V}\{\mathbf{A}\}} = \frac{S(\rho_{ij} - \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_j)}{\sqrt{S(1 - \mathbf{b}_j^T \mathbf{A}^{-1} \mathbf{b}_j)} \sqrt{S(1 - \mathbf{b}_i^T \mathbf{A}^{-1} \mathbf{b}_i)}} \quad (\text{Using } S = \text{sgn} |\mathbf{R} - \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B}|)$$

**Computation of partial correlations transformed to stabilize variance** We transform these correlations using the hyperbolic arctangent (Fisher’s transform), to stabilize the variance of these correlations, which is especially important when quantifying high correlations.

**Sparsity adjustments** In single cell data, the high amount of zeros and variability of count data means that even in PCA space, the influence of count data is highly correlated with many principal components.

To adjust for this, we adjust these correlations or partial correlations based on expected correlations or partial correlations from the same number of counts per feature (peak).

We use Fisher’s transform,  $\rho' = \text{arctanh}(\rho)$  to stabilize the variance of the correlations, then sample a 50 by 50 grid of all features against each other, stratified by percentile sparsity, to get 50x50 means and standard deviations for each bin.

Namely, we stratified all peaks by their sparsity into 50 bins and all gene accessibility scores into another 50 bins by sparsity. We group peaks and genes into these bins by sparsity in order to find sparsity bin-specific baselines from which we can remove

contributions to correlations inflated by technical issues. For each pair of bins (peak sparsity and gene sparsity), we calculate the mean and standard deviations of all the scores for peak and gene tuples that fall into this pair of bins. However, due to the computational intractability of calculating all peak-gene pairs to estimate the mean and standard deviation of peak-gene correlations is infeasible computationally, we estimate these metrics by sampling. Since Fisher transformed Pearson correlations have a known standard error [9], we perform a power calculation (Z=2 and margin of error of 0.05) to compute how many peak-gene pairs to sample for our estimates for each sparsity bin (gene, peak) pair.

Finally, we smooth these means and standard deviations across sparsity bins using a Gaussian smoothing filter over these two grids to ensure quality is uniform and that outliers from sampling do not confound analysis.

Once we set the mean and standard deviation for each peak-gene pair, we calculate z-scores for each peak-gene pair, such that the transformed and z-scored partial correlation between gene  $i$  and peak  $j$  is, for embedding  $V'$ :

$$\rho_{ij}^z = \frac{\text{arctanh}(V'_i \cdot V'_j) - \mu_{\{i,j\}}}{\sigma_{\{i,j\}}} \quad (3.14)$$

### **Computation of per-bin statistics to model sparsity effects on correlations**

Where the mean and standard deviation are the bin-wise estimates described above. We then normalize the Fisher transformed partial correlations by the mean and standard deviation derived from the bins mentioned above to produce z-scores for each of the 10 embeddings. In order to find baselines for each gene for peak-gene linking, we calculate corresponding random-peak to gene z-scores by taking 5 random peaks and correlating them to the target gene for each peak.

**Multiple ZCA-PCA correlations** Since ZCA does remove the effect of larger principal components, we have found that including PCA as well as powers of the

eigenvalues between 0 and 1 seems to provide complementary informatoin.

$$V_i'^p = \frac{V_i * S^p - \mu_{V_i * S^p}}{\sqrt{V_i * S^p - \mu_{V_i * S^p}}} \quad (3.15)$$

We find that 5 embeddings between fully ZCA and PCA via powers of the eigenvalues, from 0 to 1, encapsulate cell state well.

**Predicting link probability from correlation scores using a logistic regression model and ABC-style normalization** Finally, we train a logistic regression model to distinguish peak-to-gene z-scores against 5 random peak-to-gene z-scores per candidate peak-gene link as mentioned above. We sample the random peaks from a different chromosome to guarantee that they are not cis-acting regulators of the gene. We separately compute the peak accessibility, multiply it by the logistic regression probability, and normalize it as in ABC, accounting for distance and gene density [10]. Specifically, we use  $\frac{Activity * Contact * Probability}{\sum_{|Peak-Gene| < 1\text{Mbp}} Activity * Contact * Probability}$ , using the Hi-C power of -0.87 and scale parameter of 5.41 (parameters to the ABC model) to estimate contact, with activity defined by pseudo-bulk accessibility, and the probability as mentioned above. We call links as peak-gene pairs with scores above a cutoff of 0.001, which we determined based on testing against PLAC-seq (Nott et al. 2019).

### 3.4 Results on snRNA-seq and snATAC-seq aged human brain data

We utilized 92 post-mortem human brain prefrontal cortex (PFC) samples from the Religious Order Study (ROS) and the Rush Memory and Aging Project (MAP) including AD and age-matched controls (David A. Bennett et al. 2018). We utilize single nucleus RNA-sequencing from 414k snRNA-seq cells from matched individu-

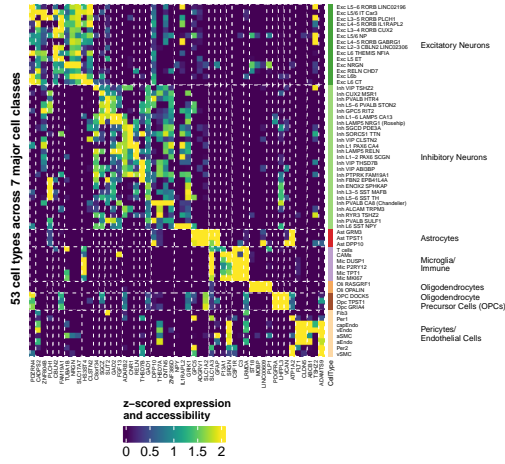
als with 7 major cell type groups and identified high-resolution cellular subtypes, including 14 excitatory, 25 inhibitory subtypes, and 22 glial and vascular subtypes.

In order to jointly resolve intra-cell heterogeneity within each major cell type and build a high-quality map of regulatory circuitry, we developed a novel cell-type-specific snRNA-snATAC integration framework that draws on distal regulatory element linking to enable high-resolution snATAC annotation. Our method iteratively refines a cell-cell graph across snRNA and snATAC by alternating between improving peak-to-gene links and using links to re-estimate gene expression from ATAC profiles in each snATAC-seq cell (Figure 2A). We integrate each major cell type separately in order to model intra-cell type heterogeneity instead of inter-cell type heterogeneity. We merged the cell-type-specific integration to build a joint cell-cell graph between modalities and transfer high resolution subtype annotations from the snRNA cells to the snATAC cells

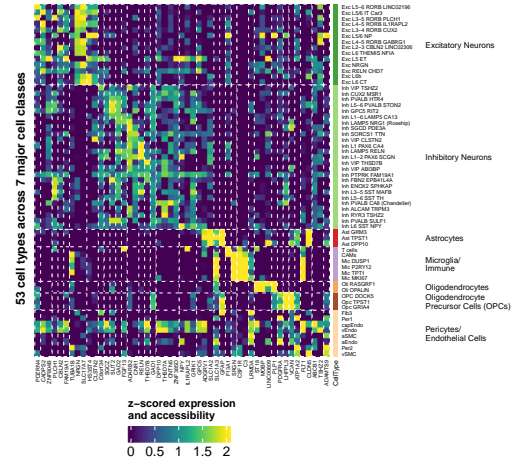
The integrated and transferred RNA and ATAC subtypes show well-aligned marker gene expression and gene-score 3-3 3-4, with much improved alignment between RNA and ATAC profiles of subtype marker genes in our integration over scGLUE [4] and Seurat [12] / Signac [19], two state-of-the-art single-cell modality integration methods based on deep learning and canonical correlation analysis, respectively.

From this integration, we calculate cell type specific peak-gene links, and evaluate against PLAC-seq (Proximity ligation-assisted ChIP-seq) data [16]. We find that area under precision-recall curves demonstrate the effectiveness of this method, as chromatin recovery of PLAC-seq interactions in four cell types compared to other existing approaches, especially in the cases of distal peak-gene interactions where a peak is not linked to the nearest gene.

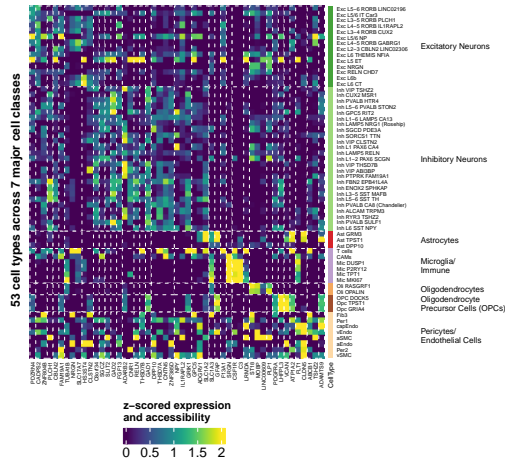
Marker expression (RNA) across subtypes in 101 genes



marker expression (ATAC) from integration across subtypes in 101 genes



er expression (ATAC) from scGLUE integration across subtypes in 101 genes



er expression (ATAC) from Seurat integration across subtypes in 101 genes

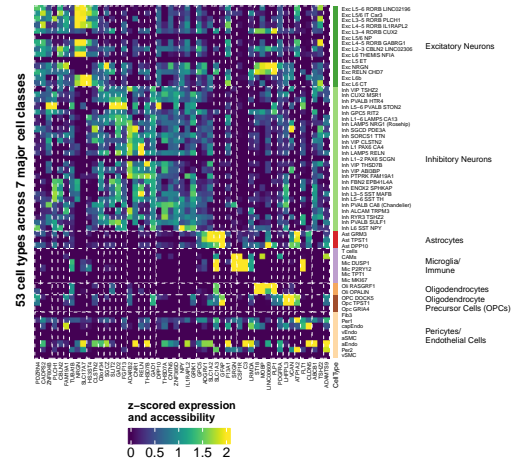


Figure 3-3: Plots of marker genes for each cell subtype, computed for each integration method (the one described), scGLUE, and Seurat/Signac

Marker gene correlation between  
ATAC and RNA across 61 cell subtypes

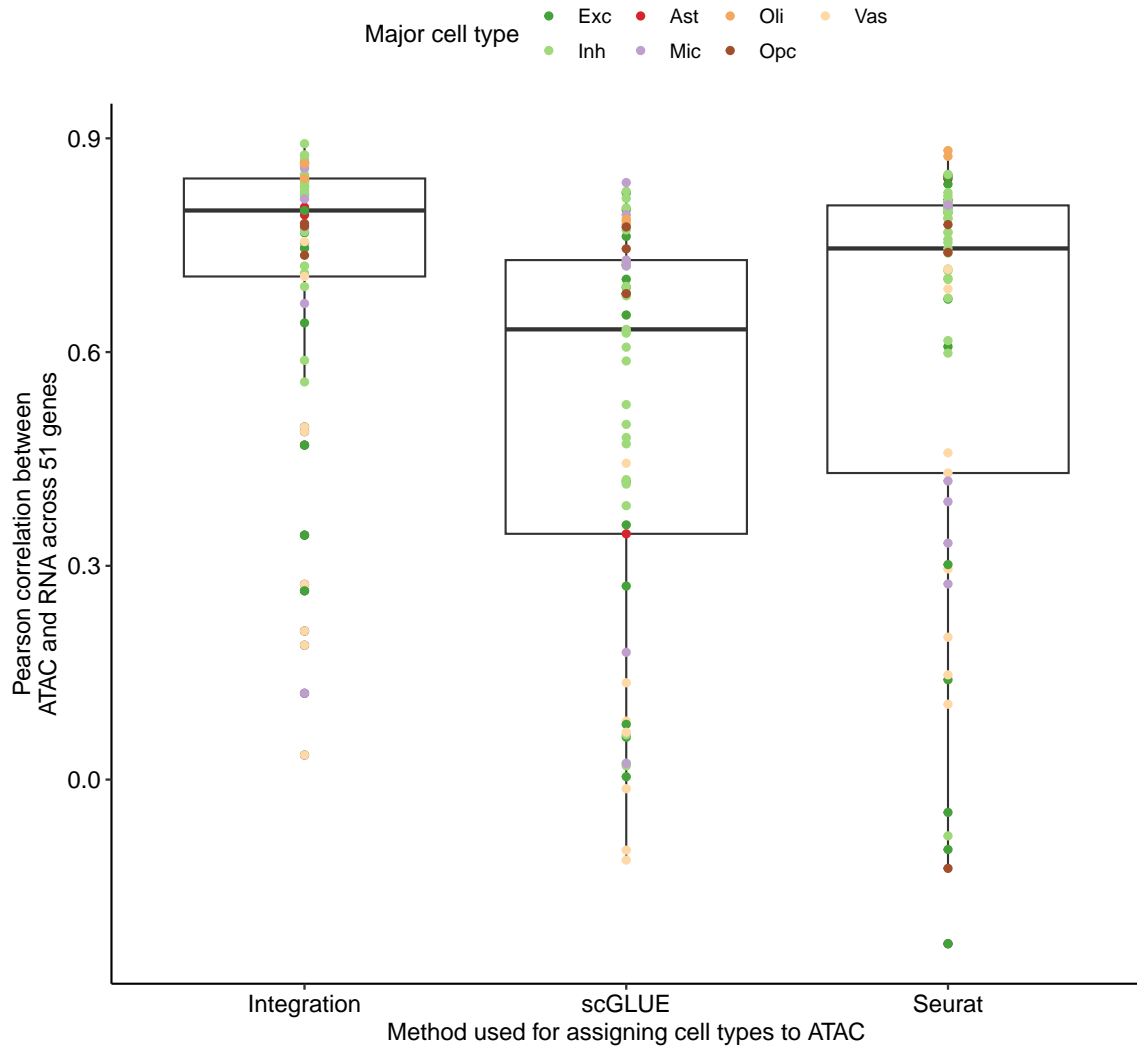


Figure 3-4: Pearson correlation coefficient across marker genes per subtype

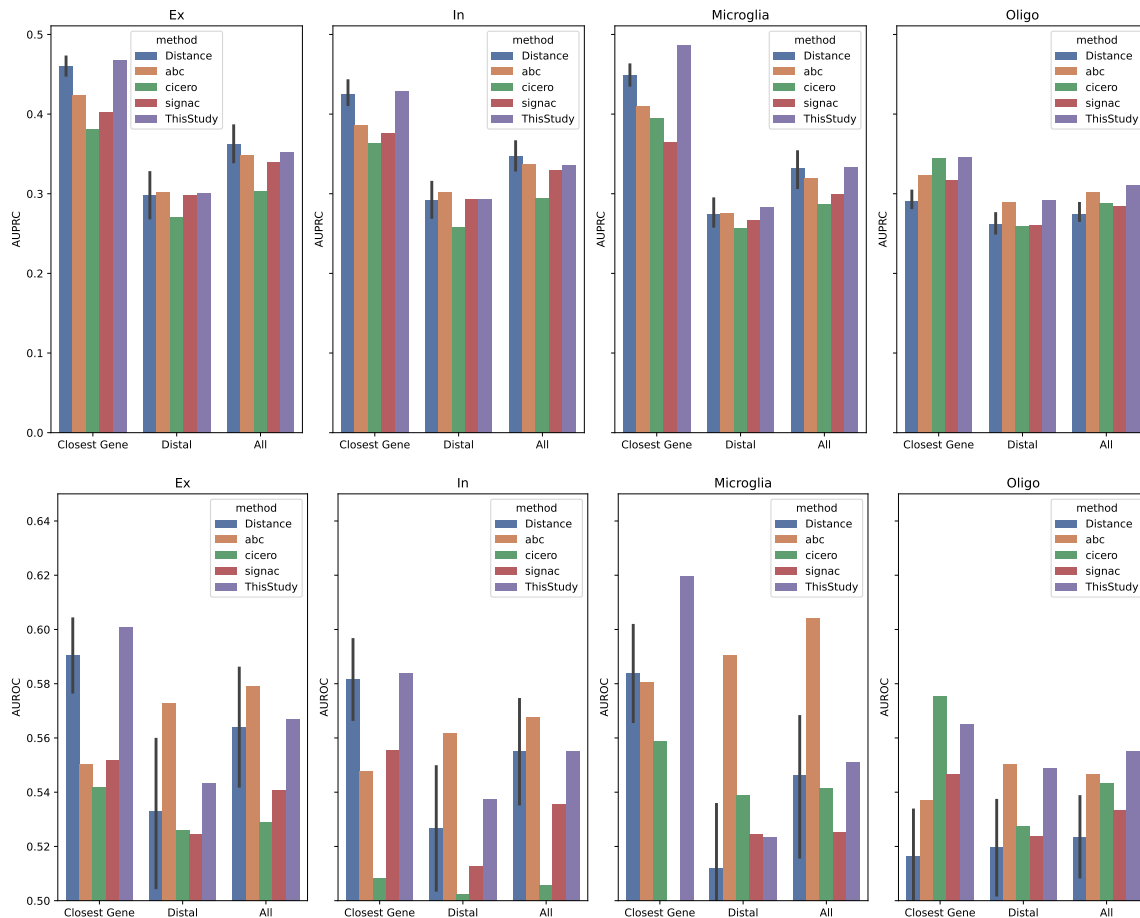


Figure 3-5: AUROC/AUPRC linking plots at the major cell type level (Ex=Excitatory neurons, In=Inhibitory neurons, Oligo=Oligodendrocytes, Microglia=Microglia and immune cells) using Nott et al PLAC-seq data. This method performs on average better than the currently used peak-gene linking methods for single cell. Notice that the scale on the AUROC plots — all the AUC values are near 0.5 due to the difficulty of this task. Additionally, note that depending on the peak set, distance from peak to gene varies in accuracy.



# Chapter 4

## Peak modules from snATAC-seq data

### 4.1 Peak modules

The peak-gene accessibility embeddings and adjustments generated for the peak to gene linking provide methodology to cluster peaks together. Moreover, we utilize a custom distance function given to the high-performance UMAP library to find 10 approximate nearest neighbors according to the z-scores. Using methodology inspired by BB-kNN (Polański et al. 2019), both top peak-peak correlations, peak-gene accessibility correlations, and gene accessibility-gene accessibility correlations are concatenated to provide a balanced nearest neighbor graph across peaks and gene accessibility. Since clustering peaks only requires high confidence correlation scores, nodes (peaks or genes) are filtered such that each embedding contains at least one correlation z-score in one of its embeddings in the 95th percentile. Thus, peaks with low connectivity do not hurt the clusters. Because these nearest neighbor graphs are calculated for each of the five co-accessibility embeddings, we utilize the multi-

plex Leiden community detection algorithm over these five graphs with a resolution parameter of 1 and a maximum community size of 2500 under a Reichardt and Bornholdts Potts model with a configuration null model. The community size restriction ensures that large modules do not dominate; however, this means that modules with similar profiles may coexist.

Additionally, we find that naively, peak modules by themselves do not aggregate into functional organization as would be expected by genomic distance. Instead, we bias these modules by incorporating a power law contact frequency normalization similar to ABC [10].

We find that since the LSI bag-of-words embedding is already unitary, we can use QR decomposition via the Gram-Schmidt method to simplify the normalization instead of re-computing SVD:

$$\begin{aligned} Q \cdot R &= \text{QR}((I + P) \cdot \text{LSI}) \quad (\text{where } P \text{ is the power-law contact sparse matrix}) \\ &= \text{QR}(\text{LSI} + P \cdot \text{LSI}) \end{aligned} \tag{4.1}$$

$$U \cdot \hat{\Lambda} \cdot V^T = \text{SVD}(R \cdot \Lambda_{\text{LSI}})$$

(Since  $R$  is upper triangular and  $\Lambda$  is a positive diagonal matrix, their product is positive)

$$E = Q \cdot U \cdot \hat{\Lambda}^p$$

The product of unitary matrices  $Q$  and  $U$  is unitary, which can be weighted ( $p$ ) by singular values  $\hat{\Lambda}$

$$\tag{4.2}$$

Thus, some weighting can be used to bias the embedding towards more cis or trans bias.

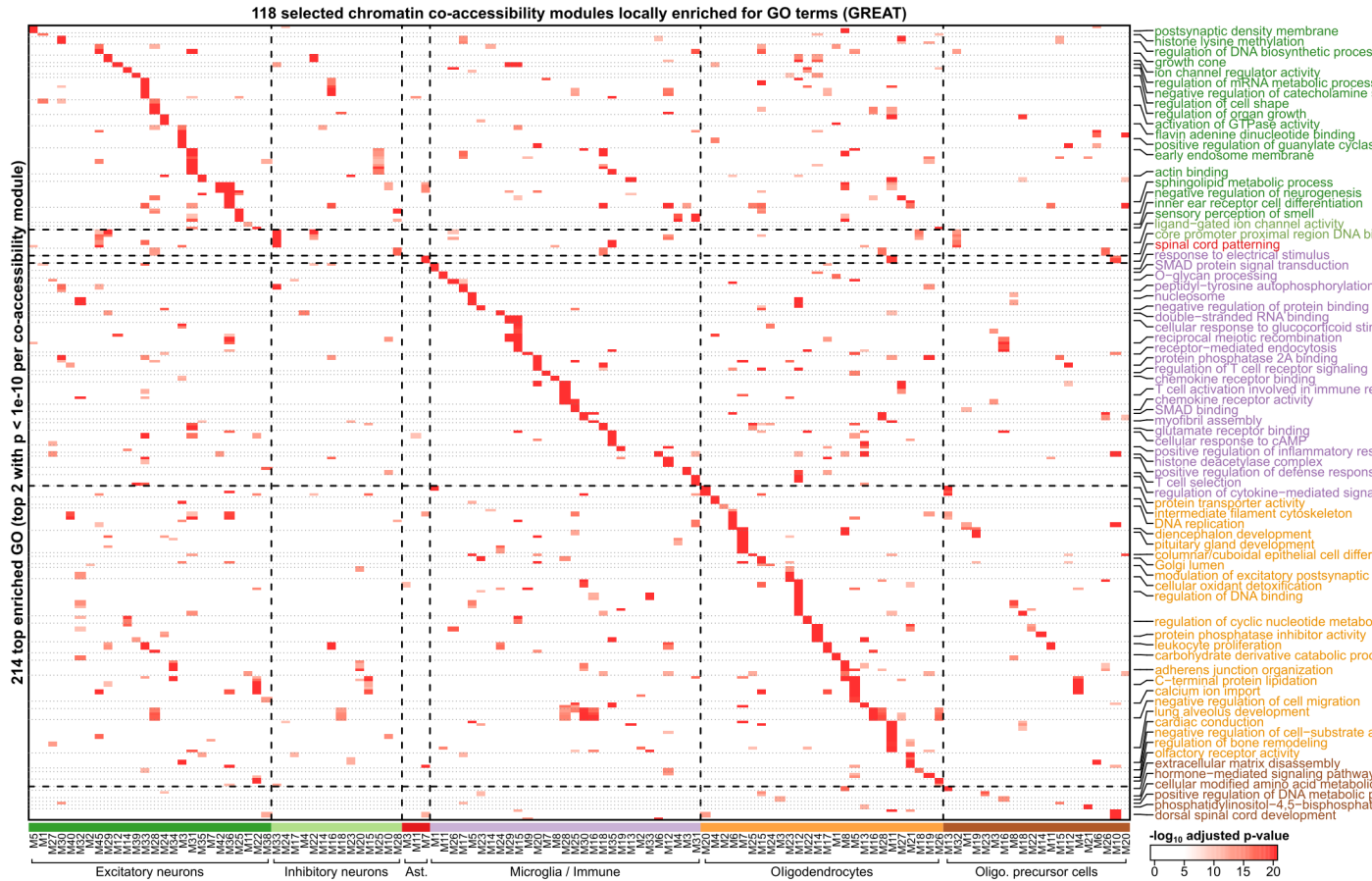


Figure 4-1: Cell type specific ATAC-seq modules capture distinct biological programs via Gene Ontology in aged human brains



# Chapter 5

## Discussion

Single cell sequencing has unlocked unprecedented search space for hypothesis generation and biological modeling.

We developed algorithms for linking peaks to genes across multiple modalities, correcting for cohort level bias, noise at the single-cell resolution. We find that using both RNA and ATAC data improves peak-gene linking tasks considerably over equivalent tasks where gene expression is estimated, even in unpaired data.

We developed an algorithm that uses this linking of peaks to genes to integrate epigenomic data and translate it into transcriptomic data, allowing for secondary analysis to be performed jointly on the meta-cells and sub-celltypes it produces.

Further, the same methodology allows for fast calculation of modules of any assay, including KNN graphs of millions of features such as in snATAC peak data. These methods allow for understanding of biological and disease mechanisms from data that could be missed by taking a pseudo-bulk approach.

The regulation of genes is still a large problem, and primary tasks including cell type identification and target gene prediction for single cell epigenomic data is still

novel. Still, experimental methods have kept on developing. The advent of single cell multiomic sequencing is removing much of the need to integrate snATAC and snRNA data, as UMIs are shared in the majority of cells between modalities. Additionally, applications of other epigenomic assays such as single cell CUT&Tag allow for finer grained regulation to be uncovered at a more specific level, but at the cost of more sparsity and less coverage.

Many of these methods will apply more generally to multiomic epigenomic sequencing technology, as the peak-gene linking, modules, and integration work on both gene and genome level information. However, much of the statistical relationships of these data and their modeling can be explored more.

However, many times integration and linking are superfluous tasks that serve no real purpose — integration of ATAC data when clusters capture distinct cell groups that correspond to functional equivalents in snRNA data does not serve much of a purpose. Linking peaks to genes in particular has been an area where lots of work has been done, but not many methods perform better than genomic distance. And at that, there isn't much of a difference at many times between genomic distance and these top-performing models. Especially in the problem space where peak, gene pairs are trained in isolation is where this problem is exacerbated. The real usage of these algorithms is diminished, as time spent optimizing AUC could be better spent understanding mechanisms.

We hope the future of this field recognizes these challenges and stuff

# Bibliography

- [1] Yael Baran et al. “MetaCell: analysis of single-cell RNA-seq data using K-nn graph partitions”. en. In: *Genome Biol.* 20.1 (Oct. 2019), p. 206.
- [2] Anthony J Bell and Terrence J Sejnowski. “The “independent components” of natural scenes are edge filters”. In: *Vision Res.* 37.23 (Dec. 1997), pp. 3327–3338.
- [3] Carles A Boix et al. “Regulatory genomic circuitry of human disease loci by integrative epigenomics”. en. In: *Nature* 590.7845 (Feb. 2021), pp. 300–307.
- [4] Zhi-Jie Cao and Ge Gao. “Multi-omics single-cell data integration and regulatory inference with graph-linked embedding”. en. In: *Nat. Biotechnol.* 40.10 (Oct. 2022), pp. 1458–1466.
- [5] Melina Claussnitzer et al. “FTO Obesity Variant Circuitry and Adipocyte Browning in Humans”. en. In: *N. Engl. J. Med.* 373.10 (Sept. 2015), pp. 895–907.
- [6] R R Coifman et al. “Geometric diffusions as a tool for harmonic analysis and structure definition of data: diffusion maps”. en. In: *Proc. Natl. Acad. Sci. U. S. A.* 102.21 (May 2005), pp. 7426–7431.
- [7] ENCODE Project Consortium et al. “Expanded encyclopaedias of DNA elements in the human and mouse genomes”. en. In: *Nature* 583.7818 (July 2020), pp. 699–710.

- [8] Jason Ernst and Manolis Kellis. “ChromHMM: automating chromatin-state discovery and characterization”. en. In: *Nat. Methods* 9.3 (Feb. 2012), pp. 215–216.
- [9] R A Fisher. “Frequency Distribution of the Values of the Correlation Coefficient in Samples from an Indefinitely Large Population”. In: *Biometrika* 10.4 (1915), pp. 507–521.
- [10] Charles P Fulco et al. “Activity-by-contact model of enhancer–promoter regulation from thousands of CRISPR perturbations”. en. In: *Nat. Genet.* 51.12 (Nov. 2019), pp. 1664–1669.
- [11] Jeffrey M Granja et al. “ArchR is a scalable software package for integrative single-cell chromatin accessibility analysis”. en. In: *Nat. Genet.* 53.3 (Mar. 2021), pp. 403–411.
- [12] Yuhan Hao et al. “Integrated analysis of multimodal single-cell data”. en. In: *Cell* 184.13 (June 2021), 3573–3587.e29.
- [13] W Evan Johnson, Cheng Li, and Ariel Rabinovic. “Adjusting batch effects in microarray expression data using empirical Bayes methods”. en. In: *Biostatistics* 8.1 (Jan. 2007), pp. 118–127.
- [14] Ilya Korsunsky et al. “Fast, sensitive and accurate integration of single-cell data with Harmony”. en. In: *Nat. Methods* 16.12 (Dec. 2019), pp. 1289–1296.
- [15] Leland McInnes, John Healy, and James Melville. “UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction”. In: (Feb. 2018). arXiv: 1802.03426 [stat.ML].
- [16] Alexi Nott et al. “Brain cell type-specific enhancer-promoter interactome maps and disease-risk association”. en. In: *Science* 366.6469 (Nov. 2019), pp. 1134–1139.
- [17] Krzysztof Polański et al. “BBKNN: fast batch alignment of single cell transcriptomes”. en. In: *Bioinformatics* 36.3 (Feb. 2020), pp. 964–965.



- [18] Roadmap Epigenomics Consortium et al. “Integrative analysis of 111 reference human epigenomes”. en. In: *Nature* 518.7539 (Feb. 2015), pp. 317–330.
- [19] Tim Stuart et al. “Single-cell chromatin state analysis with Signac”. en. In: *Nat. Methods* 18.11 (Nov. 2021), pp. 1333–1341.
- [20] Hendrik G Stunnenberg, International Human Epigenome Consortium, and Martin Hirst. “The International Human Epigenome Consortium: A Blueprint for Scientific Collaboration and Discovery”. en. In: *Cell* 167.7 (Dec. 2016), p. 1897.
- [21] V A Traag, L Waltman, and N J van Eck. “From Louvain to Leiden: guaranteeing well-connected communities”. en. In: *Sci. Rep.* 9.1 (Mar. 2019), p. 5233.
- [22] Isaac Virshup et al. “The scverse project provides a computational ecosystem for single-cell omics data analysis”. en. In: *Nat. Biotechnol.* (Apr. 2023).
- [23] F Alexander Wolf, Philipp Angerer, and Fabian J Theis. “SCANPY: large-scale single-cell gene expression data analysis”. en. In: *Genome Biol.* 19.1 (Feb. 2018), p. 15.