

StructureQTL: Novel QTL to associate SNPs and neighborhood regulatory structure

by

Aman Patel

Submitted to the Department of Electrical Engineering and Computer Science

in partial fulfillment of the requirements for the degree of

Master of Engineering in Computer Science and Molecular Biology

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

May 2020

©2020, Aman Patel. All rights reserved.

The author hereby grants to M.I.T. permission to reproduce and to distribute publicly paper and electronic copies of this thesis document in whole and in part in any medium now known or hereafter created.

Author
Department of Electrical Engineering and Computer Science
May 12, 2020

Certified by
Manolis Kellis
Professor of Computer Science
Thesis Supervisor

Accepted by
Katrina LaCurts, Chair, Master of Engineering Thesis Committee

StructureQTL: Novel QTL to associate SNPs and neighborhood regulatory structure

by

Aman Patel

Submitted to the Department of Electrical Engineering and Computer Science
on May 12, 2020, in partial fulfillment of the
requirements for the degree of
Master of Engineering in Computer Science and Molecular Biology

Abstract

Due to the effects they convey on the expression levels of certain genes, noncoding genetic regions are thought to play an integral part in the process of gene regulation and consequently the onset of several diseases. Numerous previous studies have demonstrated associations between single-nucleotide genetic changes (SNPs) and the regulatory activity of these noncoding regions. However, these studies have largely focused on single noncoding loci rather than the overall regulatory structure in a certain area, which could provide significant novel insight. We present two complementary approaches for this problem, which both involve studying the histone acetylation peak correlation matrix for a particular neighborhood. The first involves permutations and the Kolmogorov-Smirnov test, and the second relies heavily on community detection. We then demonstrate promising preliminary results on simulated and experimental data, and we identify a set of SNPs that may be attractive candidates for future study. We believe these methods, when applied on a large scale, will improve current knowledge regarding the mechanisms behind gene regulation and the causes of several important human diseases.

Thesis Supervisor: Manolis Kellis
Title: Professor of Computer Science

Acknowledgments

I am immensely grateful to Dr. Lei Hou for his constant guidance and advice throughout my MEng research as well as my preceding UROP work. I have worked closely with Lei throughout my time in the Kellis lab, and I can safely say that without his ideas and support, my project would only be a fraction of what is contained in this thesis. Lei has been the best mentor I could ask for, and I have learned so much from him in our time working together.

I would also like to thank Prof. Manolis Kellis for allowing me to join his group and work on this exciting project. I am very appreciative of his enthusiastic support for my research, both personally and financially.

Furthermore, I am very grateful to Drs. Rémy Tuyeras and Yongjin Park for expert technical advice which helped me overcome obstacles and push my project forward.

Finally, I thank my family, friends, and advisors for their constant love and support on this journey.

Contents

1	Introduction	13
1.1	Gene Regulation and Noncoding Regions	13
1.1.1	Regulatory Activity of Noncoding Regions	13
1.1.2	Detecting Active Regulatory Regions	15
1.1.3	SNPs and Regulation	17
1.2	Unsolved Problems and Potential Approaches	18
2	Data and Preliminary Methods	21
2.1	The eGTEx Dataset	21
2.2	Super-Enhancers	21
2.3	Preliminary Signal Detection	22
2.3.1	Acetylation Peak Correlation Structures	22
2.3.2	Associations with Genetic Variation	24
2.4	Limitations of Analysis	26
3	A Permutation-Based Approach	29
3.1	Explanation of Algorithm	29
3.2	Application to Super-Enhancers	30
3.3	Application to eQTLs	34
3.4	Overall Analysis of Method	37
4	A Community-Based Approach	39
4.1	A Simulation-Based Development Strategy	39

4.1.1	Simulation Process	39
4.1.2	Explanation of Algorithm	41
4.1.3	Utility of Algorithm	43
4.2	Application to Experimental Data	52
4.3	Analysis and Future Improvements	55
5	Conclusion and Future Directions	57

List of Figures

1-1	Mechanism of action of enhancers	14
1-2	Relationship between histone modifications and the activation of regulatory regions	16
2-1	Example of an acetylation peak correlation matrix	23
2-2	Examples of major and minor allele correlation matrices	25
3-1	Example SNPs detected by the permutation-based method	31
3-2	Example SNPs not detected by the permutation-based method	33
3-3	Comparison of StructureQTL and eQTL p-values for SNPs near super-enhancer regions	35
3-4	Comparison of StructureQTL and eQTL p-values at gene-level resolution	36
4-1	Examples of simulated correlation matrices	40
4-2	Community detection test results on simulated data - large differences in topology	45
4-3	Community detection test results on simulated data - small differences in topology	47
4-4	Community detection test results on simulated data - introduction of free peaks	49
4-5	Community detection test results on simulated data - single large cluster	51
4-6	Example SNPs detected by the community detection-based method	53

List of Tables

3.1	Information for selected SNPs of interest - permutation-based method	32
4.1	Information for selected SNPs of interest - community detection-based method	55

Chapter 1

Introduction

1.1 Gene Regulation and Noncoding Regions

1.1.1 Regulatory Activity of Noncoding Regions

Given that a human's genome is almost identical across various cell types, regulation of gene expression is essential to ensure that each cell only produces the proteins it needs to conduct its specialized functions. Currently, much biological research is dedicated to understanding the mechanisms by which this regulation occurs. As a result, several classes of noncoding genetic regions (regions that do not code for proteins) have been implicated in influencing expression levels of associated genes. For example, regions called enhancers and promoters both facilitate interactions between certain genes and the machinery needed to transcribe - and therefore express - them. Promoters are located directly upstream of the genes they affect, while enhancers convey their effects over a much longer range.

Figure 1-1 demonstrates the mechanisms of action of enhancers. As shown, although enhancers are located relatively far away from the gene of interest, they are able to bind transcription factors and, often aided by mediator proteins, bring these factors in closer proximity to the gene. The various transcription factors present, including those brought by the enhancer, aid the gene's promoter in recruiting the enzyme RNA polymerase, which transcribes the gene [18, 16]. Enhancer-mediated regulation

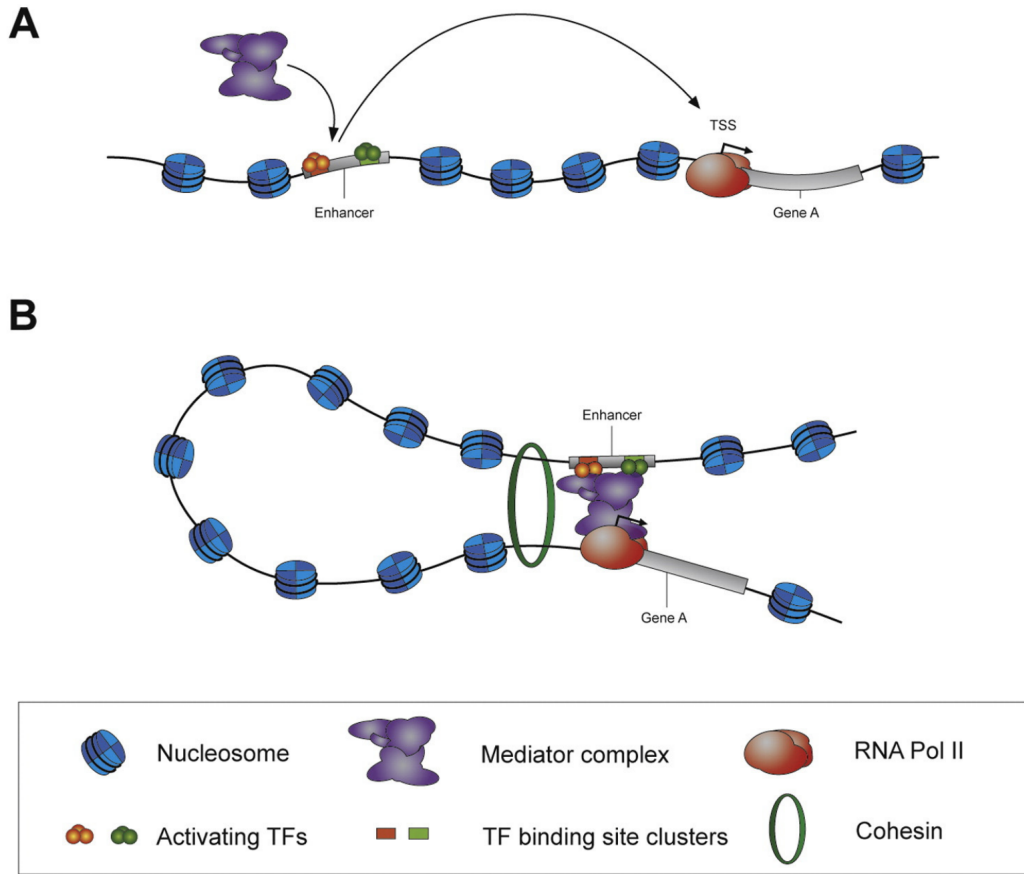


Figure 1-1: Mechanism of action of enhancers. Despite being farther away from the gene of interest, enhancers are able to bind transcription factors (TFs), and through DNA bending, bring them closer to the gene’s promoter. Figure from [16].

is vital for a variety of fundamental processes, including driving cell-type specific gene expression [11]. In short, gene regulation is mediated by a complex series of pathways, of which noncoding genetic regions are often essential elements.

It is not a surprise, therefore, that noncoding regions have been associated with a variety of complex human diseases. For example, enhancers have been connected to Type I Diabetes, Lupus, and Alzheimer’s Disease, among many other conditions [11]. Other studies have identified a role for enhancers in cardiovascular complications, particularly atrial fibrillation [15]. Finally, alterations in enhancers and promoters have been found in multiple cancers [2, 19, 25].

Thus, given the vital importance of noncoding genetic regions, it is essential to figure out as much as possible regarding the mechanisms by which they convey their effects

and, very importantly, are regulated in turn. Answering these questions can greatly augment scholarly understanding of the fundamental processes behind gene regulation as well as provide vital insights regarding the causes behind several severe human diseases.

1.1.2 Detecting Active Regulatory Regions

Since the Encyclopedia of DNA Elements (ENCODE) project, scientists have made significant progress to understand non-coding genomic regions with regulatory potential [5]. When regulatory regions are active, the DNA is open and able to be freely accessed, which means it is available to be bound by the relevant transcription machinery. In contrast, when these regions are repressive, the DNA is typically wound tightly around proteins called histones (which themselves are packaged into complexes called nucleosomes) and is therefore inaccessible to binding [20].

The process of conversion from a closed to open state typically is marked by specific modifications to the tails of histones flanking the regulatory regions in question. For active enhancers and promoters, we usually observe acetylation of the 27th Lysine residue of the H3 histone protein (denoted as H3K27ac) [20, 5]. This concept is visually illustrated in Figure 1-2 [20]. Therefore, to locate active regulatory regions throughout the genome, we simply need to search for loci at which H3K27ac signal is enriched.

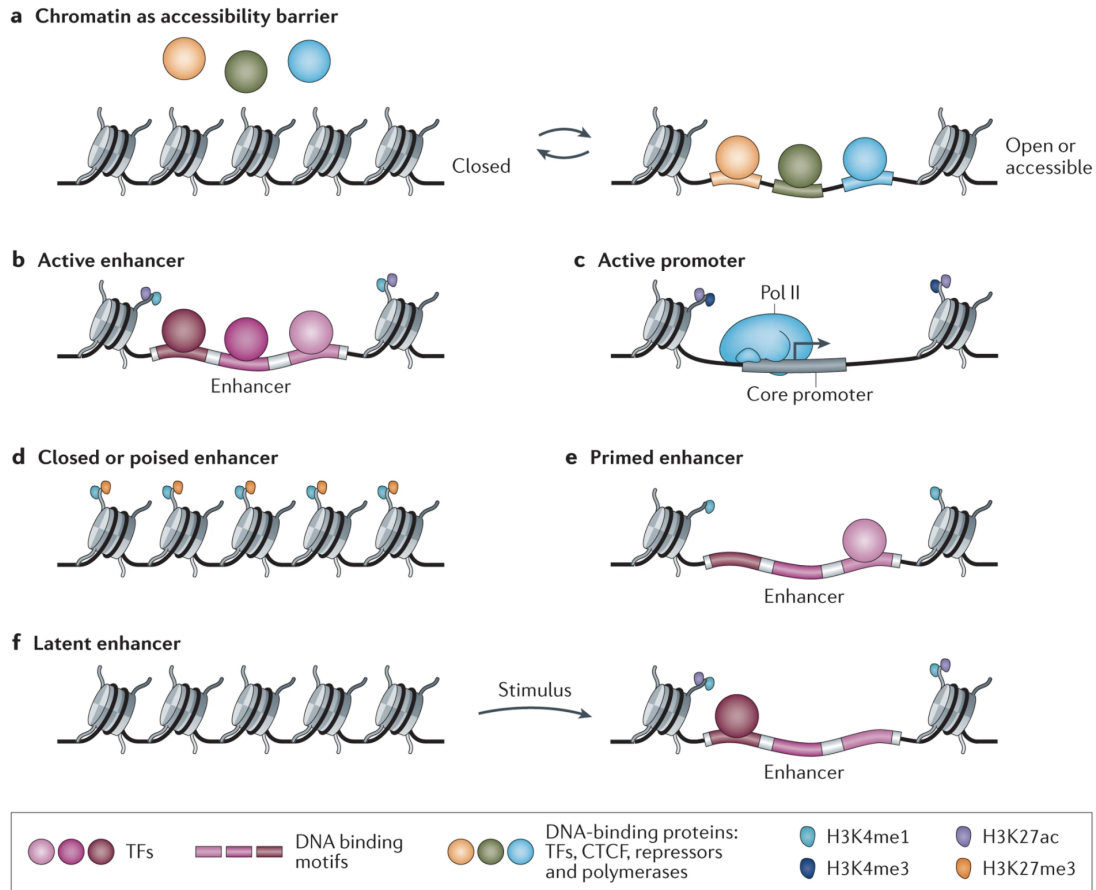


Figure 1-2: Relationship between histone modifications and DNA accessibility as well as the activation of regulatory regions. Active enhancers and promoters are usually flanked by histones with H3K27ac. Figure from [20].

In practice, this data is generated with a method called Chromatin Immunoprecipitation Sequencing (ChIP-Seq). First, molecules bound to DNA elements are crosslinked to their DNA binding sites, typically using formaldehyde. Next, cells are lysed, and the chromatin is sheared into numerous small pieces. The fragments containing molecules of interest are then collected by binding to antibodies specific to these molecules. In the case of our desired regulatory regions, we use antibodies that specifically bind to H3K27ac. Finally, the antibodies - along with the attached DNA sequences - are precipitated, allowing the DNA fragments to be sequenced [17]. Using computational alignment tools, we can map the fragments detected with ChIP-seq to their genomic locations. Using peak calling algorithms (eg. MACS), we can

determine which loci have enriched acetylation levels with respect to the background, and we assign an acetylation peak to each of these [9]. Ultimately, we define these as our active regulatory regions. The intensity of each peak is related to the frequency at which relevant region was found in the sequence data, which in turn is related to the frequency at which H3K27ac modifications are found at the locus in question [26]. Thus, the peak intensity at a particular region represents a measure for the regulatory activity in that region. This metric has become common in the field of epigenomics, and we will use it for all analysis.

In conclusion, the use of histone acetylation ChIP-seq has allowed for the capability to locate - and determine the activity of - active regulatory regions throughout the genome.

1.1.3 SNPs and Regulation

Single nucleotide polymorphisms (SNPs) refer to single-base changes in the genome which are typically found in an appreciable percentage of the population. Many SNPs have been statistically connected with expression levels of nearby genes, and many have also been associated with the presence of diseases, including diabetes, Alzheimer's, and Inflammatory Bowel Disease, among countless others. Significantly, over 90% of the thousands of disease-associated SNPs are located in non-protein coding genetic regions [10]. Furthermore, the vast majority of these SNPs are located within or nearby to enhancers [11]. Given this information, it is natural to suggest that the variants in question convey their function by somehow altering the regulatory activity of noncoding regions.

Indeed, a large number of studies have been conducted which lend credence to this hypothesis. Most notably, several genetic variants have been statistically connected to changes in histone acetylation levels at nearby loci. As stated previously, since histone acetylation is intimately associated with the presence of active genetic regulatory regions, these results serve to link genetic variation to gene regulation [22, 6].

Thus, as we seek to understand the processes governing the activity of regulatory regions, SNPs represent a very attractive target for further study. Indeed, several

unanswered questions exist regarding the specific effect conveyed by each SNP as well as the range over which this effect is manifested. Section 1.2 will describe these problems in more detail and explain the approaches this thesis will take to answer them.

1.2 Unsolved Problems and Potential Approaches

The studies described previously have effectively proven connections between SNPs and regulatory activity, but they have not exhaustively defined the full set of effects SNPs may convey. Specifically, these studies have focused on associations between SNPs and the intensity at a single histone acetylation peak. However, we can suggest another mechanism of SNP-mediated regulation. Namely, instead of affecting a specific location, SNPs can possibly affect the overall regulatory landscape in a particular region. Effects might not be statistically significant at any particular regulatory locus, but they could be when considering a larger region as a whole. This effect has already been probed with respect to different modalities, most notably regarding the effect of genotype on chromatin three-dimensional structure and subsequently gene expression [7, 24]. Thus, ample precedent exists to consider the mechanism described above to be plausible.

Because of the need to integrate data from all the acetylation peaks in a particular area, this is a much more complex problem than simply considering each locus independently. To adequately solve it, we must first develop a method to represent the regulatory structure in an area, and we must then quantify the differences in these representations between two populations (individuals with one version of a genetic variant, and individuals with the other).

A small number of previous studies have attempted to at least partially address this problem. However, all prior scholarship either answers slightly different questions or does not produce the level of statistical rigor we seek to achieve. One previous method, published by van der Wijst et al. in 2018, identifies SNPs that modulate gene co-expression, or correlations between the expression levels of two genes. This metric

indirectly involves regulatory mechanisms, and also involves neighborhood effects to a degree, but it does not achieve the resolution regarding histone acetylation that we seek to achieve [27]. The most similar study was conducted in 2019 by Delaneau et al., who presented a class of SNPs called sCRD-QTLs (structure cis-regulatory domain QTLs) as part of a larger study. They model a large region through its correlation structure, which captures the correlation between each activity peak and every other peak, and they define a sCRD-QTL as a SNP which alters this structure. Specifically, they first determine whether removal of each individual in their dataset increases or decreases the mean correlation in a particular area. Genetic variants that show an association to this phenotype are reported as sCRD-QTLs [7]. This method, while interesting, does not capture the range of associations we seek to detect. By only using the mean correlation, most of the information from the correlation structure is lost and does not factor into associations with genetic variation. In short, there could be a myriad of changes, including alterations to substructure or substantial changes to the correlations of a single peak, which do not significantly impact the mean correlation. We aim to effectively capture these differences, therefore creating a more versatile and powerful method than what is present in the current scholarship.

We therefore propose two complementary approaches to determine associations between SNPs and nearby regulatory structure. First, we present a permutation-based test which quantifies differences between peak intensity correlation matrices. Second, we present a test which models peaks as nodes of a network and evaluates differences in peak cluster assignments. We apply these methods to multiple classes of noncoding regions, and we also present a data simulation method which can be used for evaluation over a wider range of scenarios. Finally, we state the shortcomings of the methods and suggest areas for future improvement.

In conclusion, we present methodologies to identify SNPs that convey a largely unexplored type of effect on gene regulation. In keeping with previous convention, SNPs that fall under this category will be named StructureQTLs (Structure Quantitative Trait Loci). Due to the novel insights gained as compared to current variant-based analyses, we believe our methods possess great utility for future research. Overall,

we are confident that this study will add to existing knowledge regarding the ways in which the genome influences gene expression and subsequently disease.

Chapter 2

Data and Preliminary Methods

2.1 The eGTEx Dataset

All analysis for this study was conducted using the eGTEx dataset, which contains comprehensive genetic and epigenetic information for brain frontal cortex tissue from 109 individuals [21].

Through this dataset, we have access to histone acetylation levels at defined peak regions throughout the genome for various tissue types. In addition, we have information on the specific nucleotides present at SNP loci, once again over the entire genome. The dataset also provides information on which nucleotide is the "major" (more common) allele, and which one is the "minor" (less common) allele.

Thus, the dataset contains all the information necessary to allow for the identification of associations between genetic variation and the activity of regulatory regions.

2.2 Super-Enhancers

As stated previously, enhancers are essential regulatory regions involved in the promotion of gene expression. At many points in the genome, several enhancers are clustered together to form regions called super-enhancers [11]. Because of the large number of enhancers present, super-enhancers are thought to be key genomic regulatory regions. Furthermore, due to this concentration of regulatory activity, it is likely that super-

enhancer regions have numerous histone acetylation peaks, thus making them ideal for our study. Studies have also indicated that constituent enhancers within super-enhancers are often hierarchically organized, thus assigning great importance to the overall regulatory structure in a particular neighborhood [12]. In addition, it has been suggested that many enhancers which play a role in gene-expression related diseases are actually part of super-enhancers. These diseases include many that were mentioned in Section 1.1 - Alzheimer’s Disease, type I diabetes, lupus, and multiple cancers, among others. [11, 15]. Thus, any insights gained from the study of super-enhancers may be medically useful. For these reasons, we mainly focused on super-enhancer regions in our preliminary analysis.

Super-enhancers were identified using the dbSUPER database, an online catalog which provides super-enhancer genome coordinates for a variety of tissue types [14].

2.3 Preliminary Signal Detection

Before applying the algorithms which form the main contributions of this thesis, we sought to make sure that genetic variation-related changes in regulatory structure could indeed be observed in practice.

2.3.1 Acetylation Peak Correlation Structures

It was first necessary to develop a method to quantify the regulatory activity over a large region. As was done in some of the prior work explained previously, we chose to use the correlation matrix of acetylation peak intensities. That is, a certain position $[x, y]$ of this matrix would contain the correlation between the values at peak x and the values at peak y across the relevant study population.

This type of data representation gives us the advantage of being able to track associations between peaks. If two peaks are correlated very highly, it could mean that shared regulatory activity exists. If the high correlation is no longer present under a particular set of conditions (for example, the presence of a particular genetic variant), then it is possible that the shared regulation is no longer present. Thus, by utilizing

correlation matrices, we can track how each peak varies with all the others in its neighborhood. Although no data representation will be able to fully capture every aspect of such a complex feature set, this method still retains many relevant aspects. As part of our exploratory analysis, we inspected the correlation matrices at selected super-enhancer locations throughout the genome. The overall patterns reflected by the matrices varied, but a large number displayed "blocks" of high correlation, in which every peak in a block was highly correlated with every other peak. Figure 2-1 displays an example of a correlation matrix (as a heatmap) with a particularly well-defined structure.

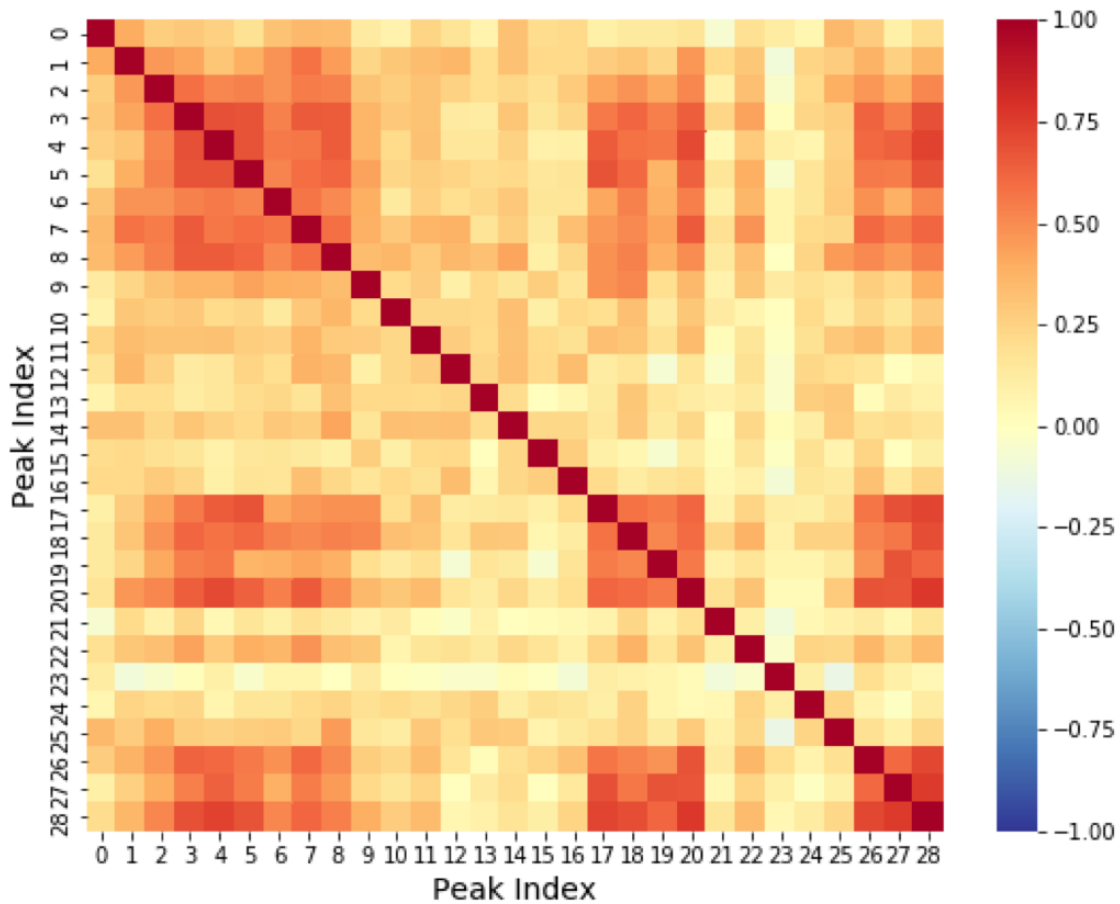


Figure 2-1: Example of an acetylation peak correlation matrix. Colors represent correlations, as depicted by the color bar. A clearly defined block structure is present.

2.3.2 Associations with Genetic Variation

We then focused on conducting rudimentary analyses to determine if associations indeed could be found and therefore if more complicated analysis would be beneficial. For each super-enhancer, we considered SNPs that were either inside the super-enhancer or no more than 50,000 base pairs from the super-enhancer's start or end points. As we were only searching for the presence of signal rather than conducting a comprehensive study, we reduced the dataset by only considering SNPs that were statistically correlated to Alzheimer's Disease with a p-value of at most 5×10^{-4} . We then applied the following procedure for each SNP that passed our filter:

1. Divide the study population into individuals containing two copies of the major (most common) allele and those containing at least one copy of the minor (less common allele).
2. For each subpopulation, calculate the correlation matrix for the peaks of the relevant super-enhancer.

Even without any quantitative testing, it became clear that differences often existed in the correlation matrices for the two subpopulations. Figure 2-2 provides three examples. In the first example, the areas of high correlation in the minor allele matrix are much less well-defined than in the major allele matrix. Furthermore, the areas surrounding islands of high correlation display much lower correlations. In the second example, blocks of high correlation are far more intense in the major allele matrix as opposed to the minor allele matrix. In the third example, the opposite effect is observed. Thus, in all three examples, we see distinct and easily visible differences between the two subpopulations.

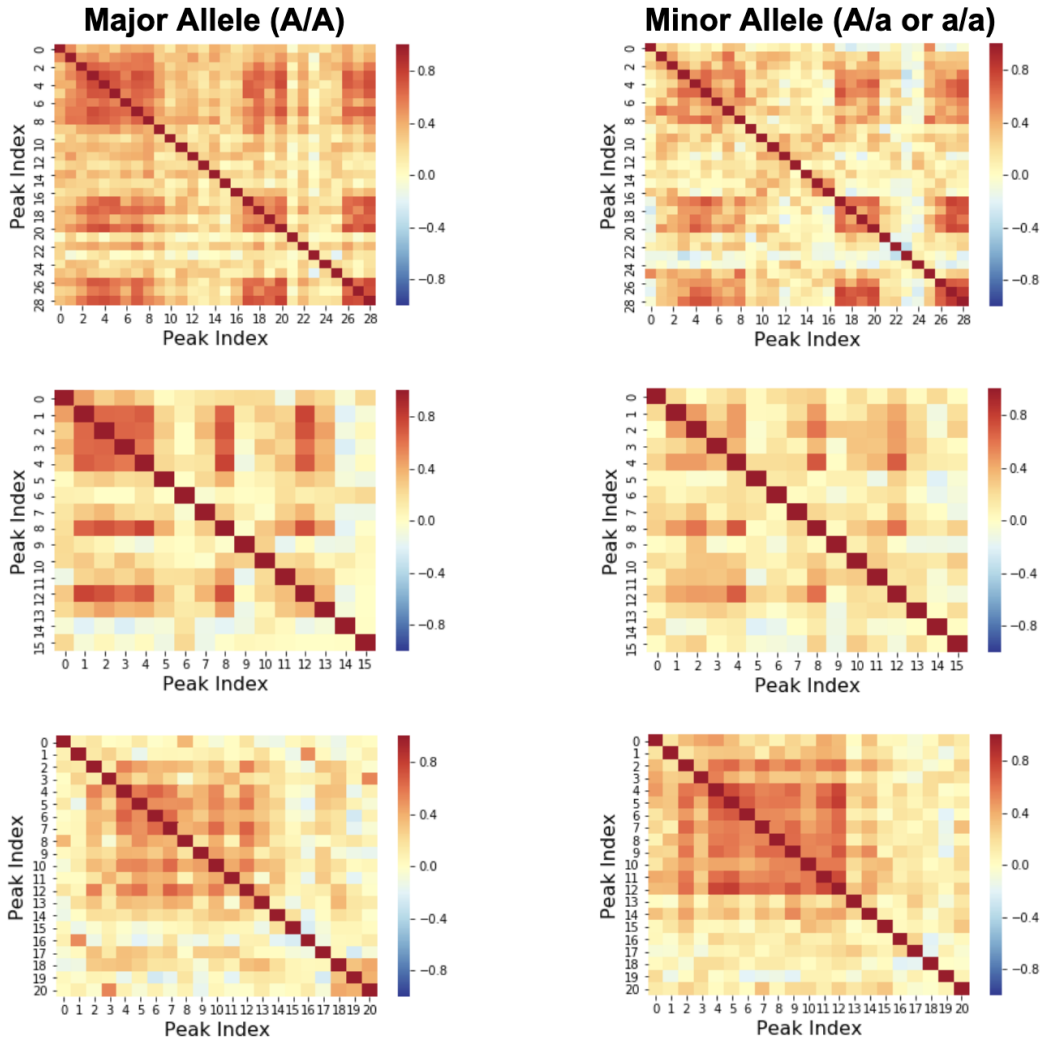


Figure 2-2: Examples of major and minor allele correlation matrices for three SNPs located nearby to super-enhancers. Major allele matrices are on the left, and minor allele matrices are on the right.

We next tested whether these matrix differences could be detected through quantitative statistical testing. For this experiment, we focused on 36 SNPs on Chromosome 2 (all near one super-enhancer) with especially strong association to Alzheimer’s ($p < 5 \times 10^{-8}$). For each SNP, we calculated the major and minor allele matrices, and we then conducted paired t-tests on corresponding pairs of correlation vectors. To clarify, a correlation vector is defined as the vector containing the correlations of one

peak with every other peak. These vectors are the rows/columns of the correlation matrix. Tests were performed with a multiple hypothesis correction, resulting in a significance threshold of 0.05 divided by the number of peaks within the super-enhancer. Each of the 36 SNPs returned multiple significantly changing correlation vectors at this threshold, with some SNPs displaying significance for as many as 19 of the 29 vectors. This provides a degree of quantifiable proof of the changes in regulatory architecture that were visually observed using Figure 2-2.

Thus, we have quantitative and visual evidence of differences in correlation structure between major and minor allele subpopulations. These differences were found at several different locations in the genome (the three examples in Figure 2-2 are all from different chromosomes), indicating this observation is not location-specific. In conclusion, the the presence of SNP-associated changes to regulatory architecture is clear.

2.4 Limitations of Analysis

The above methodology is notable for its simplicity, and it allowed for an initial determination of the presence of SNP-associated changes to regulatory structure. This, in turn, allowed us to conclude that further, more comprehensive analysis is feasible. However, it is important to note that this method lacks the resolution and versatility to represent an effective approach for a more detailed study.

The most prominent issue is that a t-test based approach only considers one correlation vector at a time. This is effective in determining how one peak changes with respect to the others, but it is not guaranteed to detect other types of changes. For example, there could be changes to the block structure of the correlation matrix, in which certain blocks expand, shrink, or change in intensity. These alterations are certainly of biological interest, yet a method which considers each vector separately - and considers relationships between a peak and the entire set of other peaks - would not be optimal to find them. In addition, changes that are only significant on a more global scale may not be detectable when only considering one correlation vector.

In short, the ideal method would detect changes over the entire correlation matrix while still retaining the ability to identify a wide variety of changes that may occur over subsets of the matrix. In Chapters 3 and 4, we develop and evaluate two approaches to solve this problem.

Chapter 3

A Permutation-Based Approach

We first present an approach which relies on comparison of biologically relevant subpopulations to those created by random permutations.

3.1 Explanation of Algorithm

After an iterative design process, the most current form of our permutation-based algorithm is as follows:

1. For a particular SNP, divide the study population into individuals containing the major allele and those containing at least one copy of the minor allele. Let p_{maj} and p_{min} be the two subpopulation sizes.
2. For each subpopulation, calculate the correlation matrix. Then, subtract the minor allele correlation matrix from the major allele correlation matrix. Convert this matrix to a vector of values, called d_{true} .
3. For a certain large value n make n random divisions of the study population into subpopulations of size p_{maj} and p_{min} . For each corresponding pair of subpopulations, calculate the difference in correlation structure as described previously, and convert to a vector. Combine the vectors from all n permutations together into one vector, called d_{perm} .

4. Run a Kolmogorov-Smirnov (KS) test to determine whether the distribution of values in d_{true} matches the distribution found in d_{perm} . If the two distributions are found to be significantly different (after application of the multiple hypothesis correction), then we can conclude that the SNP in question is indeed associated with alterations in super-enhancer regulatory structure.

By artificially generating correlation structures through random permutations, we are able to produce a background distribution for the differences in correlations between two subpopulations. By comparing this background to the differences produced by biologically relevant subpopulations, we can determine if the SNP in question is associated with changes to regulatory structure that would not be expected to arise at random.

We selected the KS test because it is the canonical test to determine if two probability distributions - for example, the biologically relevant correlation matrix difference values and the background values - are significantly different.

We next explored the utility of the algorithm on two different classes of noncoding regions, the results for which are detailed below. Our goal for this evaluation was to gauge the method's potential on experimental data and detail the scenarios over which it might prove useful.

3.2 Application to Super-Enhancers

The method was tested on the super-enhancer dataset, which was once again filtered to only include super-enhancers near SNPs with mild statistical associations to Alzheimer's Disease ($p < 5 \times 10^{-4}$).

As this is a novel type of analysis, we were hampered by the fact that no ground truth labels existed to indicate whether the SNPs in our dataset were actually associated with changes in regulatory structure. Therefore, we chose to study the method by generally categorizing the types of correlation matrix differences over which the method did and did not detect strong signal.

Overall, results were promising yet indicated significant need for further testing and

characterization. In certain instances in which global changes to correlation matrix block structure were observed, the method produced small p-values. Figure 3-1 details two examples of such changes, along with the corresponding p-values. In the first, the correlations of the large block are greatly increased in the minor allele matrix, which also contains a small block not found in the major allele matrix. In the second, we see changes in structure and intensity in virtually every block present. Furthermore, the background regions in the minor allele matrix have far more negative correlations.

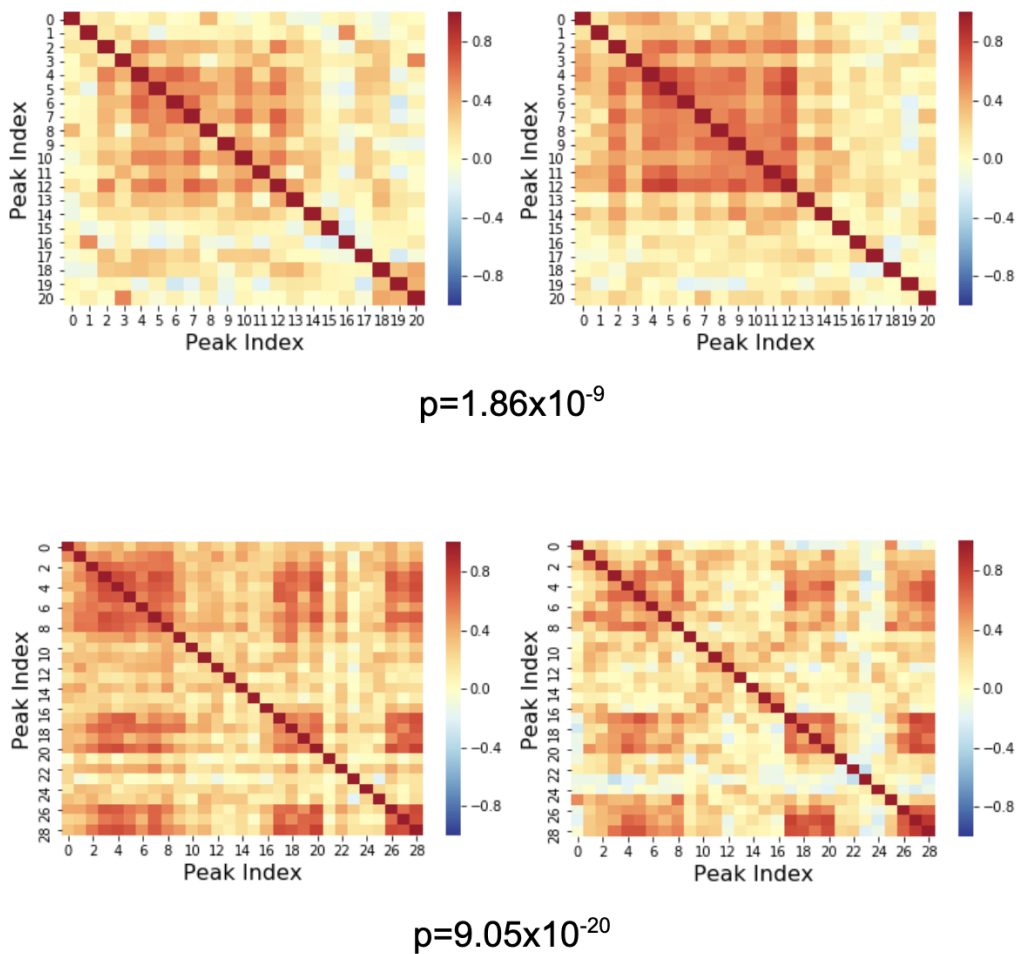


Figure 3-1: Examples of SNPs detected by the permutation-based method described in this chapter. Major allele correlation matrices are on the left, and minor allele correlation matrices are on the right.

For a subset of the SNPs we detected, we utilized the UCSC Genome Browser to locate protein-coding genes located within 50,000 base pairs from the corresponding super-enhancers [13]. A list of these genes is provided in Table 3.1, and these genes are certainly attractive candidates for further study. Specifically, several SNPs on chromosome 2 (of which two are listed in the table as examples) show strong StructureQTL signal coupled with low GWAS p-values for Alzheimer’s disease, thus indicating the possibility that alterations in regulatory structure at this locus are contributors to the presence of the disease. Indeed, the gene BIN1, which is located nearby, has been shown to have importance for the onset of the disease. [23]. This result by no means confirms the presence of any such effect, but it does identify a highly promising avenue for additional investigation.

Chr	SE Start	SE End	SNP Location	StructureQTL p-value	AD GWAS p-value	Nearby Genes
2	127050335	127153940	127136908	2.070×10^{-26}	5.126×10^{-17}	BIN1 CYP27C1
			127135380	1.396×10^{-25}	1.543×10^{-7}	
10	22599170	22656242	22558858	2.70×10^{-25}	3.391×10^{-4}	PIP42KA

Table 3.1: Information for selected SNPs that displayed strong signal when the permutation-based method was applied. Columns depict the SNP’s chromosome, the start and end points of the super-enhancer the SNP is located nearby to, the location of the SNP itself, the SNP’s StructureQTL and Alzheimer’s Disease GWAS p-values, and protein-coding genes found within 50kb of the associated super-enhancer.

However, there were situations in which the method did not detect matrix changes, two of which are detailed in Figure 3-2. In these scenarios, the method produced p-values that would likely not be considered significant after any kind of multiple hypothesis correction.

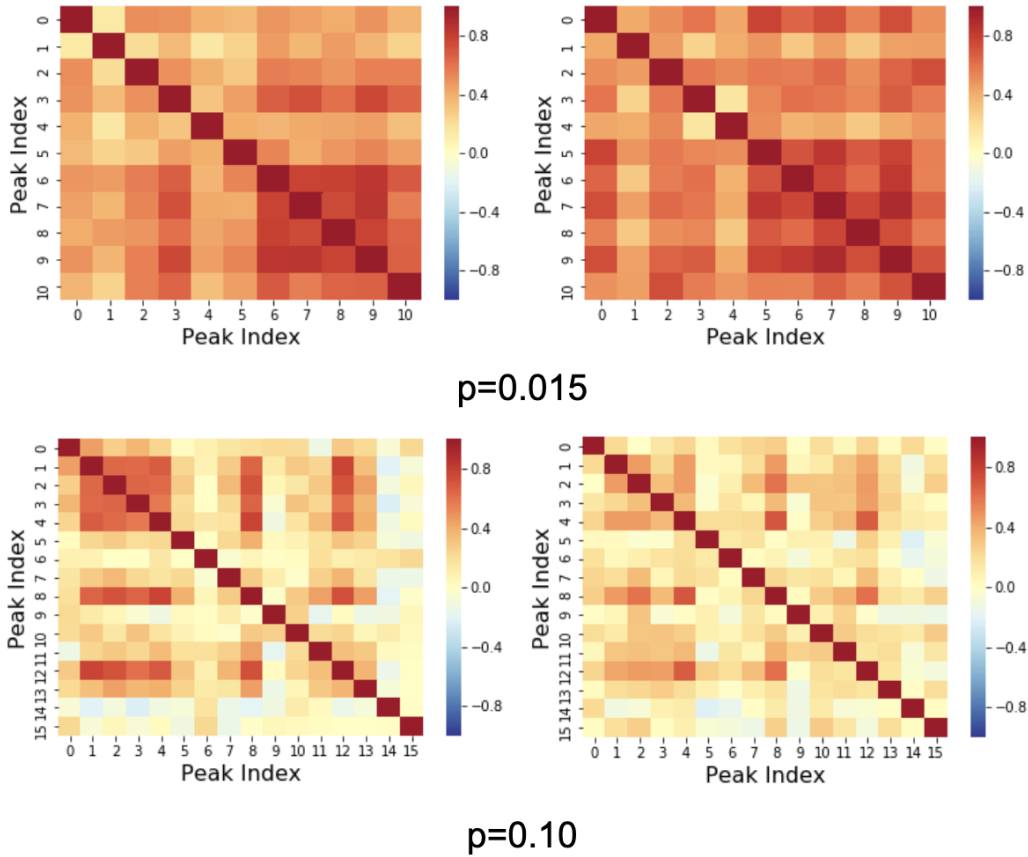


Figure 3-2: Examples of SNPs not detected by the permutation-based method even though clear differences exist between the major and minor allele correlation matrices. Major allele correlation matrices are on the left, and minor allele correlation matrices are on the right.

The differences between scenarios that produced small p-values and those that did not were not always clear. For example, the changes in the second example in Figure 3-2 seem very similar to those found in Figure 3-1 (we provide possible explanations for this in Section 3.4). However, we can note that the method was generally ineffective at detecting changes that were more local in nature, including small changes in cluster membership or local changes in cluster intensity. This is evidenced in the first example, in which the large cluster contains more members in the minor allele matrix than it does in the major allele matrix.

In short, although significant ambiguity exists, we were able to gain a degree of insight into the actions of the method for various types of changes in correlation structure. Notably, it is clear that in a subset of situations in which block structure changes are observed, the method is indeed able to detect associations.

3.3 Application to eQTLs

As a complementary analysis, the method was also applied to regions surrounding expression quantitative trait loci (eQTLs), which are defined as SNPs that have been statistically related to the expression levels of nearby genes. It is important to note that while these statistical associations have been uncovered, the mechanisms by which eQTLs actually influence gene expression have not always been elucidated.

Since histone acetylation serves as a marker for active regulation, it is possible that StructureQTL-associated changes in acetylation correlation structure could be associated with altered expression of nearby genes. Thus, we compared eQTL and StructureQTL p-values for a variety of SNPs to determine how the signals for both metrics varied with each other and if any commonalities existed.

We first conducted a high-level analysis with the same SNP cohort as in Section 3.2 - those near super-enhancers and with an Alzheimer's Disease GWAS p-value of at most 5×10^{-4} . Figure 3-3 displays the negative logarithm of the StructureQTL p-value plotted against the negative logarithm of the eQTL p-value for each of these SNPs. Note that since each SNP has a different eQTL p-value for each gene, the eQTL p-value plotted for each SNP is the minimum across all genes. As indicated in the figure, no consistent correlation exists between the two p-value metrics. We can therefore conclude that for this dataset, StructureQTL signal is largely independent of eQTL signal. That is, SNPs with the lowest StructureQTL p-values do not seem to show correspondingly low eQTL p-values.

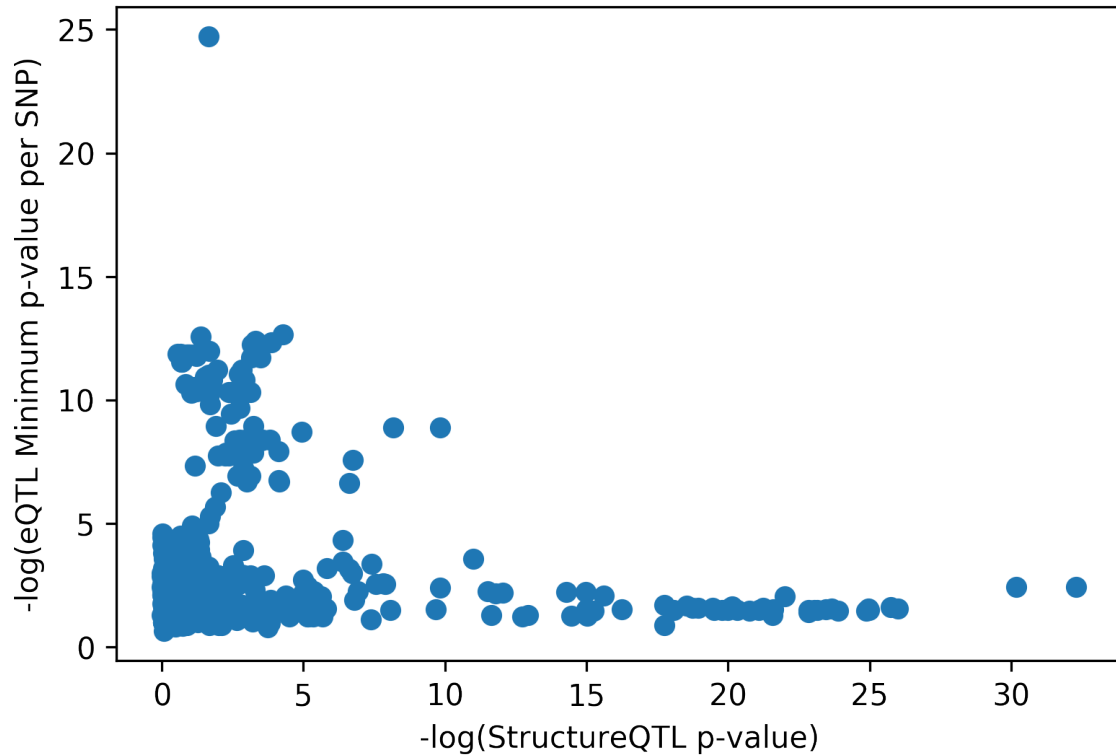


Figure 3-3: Comparison of StructureQTL and eQTL p-values for each SNP in the super-enhancer dataset. Values on the x-axis represent the negative logarithm of each StructureQTL p-value, and values on the y-axis represent the negative logarithm of the minimum eQTL value for each SNP across all possible genes.

We next increased the resolution of the analysis by considering individual genes on chromosome 19. We sought to study eQTL-rich regions, so we only considered the genes with the top 10 q-values for associations with eQTLs. For each of the genes, we conducted the following procedure:

1. Locate all acetylation peaks within a distance of 1,000,000 base pairs from the gene, in either direction. The correlation matrices will be formed from these peaks.
2. Collect the eQTL p-values for SNPs within a distance of 100,000 base pairs from the gene. Filter these SNPs to only those that have a p-value less than 0.01 (for runtime purposes).

3. Calculate the StructureQTL p-values for each SNP.

The genes depicted in Figure 3-4 are generally representative of all of genes studied in that once again, no consistent relation exists between StructureQTL signal and eQTL signal.

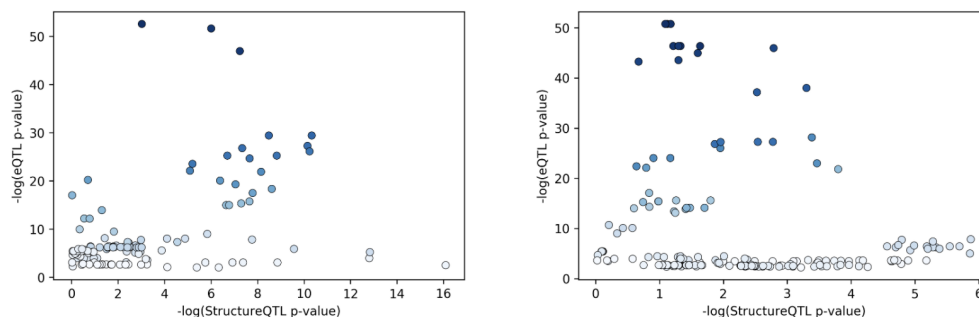


Figure 3-4: Comparison of StructureQTL and eQTL p-values for the genes LINCO1535 (left) and ZNF749 (right). Values on the x-axis represent the negative logarithm of each StructureQTL p-value, and values on the y-axis represent the negative logarithm of each eQTL p-value. Color represents R^2 with respect to the SNP with the lowest eQTL p-value.

However, despite this lack of relation, evident in both plots are cohorts of SNPs with elevated StructureQTL signal along with elevated eQTL signal. For example, in the plot for gene LINCO1535 (on the left), we see a distinct cluster with x values between 6 and 10 and y values between 10 and 30. As the p-values indicate, it is indeed possible that these SNPs may affect gene expression through alterations in regulatory structure. However, another explanation is also possible. For each plot, the color of each data point represents the R^2 value of the SNP in question with respect to the SNP with the lowest eQTL p-value. The R^2 is a measure of how well the allele at one locus is associated with the allele at another locus, and it can be thought of as analogous to the classical linear regression R^2 value. In practice, a high R^2 means that the two SNPs being compared are likely located in close proximity to each other [28]. In the plots, we see that SNPs in the aforementioned cohorts have higher R^2 values than those with equally strong StructureQTL signal but weaker

eQTL signal. Thus, a highly plausible explanation is that instead of conveying a major effect themselves, the cohort SNPs are partially linked by distance to SNPs with high eQTL activity, and as such, their own signal is elevated. Further studies are necessary to determine which explanation is ultimately true, and as such, analyzing these cohorts further represents an interesting avenue for future research.

Overall, the evidence does not suggest a consistent relation between StructureQTL and eQTL signal. In fact, we largely conclude the opposite. This lack of relation is a very interesting finding, and it is certainly worth investigating in more detail.

3.4 Overall Analysis of Method

Overall, this permutation-based algorithm produced results that are promising if not completely sufficient. In a subset of situations in which substantial alterations to correlation structure were present, the method did indeed detect these changes. Thus, we have some evidence that the method can detect the type of signal we seek to capture. Furthermore, the method produced interesting results regarding the relationship between StructureQTL and eQTL p-values.

However, it is important to acknowledge the method’s shortcomings, both in theory and in practice. Even in our inexhaustive analysis in Section 3.2, we were able to identify cases for which the method was not effective in identifying the changes to correlation structure that we intended it to capture. When we examine the Kolmogorov-Smirnov test more closely, we can potentially identify why. The KS test works by computing the differences in the Cumulative Distribution Function (CDF) of each sample at various quantiles, taking the largest difference, and using this value to compute a p-value [29]. Thus, the KS test is a holistic test, for which the entire distribution may be relevant. For example, even if certain blocks of correlation may differ, the overall matrices may be similar enough that the KS test fails to classify the distributions as significantly different. Using similar reasoning, we can also theorize cases (which were not seen in practice) in which a particularly noisy background leads to a significant result even if no real structural changes exist. Furthermore, since

the KS test only considers one difference in the CDF, the order of magnitude of the p-value produced may not necessarily reflect the overall amount of change between the two correlation matrices. This in turn would make comparisons of p-values from different SNPs very difficult. Finally, by placing subtracted correlations into a single vector, we lose all spatial data, which could certainly be relevant. Thus, as a complement to this permutation-based method, we need an approach which is optimized to reliably detect the specific types of changes we see in practice. In Chapter 4, we propose approaches that can more effectively detect changes in block membership independently of the overall composition of the relevant correlation matrix.

Chapter 4

A Community-Based Approach

We present an alternative approach to the methodology described in Chapter 3. This approach is designed to effectively and reliably detect differences in block structure as a complement to the more global approach of the KS Test-based method.

4.1 A Simulation-Based Development Strategy

We chose to develop and test this methodology on simulated data before switching to experimental data. Through simulations, we are easily able to generate a wide variety of correlation matrices, thus allowing us to readily test our methods on a comprehensive list of cases. This, in turn, will allow us to more comprehensively determine the range of scenarios over which the methods will succeed or fail.

4.1.1 Simulation Process

We simulated peak data as follows:

1. We first determine the topologies of the correlation matrices for the major and minor allele populations. That is, we assign each peak to a block or denote it as a free peak (not in any block).
2. We will now generate the peak covariance matrix. Each covariance will be generated from a normal distribution. If two peaks are part of the same cluster, this

distribution will have mean μ_{match} and variance σ_{match}^2 . If not, the distribution will have mean μ_{diff} and variance σ_{diff}^2 , where $\mu_{match} > \mu_{diff}$. Covariances will be capped so their absolute values will not exceed the variances of the individual peaks, which are artificially set as a constant value.

3. The peaks will then be generated from a multivariate normal distribution, with constant mean and the above covariance matrix. If we desire a simulated dataset with a sample size of m individuals, then we will generate m peak vectors from this distribution.
4. The peak correlation matrix can then be generated as normal.

Of course, it is unlikely that the distribution of peak intensities generated using this method matches the true distributions observed in experimental data. However, as all subsequent analysis deals with the peak correlation matrix rather than the peaks themselves, any discrepancies are far less important. Figure 4-1 displays some examples of simulated correlation matrices. As can be seen, the block structures the matrices were designated to recreate are clearly evident.

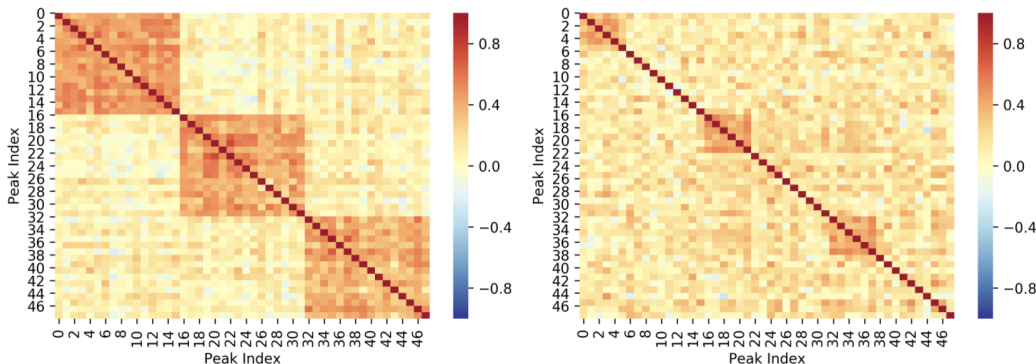


Figure 4-1: Examples of simulated correlation matrices. The matrix on the left was generated from a topology consisting of three clusters of size 16. The matrix on the right was generated from a topology consisting of three clusters of size 6, with 10 peaks between each one. The matrices recreate the topologies faithfully.

4.1.2 Explanation of Algorithm

We now describe the approach to detect differences in correlation structure. The procedure assumes simulated correlation matrices have been generated for major and minor allele populations, each with a predefined correlation matrix topology.

1. For both the major and minor allele correlation matrices, we run a clustering or community detection procedure which separates peaks into blocks (or classifies them as free peaks that do not belong to a block). Options for this procedure are described later.
2. We calculate the normalized mutual information (or a different metric measuring cluster similarity) between the two sets of community assignments.
3. For a certain large number n , we generate n simulated datasets with the same topology as the original dataset (major or minor) with the higher sample size (ie. larger number of individuals). We then run our community detection procedure on the correlation matrices of each of these background datasets. Finally, we compute the mutual information of each background community assignment with the communities of the original dataset with the same topology. That is, if the background data was produced with the major allele topology, then the communities of each background dataset will be compared to the communities of the original major allele data. (Note: if the sizes of the major and minor allele populations were equal, as was the case in all the simulations we ran, then we defaulted to using the major allele matrix for this step.)
4. We determine the fraction of these n mutual information scores that are lower than or equal to the score calculated between the major and minor allele clusterings. This value is then returned as the p-value.

Essentially, this approach tests the null hypothesis that the variations between the major and minor allele correlation matrices are *not* due to differences in topology. If (for example) the mutual information between the major and minor matrices is consistently lower than the mutual information between the major matrix and another

with the same topology, then we can likely conclude otherwise.

The procedure above is presented as a framework rather than a completely defined method since it is possible to use several different approaches at certain steps in the process. Most notably, we can apply a variety of strategies to conduct clustering or community detection on correlation matrices. Many of these are motivated by modeling the matrix as a network and instituting certain rules to preserve edges and then conduct clustering. Certain notable strategies are explained below. All of these group peaks based on the values in the correlation matrix, but they prioritize slightly different metrics.

- Conduct k -means clustering on the correlation matrix. In practice, an effective value for k is the number of blocks in the topology, plus one if any free peaks exist. For application to experimental data, an alternative method to determine k is needed as the topology is not known ahead of time.
- Set a correlation cutoff ρ , and initialize a network with peaks as nodes and with no edges initially. Create (unweighted) edges between two peaks if their correlation is greater than ρ . Finally, apply a community detection algorithm to partition the nodes into communities.
- Construct a k -nearest neighbors graph by connecting each peak to the k peaks it has the highest correlations with. We included edges only if both nodes are in each other's set of nearest neighbors, but this is not strictly necessary. Next, conduct community detection to partition the peaks.
- For each peak, calculate the Euclidean distance between its correlation vector and the correlation vector of each other peak. Next, construct a k -nearest neighbors graph by connecting each peak to the k peaks it has the lowest distance to. Once again, we included edges only if both nodes are in each other's set of nearest neighbors. Finally, conduct community detection to partition the peaks.

In all our experiments involving community detection, the Louvain algorithm was

our method of choice [3, 1]. Additionally, whenever a k -nearest neighbors graph was required, we set k equal to the nearest integer to $1 + \sqrt{a}$, where a is the total number of peaks.

It is also important to note that the process of generating n background peak datasets is designed to maximally capture the variability that can be found in correlation matrices generated from the same topology. However, as we do not know the topology ahead of time, this approach is not feasible when applying the method to real data. We therefore tested an alternate method, in which the background data is sampled (with replacement) from the original peak data. Specifically, to form one background dataset, n peak vectors are sampled from the dataset (major or minor) with the higher sample size, with n equal to the lower sample size. Although this approach is not as optimal, since the background data is not completely independent of the original, it is still a way of taking potential noise and variability into account.

4.1.3 Utility of Algorithm

We evaluated our methodology on a variety of cases to determine its ability to detect changes in correlation structure. In all of the below experiments, we created simulated regions with a total of 48 peaks, and we assigned a sample size of 100 to the major and minor allele populations. For two peaks of the same cluster, the the normal distribution to determine covariance was assigned a mean of 0.75 and a standard deviation of 0.1, and the distribution for two peaks of different clusters had a mean of 0.2 and a standard deviation of 0.2. A total of 100 background datasets were generated per test. For these experiments, we used the euclidean distance-based community detection method.

Using the above parameters, our evaluation procedure given specified major and minor allele topologies was as follows:

1. Conduct 1,000 tests comparing datasets with the major and minor allele topologies. These tests should use the approach in which background data consists of independently generated datasets. Calculate the p-value distribution, and

determine the fraction of p-values which fall below 0.05. This value can be interpreted as the overall power (at 95% confidence) over a large number of tests.

2. Conduct 1,000 tests comparing two datasets with the major allele topology to each other. Once again, background data should be independently generated. Calculate the p-value distribution, and determine the fraction of p-values which fall below 0.05. This value can be interpreted as the overall false positive rate (Type I error rate) over a large number of samples.
3. Repeat steps 1 and 2, except for these tests, create background data by sampling from the original data (as previously described) rather than generating completely new datasets.

Through these metrics, we are able to gauge the method's sensitivity as well as its propensity for false positives, which are both important factors when evaluating overall performance. Furthermore, by evaluating our alternative approach to produce background data, we gain a more accurate idea of how the method will translate to experimental data. A subset of cases the method was tested on are detailed below. Each figure displays sample major and minor allele correlation matrices as well as the p-value distributions obtained from each type of test studied.

- **Scenario 1:** We first tested our methodology on the scenario in which the major and minor allele matrices have vastly different correlation structures. Specifically, we set the major allele matrix to have three equally sized clusters of size 16, and we set the minor allele matrix to have two equally sized clusters of size 24. This scenario is unlikely to occur in real data, but we saw it as a preliminary benchmark to clear before moving on to more complex cases.

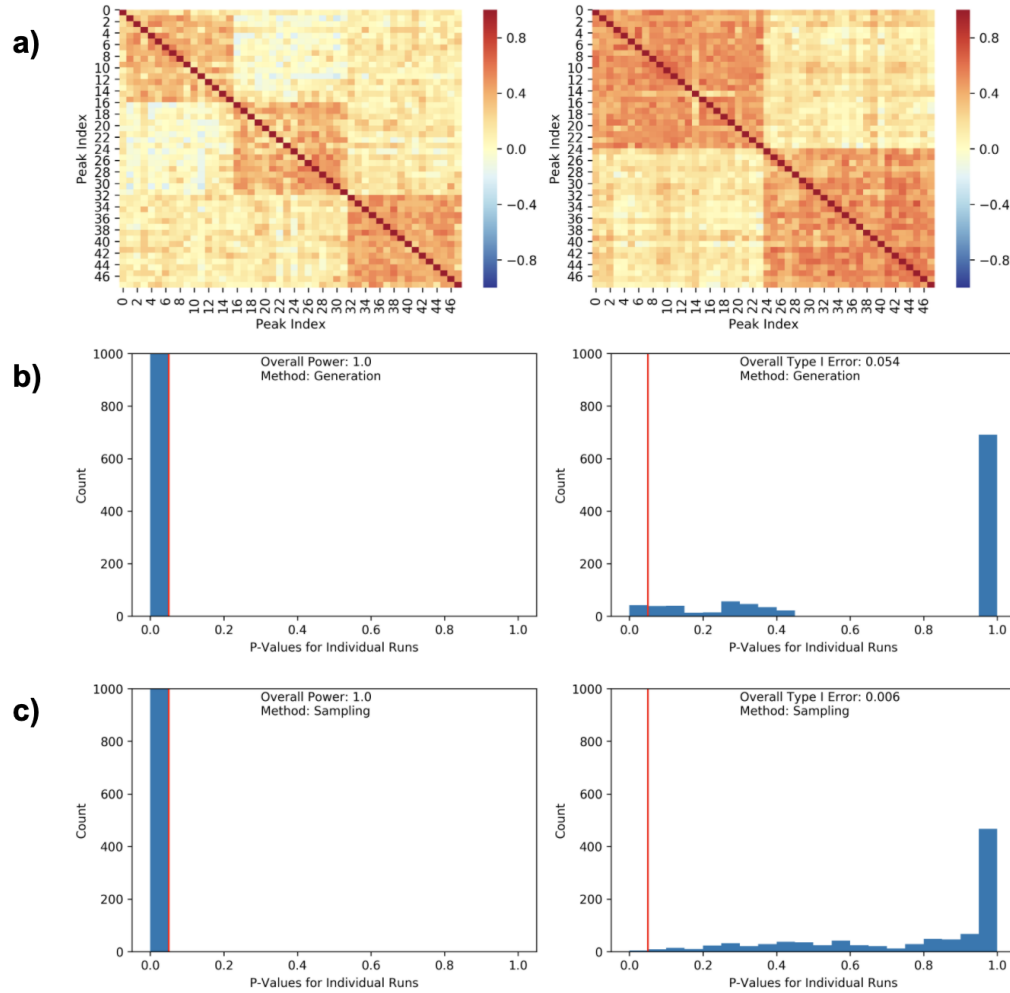


Figure 4-2: Community detection test results for simulated scenario 1. A) Example major and minor allele correlation matrices for this scenario. B) Results for the community detection test with independently generated background data. Plot on the left is p-value histogram across 1,000 tests comparing a dataset with the major allele topology to a dataset with the minor allele topology. Plot on the right is p-value histogram across 1,000 tests comparing two datasets with the major allele topology. Red lines indicate the threshold of $p=0.05$. C) Results for same experiments as in B), except with background data produced by sampling from the original data.

Both versions of the test produced extremely high powers, indicating a strong

ability to discern the differences in correlation structure between the major and minor allele populations. The "gold standard" test, using generation of new background data, produced a false positive rate of approximately 0.05, while the alternative method produced a very small false positive rate of 0.006. This information is essential for setting appropriate significance thresholds in practice, and it is reassuring that both values are relatively low.

- **Scenario 2:** We tested if our methodology could detect small changes in cluster membership. We set the major allele correlation matrix to have three equally sized clusters (of size 16), and we set the minor allele correlation matrix to have clusters of sizes 14, 20, and 14. In this situation, only four peaks change membership while the overall global structure is largely similar.

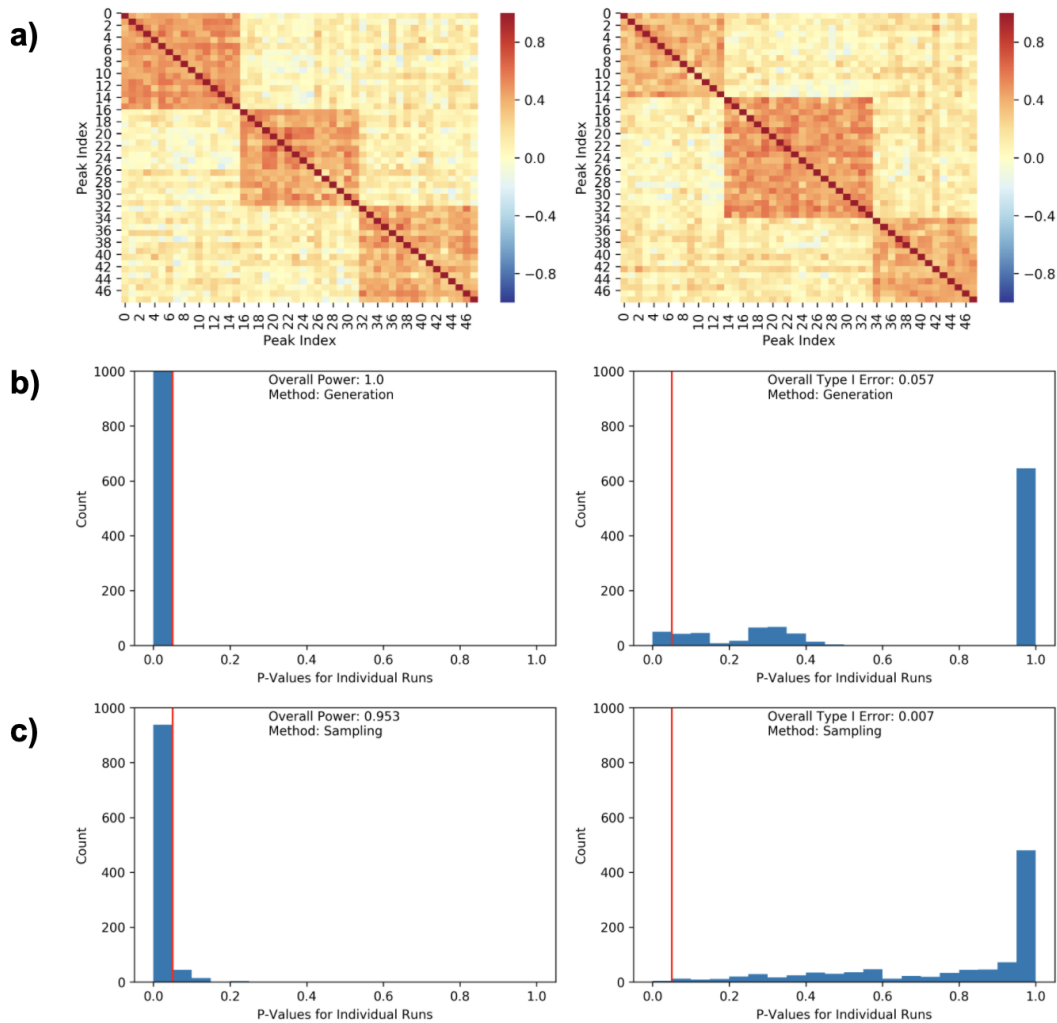


Figure 4-3: Community detection test results for simulated scenario 2. A) Example major and minor allele correlation matrices for this scenario. B) Results for the community detection test with independently generated background data. Plot on the left is p-value histogram across 1,000 tests comparing a dataset with the major allele topology to a dataset with the minor allele topology. Plot on the right is p-value histogram across 1,000 tests comparing two datasets with the major allele topology. Red lines indicate the threshold of $p=0.05$. C) Results for same experiments as in B), except with background data produced by sampling from the original data.

Once again, we see very strong power values and similar false positive rates to the previous case. These results show that when the differences between major and minor allele matrices are drastically reduced, performance remains strong.

- **Scenario 3:** We next introduced free peaks, or peaks that do not belong to any cluster, to our simulations. These are present in virtually every correlation matrix we observed from the experimental data, so we considered it essential to account for these in initial evaluation. We therefore assigned the major allele matrix 3 clusters of size 10 each, with 6 free peaks following each cluster. We assigned the minor allele matrix 3 clusters of size 6 each, with 10 free peaks following each cluster.

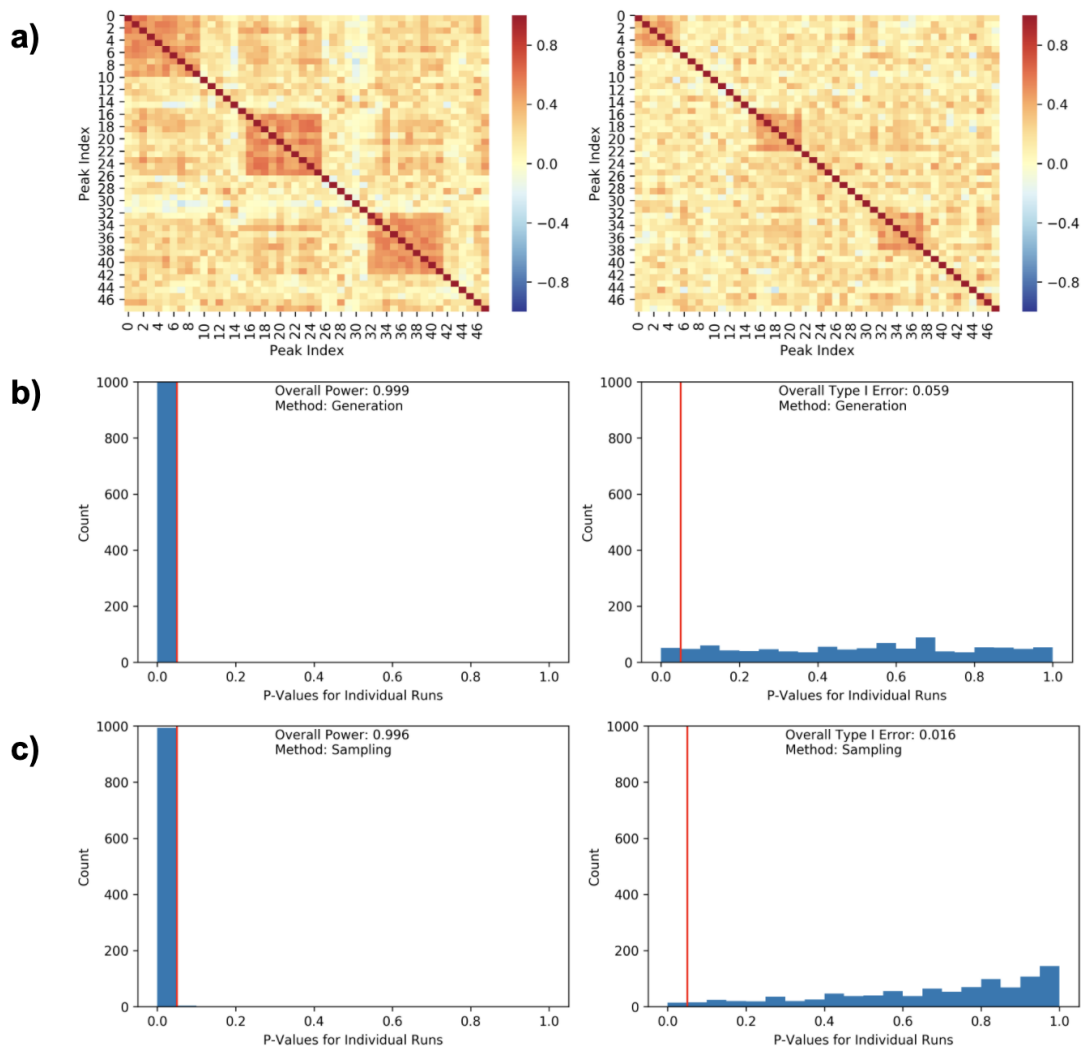


Figure 4-4: Community detection test results for simulated scenario 3. A) Example major and minor allele correlation matrices for this scenario. B) Results for the community detection test with independently generated background data. Plot on the left is p-value histogram across 1,000 tests comparing a dataset with the major allele topology to a dataset with the minor allele topology. Plot on the right is p-value histogram across 1,000 tests comparing two datasets with the major allele topology. Red lines indicate the threshold of $p=0.05$. C) Results for same experiments as in B), except with background data produced by sampling from the original data.

Despite the increased noise provided by the free peaks, which could potentially hamper community detection, the method is indeed able to effectively resolve the differences between the two populations. Power values and false positive rates remained strong.

- **Scenario 4:** We focused on a case commonly found in the super-enhancer data: a correlation matrix with one relatively large cluster and the rest of the peaks lacking cluster membership. Specifically, the major allele matrix was assigned one cluster of size 28 and 20 free peaks, and the minor allele matrix was assigned one cluster of size 33 and 15 free peaks.

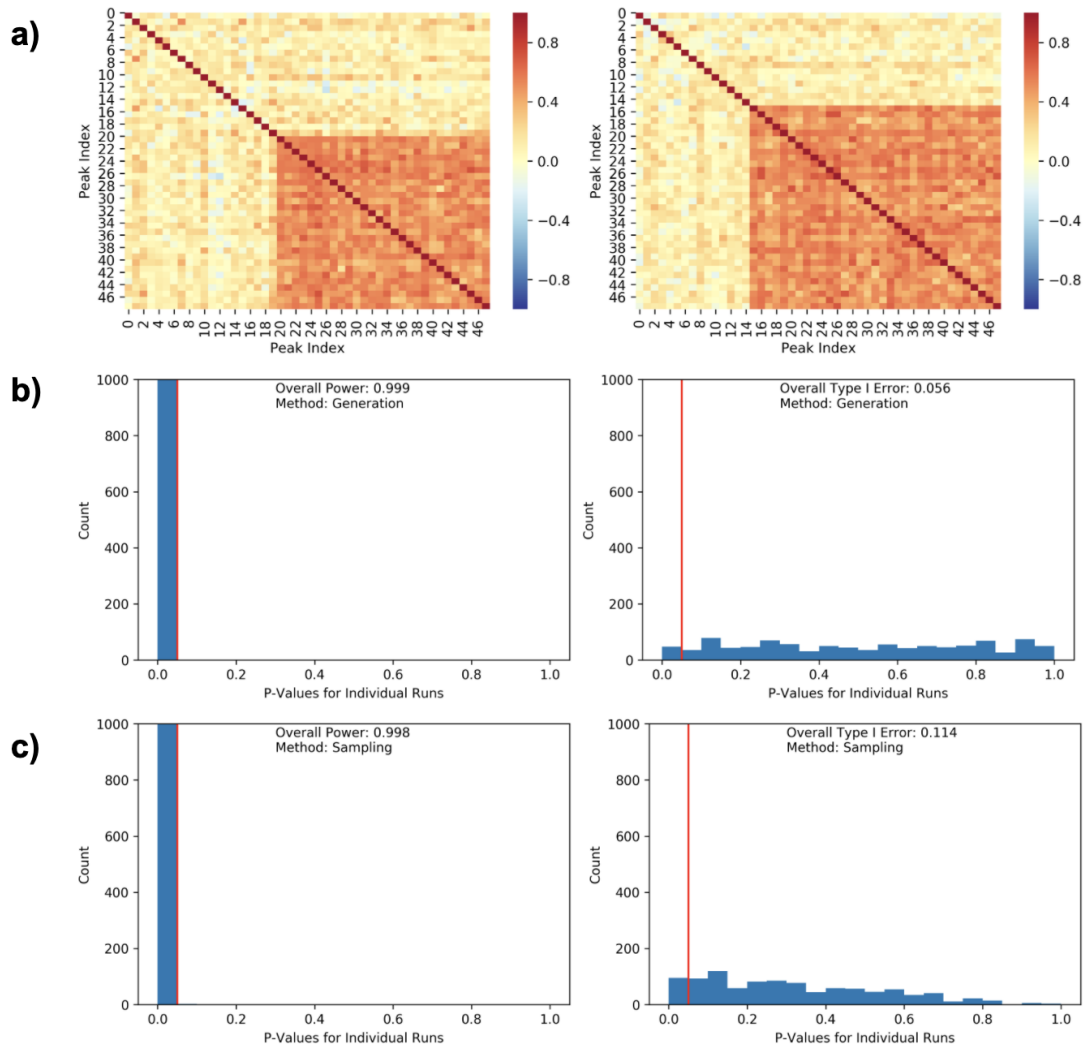


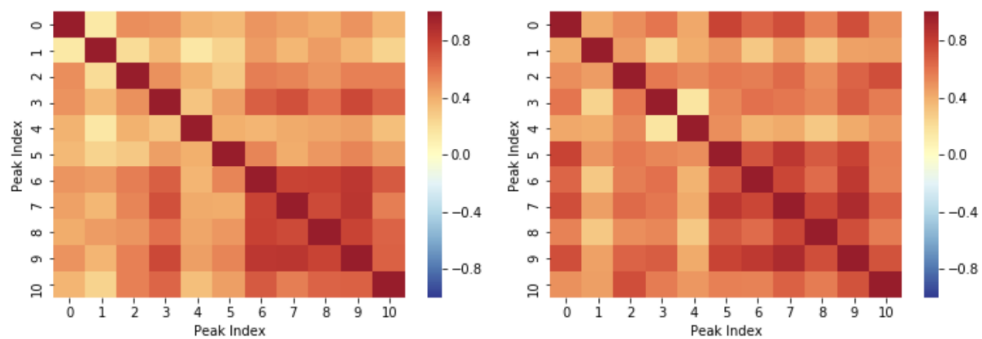
Figure 4-5: Community detection test results for simulated scenario 4. A) Example major and minor allele correlation matrices for this scenario. B) Results for the community detection test with independently generated background data. Plot on the left is p-value histogram across 1,000 tests comparing a dataset with the major allele topology to a dataset with the minor allele topology. Plot on the right is p-value histogram across 1,000 tests comparing two datasets with the major allele topology. Red lines indicate the threshold of $p=0.05$. C) Results for same experiments as in B), except with background data produced by sampling from the original data.

Once again, the test showed great power, and the false positive rate for the "gold standard" test were in line with previous results. However, the type I error rate for the test with sampling-based background data was markedly higher than before. This is possibly due to the higher noise levels, and therefore higher variability, associated with this topology. However, we do not consider this value prohibitively high, and in practice, a higher false positive rate can likely be mitigated by altering significance thresholds.

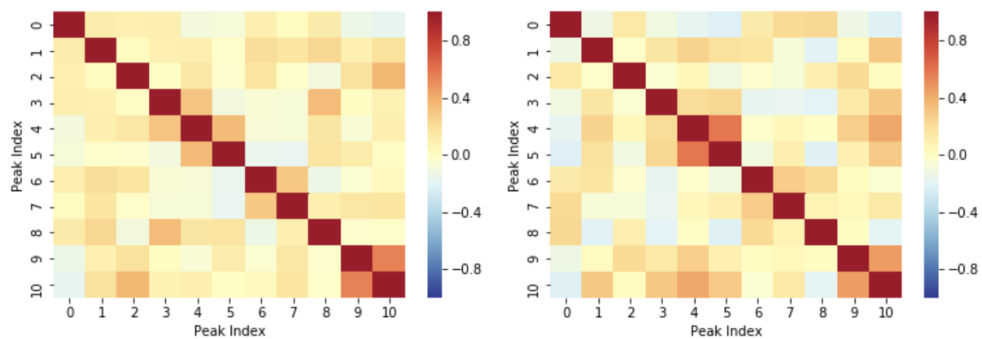
Over a wide range of cases, we find that our method can effectively detect a variety of different types of changes to correlation matrix block structure. In addition, when topologies are identical, false positive rates are relatively low. Of course, we must acknowledge that only an isolated set of correlation structures were tested above. If major and minor allele topologies are closer to each other, then it is likely that the method will have a lower power, while if they are even more different, then the method will likely perform even more strongly. We do not intend for these results to be comprehensive; we simply establish that given a reasonable amount of difference in topology, the method can indeed classify the resulting correlation matrices as significantly different.

4.2 Application to Experimental Data

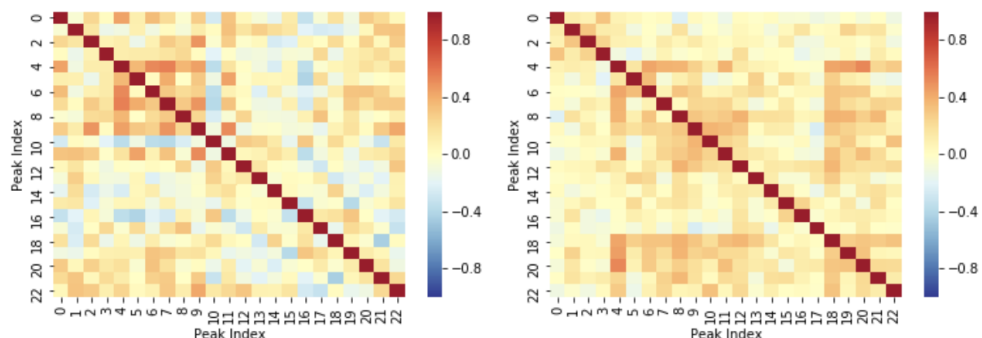
We next tested how our method translated to application on experimental data. Using the k-nearest neighbors approach based on Euclidean distance, we applied our algorithm to the super-enhancer dataset as described in Section 3.2. Specifically, we searched through the SNPs located nearby to each super-enhancer and identified those which returned a p-value of 0 with 100 background mutual information values. We then reran the tests for each of these SNPs with 10,000 background mutual information values to obtain the final p-values. Rather than apply any sort of strict significance threshold, (see Section 4.3), our approach was simply intended to determine if any signal could be detected. Representative samples are depicted in Figure 4-6, along with their p-values.



p=0.0045



p=0.0053



p=0.0096

Figure 4-6: SNPs identified by the community detection-based method. Major allele matrices are on the left, and minor allele matrices are on the right.

In all instances, there appear to be distinct differences in the community structure of the major and minor allele matrices. In the first example, peak 5 joins the large cluster in the minor allele matrix, while this is not the case in the major allele matrix.

In the second example, overlapping communities from peaks 3-5 in the major allele matrix are replaced by a single more cohesive cluster in the minor allele matrix. In the final example, the major allele matrix contains a cluster from peaks 4-9, while in the minor allele matrix, the cluster is much less cohesive and also involves peaks 10-12 and 18-22.

As the method was optimized to detect this exact class of differences, these results provide important validation of the method's performance on real data. We believe the method can be used as a standalone tool or as an effective complement to the permutation-based method described in Chapter 3.

As before, we searched the UCSC Genome Browser to locate protein-coding genes nearby to the super-enhancers associated with these SNPs [13]. Full information is noted in Table 4.1. Especially promising is the SNP located at position 27354393 on chromosome 8, which couples StructureQTL signal with a low GWAS p-value for Alzheimer's Disease. Furthermore, a nearby gene, PTK2B, has previously been linked to Alzheimer's. Thus, it is once again possible that alterations to regulatory structure contribute to the promotion of the disease [8]. As before, this region represents a very attractive target for future study.

Chr	SE Start	SE End	SNP Location	StructureQTL p-value	AD GWAS p-value	Nearby Genes
9	133107516	133147055	133058374	0.0059	1.704x10 ⁻⁵	GTF3C5 CEL RALGDS GBGT1
			133059798	0.0068	3.353x10 ⁻⁵	
			133069821	0.0045	2.112x10 ⁻⁴	
11	47143214	47180183	47222114	0.0053	1.152x10 ⁻⁵	C11orf49 ARFGAP2 PACIN3 DDB2
8	27306915	27379863	27347081	0.0096	5.256x10 ⁻⁵	STMN4 TRIM35 PTK2B
			27354393	0.0113	1.146x10 ⁻⁸	

Table 4.1: Information for SNPs that displayed strong signal when the community detection method was applied. Columns depict the SNP’s chromosome, the start and end points of the super-enhancer the SNP is located nearby to, the location of the SNP itself, the SNP’s StructureQTL and Alzheimer’s Disease GWAS p-values, and protein-coding genes found within 50kb of the associated super-enhancer.

4.3 Analysis and Future Improvements

Through simulations and studies of experimental data, we have developed and validated an effective method to detect block structure changes in correlation matrices. Simulated data allowed us to evaluate our methodology on a variety of pertinent cases, and the eGTEEx super-enhancer dataset provided promising results on real data. We have provided a variety of strategies for community detection, each of which may be optimal for a specific purpose, and the SNPs identified in our super-enhancer study are interesting targets for further analysis. We believe a logical next step would be a full genome-wide analysis of regulatory regions and a subsequent study of SNPs identified as potentially significant.

Of course, several modifications can be made to further enhance efficacy. First, the method must eventually be able to meet much higher significance thresholds. As a means for comparison, GWAS studies often do not report significance unless a p-value

does not exceed 10^{-8} . Although our results were simply for evaluation and preliminary analysis, none of our p-values even came close to these stringent cutoffs. As we ultimately seek to maintain a high power even at high significance thresholds, it is clear that general improvements to the method must be made. Many avenues are possible, but one possible change could be to focus in on the main clusters in the correlation matrix and only examine changes in community membership of the peaks in these clusters. By removing extraneous data, this will hopefully allow the method to achieve a higher confidence.

When considering stringent p-value thresholds, we face another complication. Compared to distribution-based tests like the KS-test, tests which draw a p-value by directly comparing to background values are intrinsically limited due to the number of samples necessary to obtain a certain p-value. For example, if we have a p-value threshold of 10^{-8} , then we would have to generate 10^8 background mutual information values to even have a chance of rejecting the null hypothesis at this confidence level. Of course, this process is extremely time-consuming and computationally intensive. To circumvent this problem, researchers have developed adaptive permutation methods, which provide a means to reduce the expected number of background values necessary per test [4]. In practice, such a method will likely be necessary for StructureQTL determination using community detection.

Furthermore, although the method displayed proficiency in detecting changes in community membership, it was less effective at detecting changes in cluster intensity when membership remained the same. An addition to the method which is optimized for this task could therefore prove highly useful.

Despite these areas of future work, as the results indicate, the method undoubtedly carries great promise. As this method prioritizes different types of differences than the permutation-based method presented in Chapter 3, we believe there is even greater potential in their joint use.

Chapter 5

Conclusion and Future Directions

In this thesis, we propose two complementary methods to determine associations between genetic variants and changes in histone acetylation regulatory structure. The first, which applies the Kolmogorov Smirnov Test, utilizes permutations to determine the likelihood that differences in peak correlation matrices exceed what can arise by chance given the intrinsic variability in the study population. The second, which relies on community detection or clustering algorithms, tests the likelihood that differences in correlation matrices arise from differences in underlying network topologies. Both methods have a strong theoretical basis, and both have shown promising, if early, results on experimentally derived data.

In terms of analysis, logical next steps would be to use both methods to conduct genome-wide analyses on a far more extensive set of SNPs than was studied in this thesis. This type of study would allow us to more powerfully categorize the types of regulatory structure changes which our methods detect as significant. Furthermore, we would be able to create a comprehensive catalog of SNPs which, according to our methodologies, are associated with regulatory changes. These SNPs can then be studied in a variety of ways to determine the biological reasons behind these associations. For example, it would be very telling if most of these SNPs are found in regions with heavy regulatory activity. We can also conduct causal inference studies to determine which of the identified SNPs actually convey effects on regulatory structure and which are simply implicated due to linkage with these causal SNPs. Finally, since

the landscape of gene regulation is cell-type specific, all biological insights gained from this thesis are specific to brain frontal cortex tissue. We could therefore conduct similar analyses using peak data and super-enhancer locations for other tissue types, which in turn would allow us to more readily understand the processes governing gene regulation throughout the body. In addition, our results would be applicable to the study of a much broader variety of diseases.

Of course, there is also much work to be done regarding improvements to both methods, and many of the specifics have already been discussed in previous sections. More generally, we must acknowledge that both methods were largely developed from theoretical perspectives using the principles of statistical inference. There was comparatively little optimization to account the specifics (ie. noise levels, variance) of the acetylation peak datasets the methods would eventually be applied to. For example, most modeling decisions for the community-based method were made using simulations that effectively replicated network topologies but were not specifically designed to mirror the exact peak intensity distributions found in experimental data. While we believe this approach is effective for initial iterations of method development, as presented in this thesis, future improvements must be more closely take into account the data types in question. Returning to the simulations example, we could generate more realistic peak datasets so that the resulting correlation matrices match those found in the experimental data even more closely. Furthermore, we could vary sample sizes and the number of peaks in a region to better mirror what is seen in the experimental data. Analogously, for the permutation test, we can replace the K-S test with a custom test which is optimized to identify the specific changes in distributions that we seek to detect. If changes like these are implemented, then we can more effectively bridge the gap between theory and application.

Analysis and verification were also potentially hampered by dataset size. As the eGTEx data only contains 109 individuals, we often encountered super-enhancers for which major or minor allele populations were considered too small for effective analysis. With a greater sample size, some of these super-enhancers, along with their associated SNPs, may have been flagged as significant. Thus, as the methods pre-

sented in the thesis are improved on in the future, the use of more extensive datasets seems a natural complement.

Overall, the methods described in this thesis are intuitive, theoretically robust, and aimed at addressing a pertinent need in the literature. Furthermore, preliminary results show distinct promise for applications to experimental data. As stated previously, further improvements are certainly needed, but once these are implemented, we believe StructureQTL detection could function as an effective complement to existing methods which detect associations between genetic variation and the mechanisms behind gene expression.

Bibliography

- [1] Thomas Aynaoud. Community detection for NetworkX's documentation. <https://python-louvain.readthedocs.io/en/latest/>.
- [2] Lucy W. Barrett et al. Regulation of eukaryotic gene expression by the untranslated gene regions and other non-coding elements. *Cell Molecular Life Science*, 69(21):3613–3634, Apr 2012.
- [3] Vincent D Blondel et al. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008(10), Sep 2008.
- [4] Ronglin Che et al. An adaptive permutation approach for genome-wide association study: evaluation and recommendations for use. *BioData Mining*, 7(1), 2014.
- [5] ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature*, 489(7414):57–74, 2012.
- [6] Menno P. Creyghton et al. Histone h3k27ac separates active from poised enhancers and predicts developmental state. *Proceedings of the National Academy of Sciences*, 107(50):21931–21936, 2010.
- [7] O. Delaneau et al. Chromatin three-dimensional interactions mediate genetic effects on gene expression. *Science*, 364(6439), Feb 2019.
- [8] P. Dourlen et al. Functional screening of Alzheimer risk loci identifies PTK2B as an in vivo modulator and early marker of tau pathology. *Molecular Psychiatry*, 22(6):874–883, 2016.
- [9] Jianxing Feng et al. Identifying ChIP-seq enrichment using MACS. *Nature Protocols*, 7(9):1728–1740, 2012.
- [10] Michael D. Gallagher and Alice S. Chen-Plotkin. The post-GWAS era: From association to function. *The American Journal of Human Genetics*, 102(5):717–730, 2018.
- [11] Denes Hnisz et al. Super-enhancers in the control of cell identity and disease. *Cell*, 155(4):934–947, 2013.

- [12] Jialiang Huang et al. Dissecting super-enhancer hierarchy based on chromatin interactions. *Nature Communications*, 9(1), May 2018.
- [13] D. Karolchik. The UCSC Genome Browser database. *Nucleic Acids Research*, 31(1):51–54, Jan 2003.
- [14] Aziz Khan and Xuegong Zhang. dbSUPER: a database of super-enhancers in mouse and human genome. *Nucleic Acids Research*, 44(D1), Apr 2015.
- [15] Marc. R. Mansour et al. An oncogenic super-enhancer formed through somatic mutation of a noncoding intergenic element. *Science*, 346(6215):1373–1377, 2014.
- [16] Alexandre P. Marand et al. Towards genome-wide prediction and characterization of enhancers in plants. *Biochimica et Biophysica Acta (BBA) - Gene Regulatory Mechanisms*, 1860(1):131–139, 2017.
- [17] N. D. Paauw et al. H3k27 acetylation and gene expression analysis reveals differences in placental chromatin activity in fetal growth restriction. *Clinical Epigenetics*, 10(1), 2018.
- [18] Abhijeet Pataskar et al. Deciphering the gene regulatory landscape encoded in DNA biophysical features. *iScience*, 21:638–649, 2019.
- [19] Dilmi Perera et al. Differential DNA repair underlies mutation hotspots at active promoters in cancer genomes. *Nature*, 532(7598):259–263, 2016.
- [20] Daria Shlyueva et al. Transcriptional enhancers: from properties to genome-wide predictions. *Nature Reviews Genetics*, 15(4):272–286, Nov 2014.
- [21] Barbara E. Stranger et al. Enhancing GTEx by bridging the gaps between genotype, gene expression, and disease. *Nature Genetics*, 49(12):1664–1670, Nov 2017.
- [22] Wenjie Sun et al. Histone acetylome-wide association study of Autism Spectrum Disorder. *Cell*, 167(5), 2016.
- [23] Meng-Shan Tan et al. Bridging integrator 1 (BIN1): form, function, and Alzheimers Disease. *Trends in Molecular Medicine*, 19(10):594–603, 2013.
- [24] Tuan Trieu and Ekta Khurana. A deep learning approach to predict the impact of non-coding sequence variants on 3D chromatin structure. Oct 2019.
- [25] Yuan Wang et al. Genome-wide methylome and chromatin interactome identify abnormal enhancer to be risk factor of breast cancer. *Oncotarget*, 8(27), Feb 2017.
- [26] Warren A. Whyte et al. Master transcription factors and mediator establish super-enhancers at key cell identity genes. *Cell*, 153(2):307–319, 2013.

- [27] Monique G. P. Van Der Wijst et al. Single-cell RNA sequencing identifies celltype-specific cis-eQTLs and co-expression QTLs. *Nature Genetics*, 50(4):493–497, 2018.
- [28] Naomi R. Wray. Allele frequencies and the r^2 measure of linkage disequilibrium: Impact on design and interpretation of association studies. *Twin Research and Human Genetics*, 8(2):87–94, Jan 2005.
- [29] Ian T. Young. Proof without prejudice: use of the Kolmogorov-Smirnov test for the analysis of histograms from flow systems and other sources. *Journal of Histochemistry Cytochemistry*, 25(7):935–941, 1977.