

A High-Resolution Study of the Chromatin Environment Around Regulatory Elements

by

Hayden C. Metsky

S.B. in Computer Science and Engineering and in Physics
Massachusetts Institute of Technology (2013)

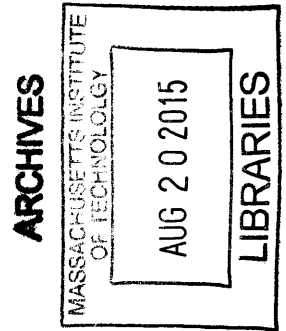
Submitted to the Department of Electrical Engineering
and Computer Science
in Partial Fulfillment of the Requirements for the Degree of

Master of Engineering in Electrical Engineering and Computer Science
at the Massachusetts Institute of Technology

September 2014

Copyright 2014 Hayden C. Metsky. All rights reserved.

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



ⁿ
Signature redacted

Author
Department of Electrical Engineering and Computer Science
September 17, 2014

ⁿ
Signature redacted

Certified by
Manolis Kellis, Professor, Thesis Supervisor
September 17, 2014

ⁿ ^{///}
Signature redacted

Accepted by
Prof. Albert R. Meyer
Chairman, Master of Engineering Thesis Committee



77 Massachusetts Avenue
Cambridge, MA 02139
<http://libraries.mit.edu/ask>

DISCLAIMER NOTICE

Due to the condition of the original material, there are unavoidable flaws in this reproduction. We have made every effort possible to provide you with the best copy available.

Thank you.

The images contained in this document are of the best quality available.

A High-Resolution Study of the Chromatin Environment Around Regulatory Elements

by

Hayden C. Metsky

Submitted to the
Department of Electrical Engineering and Computer Science

September 17, 2014

in Partial Fulfillment of the Requirements for the Degree of
Master of Engineering in Electrical Engineering and Computer Science

ABSTRACT

Chemical modifications to histones, the proteins around which DNA wraps, are believed to play an important role in gene regulation. These modifications, along with others, make up a cell’s “epigenome.” It is known that the presence of a particular combination of these modifications at a region of a cell’s genome determines, for that region, a state that carries functional significance. This work seeks to better understand the importance of not just presence, but also distribution of modifications within regulatory regions. One approach aimed at improving our understanding is to cluster regulatory regions based on information contained in signals that describe, at a high-resolution, the distribution of these modifications.

In this thesis we develop a tool, called ChromSMS, to perform this clustering in a biologically meaningful and efficient way that is versatile in handling the underlying complexities of these signals. We apply the tool to data from the NIH’s Roadmap Epigenomics Project to analyze ChromSMS and to better understand the mechanisms behind the patterns we observe. We find that ChromSMS produces meaningful clusters that are different from each other at a statistically significant level. Using ChromSMS to conduct analyses of epigenomic data, we discover strong relations between GC-content and the distribution of particular modifications. Furthermore, we uncover a small number of patterns that display high functional enrichment, and we begin to study the possible role and significance of motifs in driving these patterns. We conclude that ChromSMS can serve as a useful tool in examining regulatory regions at a high-resolution.

Thesis Supervisor: Manolis Kellis
Title: Professor

Acknowledgments

The research presented in this thesis was carried out under the supervision of Professor Manolis Kellis in MIT's Computer Science and Artificial Intelligence Laboratory.

I offer my profound gratitude to Manolis, who has supported me throughout this research with his knowledge, ideas, and vision. I am thankful for his guidance in this thesis, while also allowing me the chance to pursue my own approaches. I am especially grateful for his sustained enthusiasm for the subject, both as an advisor and a teacher.

Much of my progress in this thesis would not have been possible without Wouter Meuleman. I would like to thank Wouter for all that he has taught me, as well as for his unwavering patience in answering all of my questions. He has provided fantastic direction in my research and has been an amazing mentor to work with.

There are many others in the Kellis lab who have repeatedly offered their help and thoughts, and I am incredibly appreciative of all of them.

I thank Anshul Kundaje, who initially prompted my interest in this project two years ago. I also thank Manuel Garber and Mitch Guttman, my first mentors in computational biology, for their generosity in advising a young student and for contributing to the growth of my curiosity in the field.

Finally, I thank my parents for their guidance and their continual support.

Contents

1	Overview	15
1.1	Processing Epigenomic Data	17
1.2	ChromSMS: the Chromatin Signal Magnitude and Shape clustering tool	19
1.3	Applying ChromSMS to Analyze the Chromatin Environment	20
2	Epigenomics Background	21
2.1	The Structure of Chromatin	22
2.2	The Landscape of Histone Modifications	24
2.3	Reading Epigenomic Signals	26
2.4	ChromHMM and Finding Patterns of Activity	28
3	Past Work on Uncovering Patterns in the Chromatin Environment	31
3.1	Survey of Past Results Related to This Thesis	31
3.2	Outline of Two Computational Techniques That Seek Patterns in the Chromatin Environment	33
3.2.1	Binarizing the presence of each modification	33
3.2.2	Modeling DNase magnitude and shape	33
3.3	CAGT: The Clustered AGgregation Tool	34
3.3.1	The procedure	34
3.3.2	Limitations of what it accepts and analyzes	35
4	ChromSMS: the Chromatin Signal Magnitude and Shape clustering tool	39

4.1	Overview of its Design as a System	40
4.2	Preprocessing Input	41
4.2.1	Selecting regions to avoid overlap	41
4.2.2	Aligning regions and extracting signals	42
4.2.3	Filtering signal profiles with low complexity	44
4.3	Clustering Signal Profiles	45
4.3.1	Dissimilarity measurements between two profiles	45
4.3.1.1	Varied-length mode	45
4.3.1.2	Multiple-mark mode	47
4.3.2	Implementing agglomerative hierarchical clustering	50
4.3.2.1	Flipping regions to maintain orientation	52
4.3.2.2	Terminating the merging process	53
4.3.2.3	Other relevant clustering algorithms	55
4.3.3	Reporting results and creating visualizations	56
4.4	Implementation Details and Performance	57
5	A Review of Discovered Shapes and Analyses Within These Shapes	59
5.1	Data Provided as Input to ChromSMS	60
5.2	Review of Output from Varied-length Mode	61
5.3	Review of Output from Multiple-mark Mode	62
5.4	Associations Between Directionality and Positions of Nearby Genes	65
5.4.1	Directionality in promoters	65
5.4.2	Directionality in enhancers	68
5.5	Functional Enrichment of Discovered Shapes	68
5.6	Correlations Between Signal Profiles and GC-content	71
5.7	Motif Enrichment Within Discovered Shapes	74
5.7.1	Method for determining enrichment	74
5.7.2	Analyzing intra-shape motif similarity	75
5.7.3	Findings	76
5.8	Discussion	77

6	Analyses Across Discovered Shapes	81
6.1	A Statistical Test for the Comparison of Shapes	81
6.1.1	Method: measuring similarity between two shapes	82
6.1.2	Method: evaluating significance of difference between two shapes	83
6.1.3	Findings	84
6.2	Comparing Motif Enrichment Across Shapes	84
6.2.1	Comparison of motif clusters: sequence	86
6.2.2	Comparison of motif clusters: dense shapes	88
6.3	Discussion	89
7	Further Applications of Ideas Behind ChromSMS	91
7.1	Detecting Regulatory Regions by Searching for Shapes	91
7.1.1	Dataset we use and constructing reference shapes	92
7.1.2	Generating emissions for the HMM	93
7.1.3	Determining hidden states from the HMM	93
7.1.4	Discussion	94
7.2	Querying for Region Near Neighbors via Locality-Sensitive Hashing .	94
7.2.1	Comparing regions with Hamming distance	95
7.2.2	Comparing regions with Euclidean distance	96
7.2.3	Results: query time	96
7.2.4	Results: accuracy	97
7.2.5	Discussion	97
8	Conclusions	99
A	Tables	101
B	Figures	107

List of Figures

1-1	Example of the information to be processed	18
2-1	View of DNA wrapping around histones	22
2-2	Distribution of histone modifications around genes	25
2-3	Heatmap of enhancer activity clusters	29
3-1	Summary of CAGT's procedure	35
4-1	Overview of the design of ChromSMS	40
4-2	Example of selecting regions to avoid overlap	43
4-3	The consequence of naïvely merging signal profiles	53
4-4	Curve of minimum distances between clusters while merging	54
5-1	Example 1 of varied-length mode output	62
5-2	Example 2 of varied-length mode output	63
5-3	Example 1 of multiple-mark mode output	64
5-4	Example 2 of multiple-mark mode output	65
5-5	Relation between TSS and directionality within shapes	67
5-6	Two shapes that display a high degree of functional enrichment	71
5-7	Correlation between histone modification patterns and GC-content	73
5-8	Counts of enriched motifs in shapes	77
6-1	Patterns in the enrichment of motifs across 242 shapes	85
6-2	Similarities across motif clusters	87

B-1	Examples of active plots, control plots, and heatmaps for shapes output by ChromSMS on two enhancer activity clusters	107
B-2	Example of top functional associations in a shape	110

List of Tables

5.1	Summary of significant functional associations by ontology	70
7.1	Summary of regions behind 10 hidden states	94
7.2	Results for query points using Hamming distance	98
7.3	Results for query points using Euclidean distance	98
A.1	Details about each of 19 clusters of motif instances	101

Chapter 1

Overview

In 1942 the biologist Conrad Waddington first proposed the name **epigenetics** to describe studies that attempted to understand the causal mechanisms at work when genes in the **genotype**, or hereditary information of an organism, bring about effects in its **phenotype**, or observed traits [40]. Even after biology's renaissance of the past seven decades, the related term **epigenomics** still encompasses, at its core, the same idea. The field serves to bridge genotype and phenotype by helping us better understand how observed traits change even with no change in the underlying DNA sequence [19]. In particular, epigenomics is about the mechanisms by which stable and heritable changes in phenotype can occur without alterations to the base pairs of DNA.

Perhaps the most stunning consequence of epigenomics is the explanation it offers for cellular differentiation [19]. Although most cells in an organism share the same underlying DNA sequence, they certainly display expression of different genes and exist for different functions. The mechanisms studied in epigenomics can describe such phenomena. Indeed, unlike the genotype, the states that underlie epigenomic phenomena are both reversible and transient and are influenced by an organism's aging and its environmental factors [13].

Recently there has been a tremendous growth of interest in the field. This is likely due in part to our better understanding of the great importance that epigenomics has in biology, but also to amazing advances in our technological ability to gather epige-

omic information at a genome-wide scale [13]. Such a swell of research and interest in the emerging field of course enhances our understanding of basic biology. But its unique role also allows us to further understand how biological phenomena change in individuals at different developmental stages and in individuals with different diseases, which can offer direct impacts on medicine. As James Watson noted in 2003, “the major problem, I think, is chromatin. . . you can inherit something beyond the DNA sequence. That’s where the real excitement of genetics is now” [19].

Today research in the field is largely focused on the mechanisms that govern the structure of **chromatin**, which is the system in a cell’s nucleus consisting of DNA and the proteins (called **histone** proteins) that work to organize DNA. The overall structure of chromatin can alter gene expression and thus phenotype, so its study is a natural focus of epigenomics. In particular, it is known that covalent and non-covalent chemical modifications to both DNA and histone proteins alter chromatin structure [19]. Hence, much research is being done to better understand the mechanisms behind these modifications as well as their effects.

The National Institutes of Health’s Roadmap Epigenomics project is a significant, ongoing effort to produce and analyze data that can aid studies in epigenomics [7]. The project takes advantage of the increased throughput and sensitivity offered by next-generation sequencing and array technologies [13]. In doing so, it offers a catalog of epigenomic information — such as data about modifications of DNA and histones — for a variety of human cell types at a genome-wide scale.

This thesis fits into the larger endeavor of understanding the epigenomic information carried by histone proteins. In particular, our goal is to study, at a high-resolution (roughly 10 base pairs), this information at **regulatory regions** of the genome, which are segments involved in increasing or decreasing the expression of genes. Of particular interest in this thesis are signals, provided by the Roadmap project, that describe the enrichment of various types of histone modifications across the genome.

The paramount problem considered in this thesis is that of finding distinct clusters of regulatory regions such that the regions in each cluster display similar patterns in how they carry epigenomic information in a particular cell type. The primary

objective of this thesis is to develop a tool that addresses this problem in a biologically meaningful and computationally efficient way. Given that the provided signals are noisy at a high-resolution and that the number of regulatory regions we use is large, this is a challenging problem from both the biological and computational fronts. The other objective of this thesis is to use this tool to conduct an analysis of the epigenomic landscape in chromatin around regulatory elements.

In the remainder of this chapter we continue to provide motivation for and an overview of this thesis. In Chapter 2 we offer a brief background in the biology of epigenomics, which may aid the reader in understanding the remainder of this thesis. Then, in Chapter 3 we survey existing methods for solving the previously mentioned problem and discuss the biological results obtained from these methods. Next, in Chapter 4 we describe the tool developed in this thesis, called ChromSMS, by discussing the framework behind the tool and some details of its implementation. In Chapters 5 and 6 we describe and discuss the biological insights we gather from the application of ChromSMS. In Chapter 7 we provide two further applications of the ideas behind ChromSMS, each of which was briefly explored as part of this thesis. Finally, we offer concluding remarks in Chapter 8.

1.1 Processing Epigenomic Data

The primary type of data considered in this thesis is a signal that conveys the patterns in which a particular cell type carries a particular chemical histone modification. The signal spans the entire genome (namely, the epigenome) and its magnitude at each base pair is a measure of the enrichment of a modification on histones found at that position. Figure 1-1 shows an example to demonstrate this idea. It displays eight signals over a 2,000 base pair (bp) region of the genome. There is a clear diversity in the magnitudes and shapes of these signals. For example, one signal (black) displays symmetry and is enriched at the middle of the region. Another (yellow) also displays symmetry, but shows a drop in enrichment at the middle. Two others (blue, purple) display asymmetry and are enriched more on one side of the region than the other.

Relative to these four, the other four signals display little enrichment. Section 2.3 describes more precisely the data stored by these signals and what they measure.

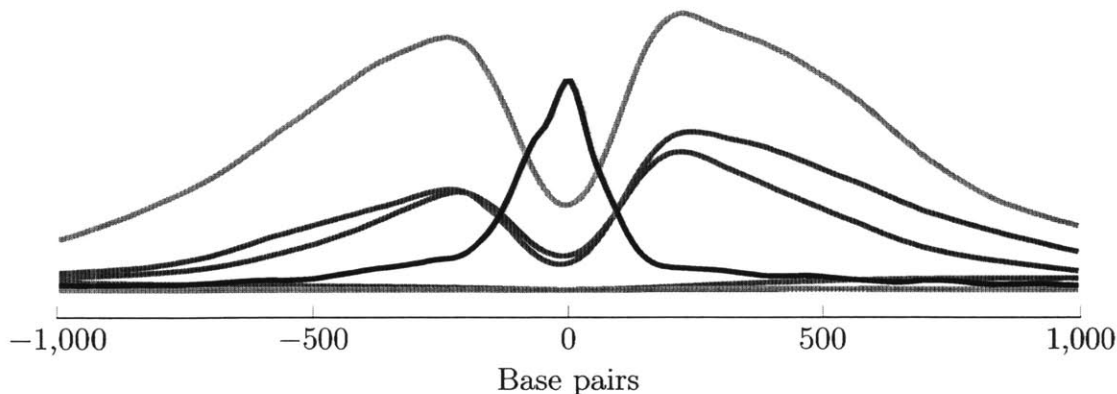


Figure 1-1: An example of the type of information processed in this thesis. The ticks on the axis give the number of base pairs offset from an anchor point. We show eight signals, each of which describe a “mark.” Seven marks correspond to a different chemical modification to histone proteins, and the other (black) corresponds to a measure of a region’s accessibility to proteins. The signal magnitude is a measure of the enrichment of a mark at a position in the genome. The signals displayed are the average over observed signals at 100 regions of the epigenome, each of which is 2,000 bp long and centered around a promoter, which is defined in Section 2.2.

Performing computations on this type of epigenomic data is, in many ways, different than working with typical genomic data. For example, rather than assigning to each position in the genome one of four base pairs, we assign it a value from some (in theory, unbounded) range. Moreover, the value itself is continuous. These features present a number of challenges to addressing the principal problem of this thesis. The first is that the signal may contain noise and, namely, a small region (or window) of the signal may convey little information, which presents a problem for high-resolution studies. Second, intrinsically a signal with continuous values from some large range may exhibit very different patterns in different regions.

Furthermore, there are many different types of histone modifications at any given region (we study on the order of 10). Thus, a region of the epigenome is described by many signals. As Figure 1-1 displays, within one small region these signals may take on very different patterns — i.e., both different shapes and different magnitudes. This is yet another challenge.

An additional complication arises from the notion that similar regulatory regions may be of considerably different lengths. This necessitates comparing signals extracted from different sized windows.

Indeed, these unique challenges, as well as others, make it difficult to efficiently find meaningful clusters of regulatory regions based on the patterns in which they carry epigenomic information. A tool to address this problem while tackling such challenges must draw on a wide collection of techniques, including ideas from signal processing. ChromSMS is a tool that strives to do this.

1.2 ChromSMS: the Chromatin Signal Magnitude and Shape clustering tool

There do exist tools that aim to solve the problem posed above: clustering regulatory regions based on the patterns in which histones in those regions carry epigenetic information. One such tool is the Clustered AGgregation Tool, or CAGT [28]. CAGT is discussed in more detail in Section 3.3. While CAGT has yielded novel insights, it does not account for some of the features present in epigenomic data. For example, when clustering regions it does not account for the patterns given by different types of histone modifications; instead, it considers just one modification. Similarly, it does not tackle the complication that regulatory regions may be of different lengths.

ChromSMS builds on the ideas present in tools like CAGT, but accounts for many of these challenges. It aims to find clusters of similar regulatory regions by considering both the magnitude and shape information encoded by signals for various histone modifications in those regions. By appropriately handling the complications intrinsic in this biological data, our goal is for ChromSMS to better capture patterns than existing methods.

Of course, accounting for these challenges introduces significant computational complexity. As a software package, ChromSMS aims to account for them in a way that — while surely requiring more resources than CAGT — is still usable, practical,

and scalable to large input. This involves making certain assumptions about the data, adapting existing algorithms to this use case, making use of distributed computation, and performance engineering of the software.

1.3 Applying ChromSMS to Analyze the Chromatin Environment

The second objective of this thesis is to apply ChromSMS to a real dataset, namely one provided by the Roadmap Epigenomics project. We use signal tracks for various histone modifications across 127 cell types, and we are provided with collections of regions that are classified as putative promoter or enhancer regulatory regions. Section 5.1 describes in more detail the input that we use.

Using ChromSMS we verify the key results of [28], including the idea that histone modifications are deposited in a remarkably asymmetrical manner about the sites to which chromatin-binding proteins bind within regulatory elements. We find that the patterns discovered by ChromSMS are distinct from one another at a significant level and hold intrinsic biological meaning.

We uncover a limited number of patterns that display a high degree of functional enrichment, and we also discover a strong correlation between GC-content and the patterns of particular histone modifications; these studies are described in Sections 5.5 and 5.6, respectively. In an effort to better understand the mechanisms behind the patterns we observe, we also investigate the relations that our discovered patterns have with sequence motifs, as discussed in Sections 5.7 and 6.2.

From our analyses we conclude that ChromSMS can serve a helpful role in studying the chromatin environment at a high-resolution around regulatory elements.

Chapter 2

Epigenomics Background

In eukaryotic cells the genome must be densely compacted, in some reproducible three-dimensional structure, to fit inside the small space of a nucleus [43]. This compaction is aided by a complex of proteins that, together with DNA, form **chromatin**. Although the compaction must be dense, it must also be highly organized in a way that allows access to DNA by proteins involved in transcribing, replicating, and repairing it. To regulate such access, there are numerous mechanisms that alter the properties of chromatin. Some of these mechanisms are encoded in a stable, heritable form and are the focus of the field of epigenomics [13].

In this chapter we outline the structure of chromatin and describe, at a high-level, some of the well-studied mechanisms that alter it. We focus on one phenomenon, that of histone modifications, that has a profound association with gene regulation and is the main subject of this thesis. Next, we describe the biological experiments used to gather data about histone modifications and discuss the data itself, as well as the steps by which it is processed. Finally, we provide an overview of a powerful tool, called ChromHMM, that uses such data to discover, *de novo*, regulatory regions of the genome. We also discuss how its output has been used to cluster regions based on their activity patterns across cell types; these regions and activity-based clusters are used heavily as input to ChromSMS, the tool developed and applied in this thesis.

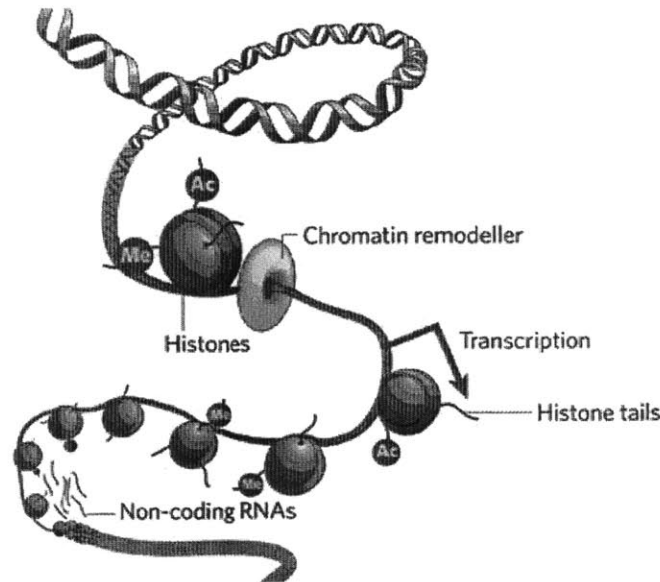


Figure 2-1: A portion of chromatin is shown, in which a segment of DNA wraps around histone proteins that make up nucleosomes. Source: This is a cropped portion of Figure 1 from [24].

2.1 The Structure of Chromatin

At its core, chromatin is a repeating structure with a fundamental unit known as the **nucleosome** [31]. A nucleosome consists of eight **histone** proteins around which 147 bp of DNA wraps, which equates to about 1.65 turns. Figure 2-1 displays this structure. Because of the close interaction between the DNA and histone proteins, the nucleosome is a stable complex that also plays a central role in regulating access to DNA. This regulation is achieved through mechanisms that alter DNA and nucleosomes.

One notable mechanism, known as **DNA methylation**, involves a direct chemical modification to DNA itself in which a methyl group is added to a DNA nucleotide, typically cytosine in a CpG dinucleotide. Experiments demonstrate that DNA methylation has the effect of stably repressing transcription [9]. This includes repression of “invading DNA,” or sequences embedded in the genome that may prove fatal. Furthermore, there is strong evidence that DNA methylation is fundamental to development and differentiation of cells. Indeed, knockout experiments in mice have clearly demonstrated the necessity of DNA methylation for survival. DNA methylation patterns

also play a pivotal role in diseases such as cancer, where tumor suppressor genes may be silenced due to unusually high methylation [8]. Methylated bases function by physically blocking regulatory proteins from binding to DNA, as well as by attracting a complex of other proteins that mediate repression.

Along with direct modifications to DNA, there are other well-studied mechanisms for regulating access that involve altering nucleosomes and their histone proteins. For example, there exist macromolecular assemblies that use energy from ATP to physically alter chromatin structure and are involved in many cellular processes, including repositioning and exchanging nucleosome to aid in transcription [20]. Additionally, covalent posttranslational modifications to the tails of histones that form nucleosomes are associated with regulation. These are known as **histone modifications** and have been studied extensively since the discovery in 1964 that histones associated with active genes are acetylated at a level beyond what would be expected [43]. Since that study, over 100 histone modifications have been discovered and described. We provide a more extensive review of histone modifications in Section 2.2. These phenomena, which restructure chromatin, are part of a process referred to as **chromatin remodeling**.

In altering the chemistry and structure of chromatin, these mechanisms are believed to modulate gene expression and thus have far-reaching consequences. Two of them — DNA methylation and histone modifications — clearly carry information that is encoded in the type, level, and genomic position of the chemical modification. In particular, these types of modifications are said to be carriers of **epigenomic information**.

There is strong evidence that some of this information is heritable. Patterns of DNA methylation are retained through cell division by a process in which an enzyme methylates cytosine bases in the newly synthesized strand of DNA using methylation on the template strand [8]. Similarly, a subset of histone modifications appears to show inheritance in yeast and multicellular organisms, as mediated by multiple protein complexes, although the details behind this inheritance remain unclear [8]. Hence, through these chemical modifications, cells can inherit information and properties

that are not encoded in the base pairs of DNA.

The term **epigenome** refers to information carried by DNA methylation and all histone modifications, even though, as far as we know, only a subset of the histone modifications demonstrate inheritance [8]. In the field of **epigenomics** we study the epigenomes of particular cell types, as well as how they differ across cell types and individuals.

2.2 The Landscape of Histone Modifications

There are four families of histone proteins — H2A, H2B, H3, and H4 — and each nucleosome consists of two of each [27]. The histone proteins H3 and H4 have a long tail that begins at the N-terminal end and modifications (e.g., the addition of a methyl or an acetyl group) occur to particular amino acids in this tail. The naming of a histone modification provides the name of the modified histone, the amino acid that is modified (typically Lysine, abbreviated K), the position of the amino acid from the N-terminal, the type of modification, and the number of modifications. (When the number of modifications is omitted, there is just one.) For example, H3K4me3 denotes three methyl groups (i.e., tri-methylation) added to the 4th amino acid of the H3 protein. Similarly, H3K9ac denotes acetylation of the 9th amino acid of the H3 protein.

Histone modifications are strongly associated with fundamental regulatory processes. Indeed, the term **histone code** has been used to describe the notion in which the information carried by histone modifications may be involved in exhaustive regulatory mechanisms [23]. In establishing this “code,” many studies have linked modifications to the states of types of **regulatory regions**. These regions are segments of DNA that are involved in the mechanisms behind increasing or decreasing the expression of particular genes. They include, among other types, enhancers and promoters. **Enhancers**, which can be up to millions of base pairs from their target gene(s), are bound by proteins, called **transcription factors**, that facilitate the transcription of a target gene, thereby increasing its expression. **Promoters** are just

upstream of a gene and are connected with the initiation of its transcription; they too are bound by transcription factors. Figure 2-2 displays the connections between modifications and regulatory states at promoter regions.

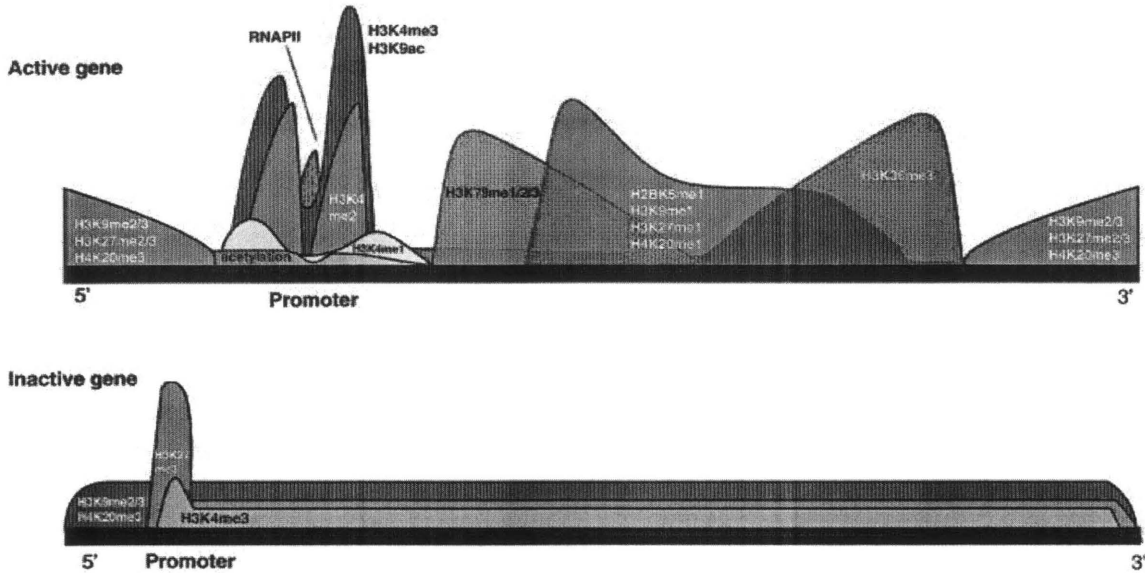


Figure 2-2: A visualization of the distribution of histone modifications around active (top) and inactive (bottom) genes. The patterns display enrichment of different histone modifications at a promoter and on the gene body. Not all of these modifications will be explored in this thesis. Most importantly for this thesis, note that H3K4me3 and H3K9ac are often found highly enriched in the promoters of active genes, and that H3K36me3 is often found enriched in the gene bodies of active genes; H3K27me3 is found enriched at the promoters of inactive genes, and H3K9me3 is found enriched in their gene bodies. Although this figure only shows the region around promoters and genes, similar correlations exist between histone modifications and enhancers. The modifications H3K27ac and H3K4me1 are typically found enriched in active enhancers [42]. Source: This figure is an edited copy of Figure 1 from [5].

While correlations between modifications and regulation are clear (e.g., a modification appears on a gene only when it is transcribed), it can be difficult to prove the causative nature of a modification, especially in humans [27]. Rather, it is often assumed with some degree of confidence based on the extent of analysis performed on the modification. In cases where causality can be confirmed, the modifications function by both directly altering the high-order structure of chromatin and by attracting particular protein complexes to DNA, which can in turn alter the structure or affect levels of transcription.

It is important to note that a naïve picture of histone modification patterns might be static, but in reality it is highly dynamic [27]. Detailed patterns differ not only across cell types and individuals, but also across just minutes of time. Hence, an aggregate analysis of many histone proteins under particular conditions provides information about only a subset of the diverse patterns.

2.3 Reading Epigenomic Signals

Over the last decade there has been a substantial advancement in the technology used for analyzing the incidence of histone modifications across an entire epigenome. Recently, the frontier of this technology has been ChIP-seq, which combines a procedure known as chromatin immunoprecipitation with massively parallel sequencing [35]. The powerful technique can be used to locate the sites at which particular proteins associate with DNA; these include modified histones. First, formaldehyde is introduced into the cells of interest, which induces cross-linking between DNA and the proteins bound to it. The chromatin is then fragmented using sonication. Antibodies — which were grown beforehand to recognize the protein of interest (e.g., a histone with certain modifications) — are introduced to the chromatin solution. This is then used to separate the protein of interest from the remainder of the sample. The cross-linking is reversed so that the DNA linked to the protein can be separated and purified. Finally, the DNA fragments are sequenced with massively parallel next-generation sequencing. Intuitively, the positions of the DNA fragments in the genome should provide information about the presence of the protein, or histone modifications, at positions of the epigenome.

The DNA fragment reads are mapped to a reference genome to obtain a locus for each read. Bowtie is a tool that uses a Burrows-Wheeler transformation to perform this mapping efficiently [30]. Indeed, after preprocessing of the reference genome, it can map each read in $O(m)$ time and $O(n)$ space where m is the size of the read and n is the size of the genome.

Of course, there are inherent biases in a ChIP-seq experiment and, consequently,

a control experiment is typically performed in which ChIP-seq is run without an antibody. Many quality control metrics are typically evaluated on the outcomes of these experiments [29]. Given read counts, it is possible to generate signals corresponding to the true and to the control experiments; we refer to these respectively as ChIP and control mapped reads, or **tags**. (In fact, by estimating the length distribution of DNA fragments from the experiment, it is possible to extend each read to the expected length of the DNA fragment, and these extended tag counts can provide a more accurate signal track of binding positions [26].)

With the signals from the ChIP and control experiments, one can generate an **enrichment signal** that describes the degree to which histone modifications are observed at a particular location beyond what is expected by chance. We outline two methods for doing this [44]. Consider a small window r in which there are λ_r control tag counts and k_r ChIP (true) tag counts. One sensible enrichment signal value $s_{\text{fold}}(r)$ at r is an evaluation of the fold-change between the two:

$$s_{\text{fold}}(r) = \frac{k_r}{\lambda_r} \tag{2.1}$$

Alternatively, one can model the distribution of tags along the genome with a dynamic Poisson distribution and then compute a p -value for enrichment in r as $\Pr[\text{tag count} \geq k_r]$. Thus, an enrichment value $s_{\text{pval}}(r)$ can be obtained by taking the negative logarithm of this p -value:

$$s_{\text{pval}}(r) = -\log \Pr[\text{tag count} \geq k_r] \tag{2.2}$$

$$= -\log \sum_{i=k_r}^{\infty} \frac{\lambda_r^i e^{-\lambda_r}}{i!} \tag{2.3}$$

The enrichment signals input to ChromSMS can be computed like s_{fold} , s_{pval} , or another technique. The data used in this thesis is, unless otherwise stated, based on p -value signal tracks, i.e., computed like s_{pval} . A signal is said to describe one **chromatin mark**, which may be a histone modification, a particular protein, or some indicator of nucleosome positioning. We refer to the signal(s) within and around a

regulatory region as a **signal profile** of the region.

Although the described values offer a measure of enrichment and are directly useful for the goals of a tool like ChromSMS, it is often useful to perform “peak calling,” in which enrichment in regions of the epigenome is determined at a binary level. It is nontrivial to establish universal thresholds (e.g., fold-change or p -value cutoffs) for this, but there do exist methods, outside the scope of this thesis, for peak calling [26].

2.4 ChromHMM and Finding Patterns of Activity

A beautiful application of ChIP-seq data on histone modifications is ChromHMM [16]. The tool seeks to discover chromatin states on a genome-wide scale. A **chromatin state** is a biologically meaningful combination of chromatin marks (including histone modifications) that is present throughout various parts of the epigenome. For example, there are particular chromatin states that underlie promoter regions of the genome, particular states that underline enhancer regions, and so on. ChromHMM learns these states throughout the genome, without requiring existing annotations, with a multivariate Hidden Markov Model (HMM). In particular, given data on M chromatin marks, the HMM “binarizes” each 200 bp window of the genome such that the emission at a bin is a vector $[e_1, e_2, \dots, e_M]$ where $e_i \in \{0, 1\}$ indicates whether the bin is enriched with chromatin mark i . Each hidden state is then encoded as a vector $[X_1, X_2, \dots, X_M]$ where X_i is a Bernoulli random variable indicating the probability of mark i being present in the state. After training (e.g., with the Baum-Welch algorithm) to learn emission and transition probabilities, ChromHMM determines the maximum likelihood chromatin state behind all loci of the genome.

Applied to human T cells, this has led to the discovery of 51 distinct chromatin states [15]. They were manually characterized and assigned to classes of genomic elements; the states can be grouped into promoter-associated, transcription-associated, active intergenic, large-scale repressed, and repeat-associated classes. The tool has been enormously powerful in studying the epigenomic landscape of human and other genomes.

In the context of this thesis, the tool is important because of its ability to discover, *de novo*, regulatory regions of the genome. For example, active promoter regions can be discovered by searching for regions with promoter-associated chromatin states. Moreover, by applying ChromHMM to a concatenation of ChIP-seq data across cell types, it is possible to learn chromatin states jointly across cell types such that the state meanings are shared. These learned states can then be used to characterize the activity profiles of each regulatory region. In particular, this can be done by considering a vector for each region that describes the chromatin states underlying the region in the different cell types. With these vectors, as well as information about chromatin accessibility (DNase I hypersensitivity), it is possible to cluster the regions for a particular type of regulatory element based on their functional enrichments, or activity across cell types. We will refer to each of these clusters as **activity clusters**. Figure 2-3 shows 474,004 enhancer regions grouped into 47 activity clusters based on their profiles across 127 cell types.

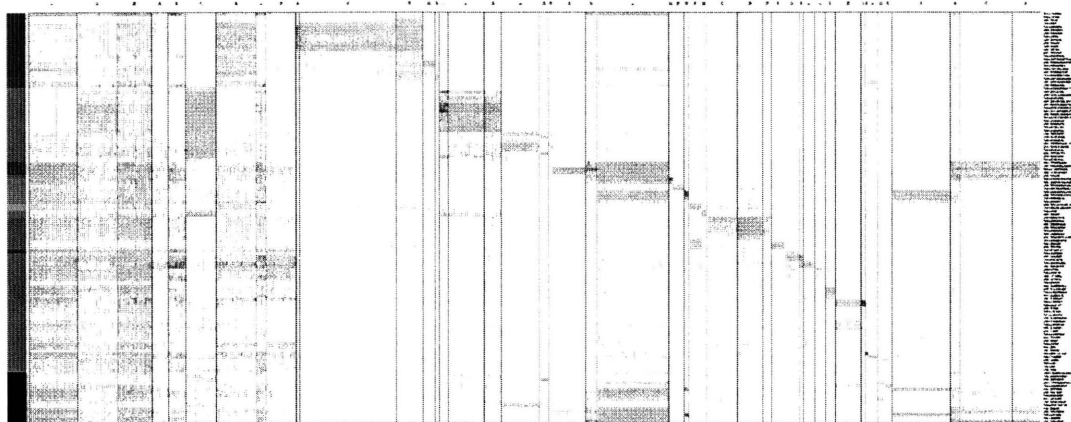


Figure 2-3: This heatmap shows functional enrichment of 474,004 (putative) enhancer regions across 127 cell types, and was generated with a 15-state model. The cell types, clustered, are along the rows and the enhancer regions, grouped into 47 clusters, are along the columns. This figure, as well as the activity clusters, were provided by Wouter Meuleman and were generated as part of the NIH's Roadmap Epigenomics project.

These activity clusters are a part of the input to ChromSMS, as described in Chapters 4 and 5.

Chapter 3

Past Work on Uncovering Patterns in the Chromatin Environment

A variety of work has been performed that is aligned with the objective of this thesis. Much of it is focused on developing and applying biological methods, or performing integrative analyses, to assess the relations, at a genome-wide scale, between transcription factors (TFs) and nucleosome occupancy [12, 18, 39, 41]. In this Chapter we discuss two computational analyses that have been performed and relate on some level, in their goal or approach, to this thesis. We also discuss, in detail, the methods behind a tool that served as inspiration for ChromSMS. But we begin by providing a survey of the biological results obtained from these studies.

3.1 Survey of Past Results Related to This Thesis

With over 100 known and characterized histone modifications, there is the potential for enormous combinatorial complexity. As expected, studies have uncovered large numbers of histone modification patterns [42]. However, many of the modifications in these patterns cluster together and are well-correlated, leading to the idea they act cooperatively and that there are a relatively small number of classes of patterns that uniquely determine expression [42]. Indeed, it has been demonstrated across many human cell lines and other species that just a few principal components can account

for much of the variance seen in histone modification patterns [43].

Moreover, in support of the theme of a limited number of major patterns, a study using the Clustered AGgregation Tool (**CAGT**) found a small number of spatial arrangements (or “shapes”) present in the patterns of histone modifications around the binding sites of many TFs [28]. While these shapes displayed substantial heterogeneity, they were shared by repressive and activating modifications alike. To recognize this, the authors constructed “meta-clusters” based on running CAGT on each mark on the full list of all binding sites of available TFs. They then ranked the clusters found for each modification based on the number of binding sites in the cluster and compared the rank order across modifications.

Perhaps most remarkably, the application of CAGT revealed substantial asymmetries in the patterns of histone modifications around binding sites for different TFs, different modifications, and different cell lines [28]. Indeed, this was apparent across modifications associated with promoters, enhancers, and gene bodies. Furthermore, it was found that all modifications showed asymmetry around a majority of binding sites proximal to transcription start sites as well as a majority distal from transcription start sites. These asymmetries were apparent in nucleosome positioning signals as well. The only chromatin mark displaying clear symmetry was one representing **DNase I hypersensitivity**, which indicates regions of chromatin that are loosened, exposed, and thereby accessible to the DNase I enzyme and to other proteins.

Such asymmetry was confirmed by a study using a technique known as protein interaction quantitation (**PIQ**) [37]. The authors of the technique define a pioneer TF to be one that binds to closed chromatin and, upon doing so, opens it and recruits other TFs to bind nearby. They discovered a set of directional pioneer TFs that display an ability to open chromatin asymmetrically — i.e., chromatin opens more strongly on one side of their motif than on the other.

3.2 Outline of Two Computational Techniques That Seek Patterns in the Chromatin Environment

3.2.1 Binarizing the presence of each modification

Consider the problem of finding disjoint sets of regulatory regions such that the regions in each set all share a common pattern for the way they carry epigenomic information in a cell type. While poorly defined (the words “share” and “pattern” are imprecise), it encompasses the general goal of ChromSMS. A first, natural approach to this problem is to simply assign each regulatory region a binary string in which each element represents either the presence or absence of a modification in the region. The string then characterizes a combinatorial pattern for the region, and regions with identical or similar patterns can be grouped together. Indeed, this was the approach taken in a comprehensive analysis of the combinatorial patterns of histone modifications in the human genome [42].

The study used a universal p -value threshold for establishing presence, and computed the p -value using tag counts with the background being a Poisson distribution. It is worth noting that this is nearly identical to the approach taken by ChromHMM in computing emission vectors, and that the information carried by the binary strings assigned to regions are, in some ways, similar to the information of the hidden chromatin states that ChromHMM learns. The approach is certainly suitable to an analysis of genome-wide combinatorial patterns and complexity — which was the goal of the study — but lacks an ability to uncover other patterns that may be of interest. Namely, one might be interested in patterns that are only noticeable on smaller scales (i.e., within a regulatory region). A binarization of enrichment is, in general, unable to capture differences in regions present at a high spatial resolution.

3.2.2 Modeling DNase magnitude and shape

One way to capture high-resolution differences is to consider the magnitude and shape offered by the full enrichment signal at regulatory regions. A recently developed com-

putational method that models the magnitude and shape of an enrichment signal is protein interaction quantitation (PIQ) [37]. PIQ evaluates a signal of DNase I hypersensitivity. Whereas each ChIP-seq experiment profiles only one TF to yield its binding locations, a DNase-seq experiment runs independently of TFs and yields a signal providing DNase profiles. From these DNase profiles and a list of motifs at candidate binding sites, PIQ identifies binding sites for hundreds of TFs at once. PIQ first selects candidate TF binding sites from its list of motifs, and then constructs a background model in which DNase-seq reads arise from a Gaussian process. After, PIQ iteratively builds a model of both TF binding estimation and expected DNase I hypersensitivity; it uses TF binding profile magnitudes and shapes to improve its expected DNase I hypersensitivity model by including its estimated binding events, and then refines its estimated binding events using the expected DNase I hypersensitivity model.

Compared with other methods that have the same aim of inferring TF binding from DNase data, PIQ shows high predictive accuracy. Although PIQ does not rely on ChIP-seq or incorporate histone modification data, it does demonstrate the potential power of a tool that can take advantage of the magnitude and shape encoded in histone modification enrichment signals.

3.3 CAGT: The Clustered AGgregation Tool

The Clustered AGgregation Tool (CAGT) is a methodology for discovering signal patterns from a chromatin mark (e.g., from histone modifications or nucleosome positioning) around functional elements like transcription start sites and the binding sites of TFs [28]. We describe the steps taken by CAGT to discover these patterns, and comment on some of its limitations.

3.3.1 The procedure

CAGT accepts as input a genome-wide enrichment signal for a particular chromatin mark as well as a set R of genomic regions. These regions are each made up of L

base pairs and may be, for example, windows around TF binding sites. Figure 3-1 summarizes the steps used by CAGT. The output consists of a set of clusters of regions. By plotting the aggregate signal over the regions in each cluster (accounting for “flipping”, as mentioned in the final step), the output may be used for discovering common signal patterns. Furthermore, the clusters may be used for performing desired analyses on sets of regions that share common signal profiles.

1. **Preprocessing and filtering.** For each region $r \in R$, CAGT extracts a signal vector at r of length $\frac{L}{W}$, such that each signal value is an enrichment averaged over a W bp bin. W parameterizes the resolution, and $W = 10$ is typically used. CAGT places regions into either a low or high “magnitude category.” The low category consists of regions whose signal variance falls below a threshold or whose 99th percentile signal value falls below a threshold, and these are considered to be regions with little functional signal. This category is filtered out, and only the regions in the high category proceed.
2. **Standardization and first phase of clustering.** CAGT standardizes the signal in each region to have zero mean and unit variance; the intention is for CAGT to capture patterns in signal shape rather than magnitude. Using the k -means or k -medians algorithm, CAGT clusters the regions based on their signal profiles. It uses as its distance function the Euclidean metric or one minus the Pearson correlation coefficient between two signals.
3. **Second phase of clustering.** Many of the k clusters from the above step have small internal variance but, because k is typically set to a large constant, there are many of them. Hence, CAGT performs agglomerative hierarchical clustering on the k clusters, and stops once the nearest distance between two clusters exceeds a parameter. In merging clusters, CAGT “flips” them as needed (i.e., if they are mirror images). This is done because the input regions may not have an explicit genomic directionality (e.g., if they are binding sites in enhancers). ChromSMS makes use of the algorithm for clustering with flips, and its mechanics are detailed in Section 4.3.2.1.

Figure 3-1: A summary of each of the steps taken by CAGT for discovering signal patterns [28].

3.3.2 Limitations of what it accepts and analyzes

CAGT has been successful in revealing unforeseen phenomena [28]. However, it has limitations in the type of data it accepts and the assumptions it makes about this

data. Here we recognize and describe three of these limitations.

Standardization. CAGT ignores the magnitude of signal profiles when uncovering patterns. Although it filters out regions with low signal intensity, it clusters on the remaining profiles after standardizing them. It may be the case that revealing information is encoded in their magnitudes. Differences in magnitude across signal profiles may display more detail about how they compare. Furthermore, variance within a signal profile may encode useful information. Even if signal profiles were grouped into a larger number of “magnitude categories” before standardization, and then clustering were done in each category separately, problems would persist. Indeed, it could be the case that differences across shape capture more of the variance between regions than differences in magnitude, and clustering within separate categories would fail to reveal this.

Requiring fixed-length regions. CAGT only clusters regions that have a fixed number of base pairs. But it may be that signal profiles that are similar — in both standardized shape and magnitude — are stretched over different length scales. This may present problems even over small differences in scales, and could amplify as the length of considered regions increases. Due to length-scale differences, a Euclidean or correlation-based metric may assign very different distances to signal profiles that appear similar.

Accepting only one mark. CAGT clusters profiles containing a signal from just a single chromatin mark. However, one may be interested in comparing regions based on information contained in signals from more than one mark. Moreover, it might be difficult to analyze results across many runs of CAGT on the same regions, each for a different mark. Including more than one signal in the profile of a region may reveal new combinatorial differences and complexities that individual runs cannot capture.

Despite many fascinating and surprising observations, there is still much to the complexities of the chromatin environment that is not understood or known. Because

of the incredibly vital and diverse functional role of regulatory elements, gathering a deeper and more complete understanding of the landscape around these elements is a pivotal task in epigenomics. Improving on the limitations of CAGT, as well as expanding what the tool offers to allow for more robust analyses of the chromatin environment around regulatory elements, can help improve our knowledge and understanding.

Chapter 4

ChromSMS

The Chromatin Signal Magnitude and Shape clustering tool

The primary goal in this thesis is to develop a computational tool that efficiently finds biologically meaningful clusters of regulatory regions based on the underlying information in their chromatin mark enrichment signals. We call this tool the **Chromatin Signal Magnitude and Shape clustering tool (ChromSMS)**.

ChromSMS requires input and provides output that is similar in structure to CAGT, and shares its fundamental approach. But the tool is more powerful and robust in that it offers solutions for the three major limitations discussed in Section 3.3.2, as well as for other more minor ones.

Making these improvements introduces considerable computational burden. Furthermore, a central goal from the outset was for ChromSMS to be easily used and modified by researchers. For this reason, we selected to implement it in Python; however, as the system and available features grew, this choice presented notable performance challenges. To operate efficiently and in a reasonable wall clock time, we spent substantial time refining ChromSMS's algorithms and heuristics, performance engineering its software, and introducing distributed computation when applicable.

In this Chapter we describe and analyze ChromSMS's design as a system, its algorithms, and its implementation.

4.1 Overview of its Design as a System

ChromSMS supports two primary modes. One, which we call the **varied-length** mode, is designed to capture spatial patterns; it standardizes signal profiles and explicitly allows for input regions of different sizes. The second, which we call the **multiple-mark** mode, is designed to robustly process all the available chromatin mark information in regulatory regions; although it requires signal profiles from a fixed-length window, the mode supports the consideration of information from many marks at once and also considers information from signal magnitude. These modes are discussed in Section 4.3.1.

For both modes, a run of ChromSMS can be divided into two main phases: (1) preprocessing and (2) clustering. Figure 4-1 outlines the flow of information between the major components of ChromSMS.

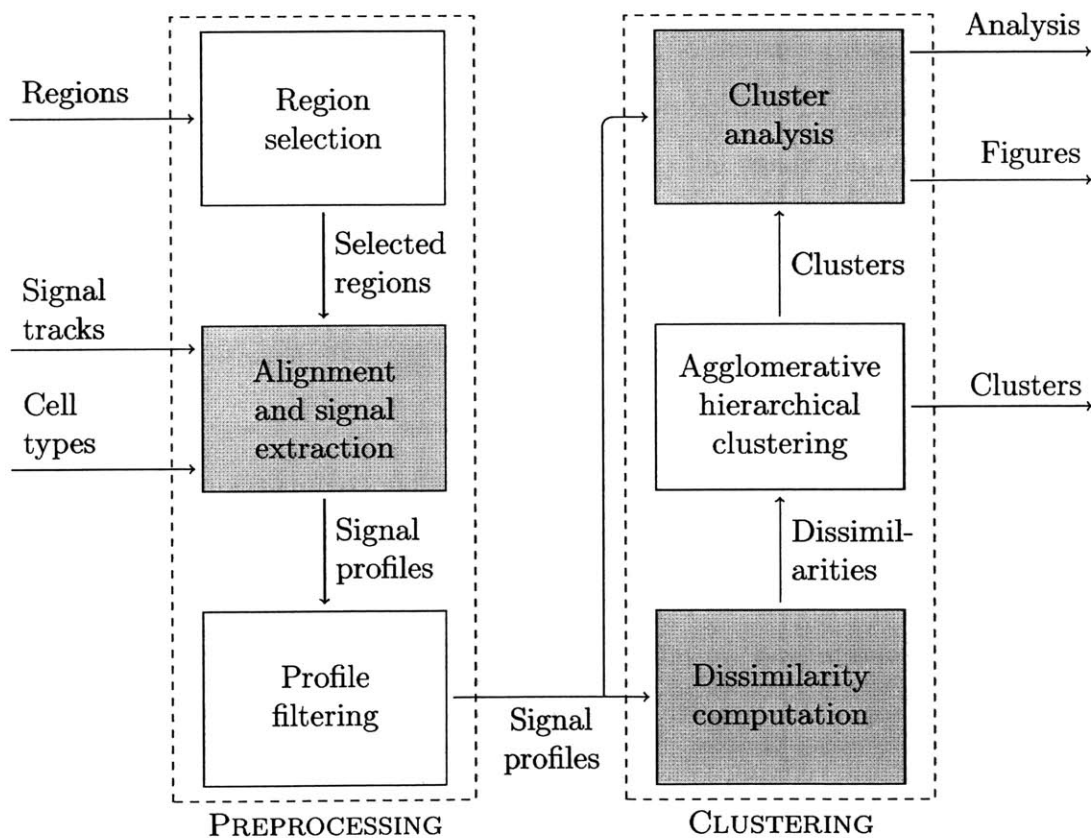


Figure 4-1: A model of the ChromSMS pipeline. The gray components require relatively high computational resources and are designed to be run in a distributed environment if desired.

Regions, signal tracks for desired chromatin marks, and a list of cell types are input to ChromSMS. Regions are first selected to ensure they are nonoverlapping, as described in Section 4.2.1. Then, their enrichment signals (or profiles) for the desired chromatin marks are extracted from the corresponding input cell types, and are aligned around a common anchor point, as described in Section 4.2.2. ChromSMS then filters out low complexity regions based on their profiles, as outlined in Section 4.2.3.

Next, using the signal profiles, pairwise dissimilarities are computed between regions in a manner that depends on the selected mode, as detailed in Section 4.3.1. ChromSMS efficiently performs agglomerative hierarchical clustering while accounting for the phenomenon of “flipping,” explained in Section 4.3.2. (In fact, in varied-length mode ChromSMS supports a centroid-based clustering approach that bypasses the need for pairwise dissimilarity computation. Similarly, like CAGT, it also supports a first phase of k -means or k -medians clustering, which was left out of the model in Figure 4-1. However, these options appear to add little value, as explained in Section 4.3.2.3, and were not used in generating the results presented in Chapters 5 and 6.)

Finally, ChromSMS computes inter- and intra- cluster statistics. It also generates the figures referenced in Chapter 5 to aid in visualizing the results. These steps are outlined in Section 4.3.3. Surprisingly, in practice this last step is usually the most computationally expensive component.

4.2 Preprocessing Input

4.2.1 Selecting regions to avoid overlap

It is possible that the input regions overlap, i.e., that some of them cover the same portions of the genome. For example, this may be due to two or more nearby TF binding sites whose regions are classified as separate regulatory elements. Additionally, as discussed in Section 4.2.2, ChromSMS chooses to consider a window of the

genome that surrounds, and is at least as large as, each input region in order to more fully capture signal patterns. This increases the number of overlapping regions, and can introduce ones even if there were originally none, because the surrounding windows may overlap. Overlapping regions present a problem because they would lead ChromSMS to consider identical portions of the genome more than once, and thus weigh these more heavily in a search for patterns. Moreover, information from the same portion of the same signals might become grouped into distinct clusters, making it difficult to conduct a comparison.

To correct for this, ChromSMS selects input regions such that the selected ones are representative of all the input regions and, even after considering surrounding windows, do not overlap. It extends each region, equally on each side, by enough so that the extended region covers 2,000 bp. Then it greedily places every extended region into exactly one set, each of which consists of a contiguous “stretch” of one or more overlapping extended regions. From each set, it selects a single region by considering at each region the DNase I hypersensitivity (DNase) signal¹, an indicator of chromatin accessibility and a strong marker of regulatory regions. In particular, it selects the region whose peak has the greatest magnitude. (If a DNase signal is not provided as an input signal track, some “peak summit” value must be provided for each input region.) Figure 4-2 displays an example of the result of this process.

4.2.2 Aligning regions and extracting signals

There is a considerable degree of computation that must be performed to obtain signal profiles for the selected regions.

First, recognize that there is a separate signal track provided for each chromatin mark. An additional input to ChromSMS is a list of one or more cell types whose signal information is desired. These cell types are often ones in which the input regions are considered “active,” as detailed in Section 5.1. When extracting the signal

¹Note that there is a separate DNase signal corresponding to each cell type. The signal considered here is the average signal across the input cell types. This follows the same procedure for signal extraction outlined in Section 4.2.2.

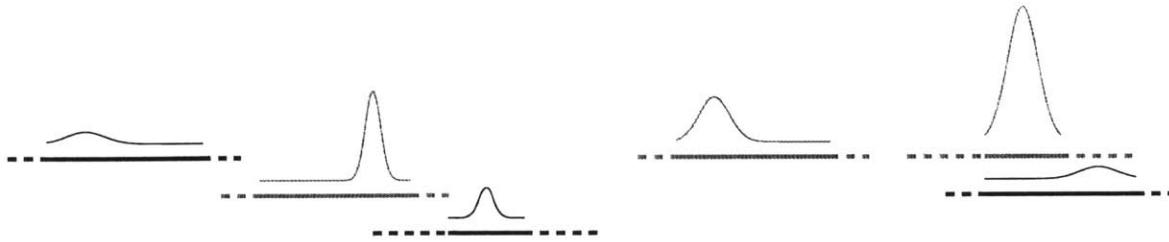


Figure 4-2: Six regions are shown ordered from left to right by genomic position. Input regions (thick solid line) are extended (dashed) on each side and every extended region is grouped into one set of overlapping regions. Then, the DNase I hypersensitivity signal (thin plot above region) is considered at each unextended region. Exactly one region from each set (red) is selected based on the magnitude of the regions' DNase peaks.

profile for a particular chromatin mark in a region, ChromSMS is actually computing an average of the signal tracks, across the provided cell types, for the mark in the region by taking the mean for each signal value. (For fold-change signal tracks, it also takes the log-transform of the mean — with an added pseudo-count — to reduce the effect of outlier values.)

But before all signals are extracted from the input tracks, the selected regions are first run through a process called **alignment**, which is also used by CAGT. To effectively compare two signal profiles, they must be spatially positioned around some common feature called an **anchor point**. This is typically an indicator of a TF binding site, such as a motif location. If an anchor point is not provided as input for each region, ChromSMS uses as an anchor point the position of a DNase peak. Similar to the process of selecting regions described in Section 4.2.1, ChromSMS averages the DNase signal across the desired cell types within each region and looks for the position of its greatest magnitude. The region is then shifted so that this anchor position is its center. In varied-length mode, a region whose coordinates are originally (s, e) is modified to have coordinates $(c - m, c + m)$ where $c \in [s, e]$ is the anchor position and $m = \max(e - c, c - s)$; this enables it to capture the landscape around the center while still keeping a sense of how large a window is intended to be captured. In multiple-mark mode, the region is modified to be a fixed 2,000 bp region centered at this anchor point.

At each region, ChromSMS extracts a signal profile from the input tracks. In

varied-length mode, this signal profile is made up of only the signal from the single desired chromatin mark. In multiple-mark mode, this signal profile consists of signals from all specified chromatin marks. The values in each 10 bp bin of the signal are averaged to one value; hence, ChromSMS operates at a 10 bp resolution.

It is worth noting that this extraction is more computationally intensive than it may seem. Indeed, consider 10,000 selected input regions in a multiple-mark run with 10 marks and 100 specified cell types. If the signal tracks are provided at a single bp resolution, given that each region consists of 2,000 bp this extraction requires reading from disk 20 billion values in total. Such input sizes are not atypical. However, the process can be easily distributed, and ChromSMS does perform it in a distributed manner if desired.

4.2.3 Filtering signal profiles with low complexity

Following a technique from CAGT, ChromSMS filters out regions based on their signal profiles. In particular, it seeks to remove regions whose signals indicate that they have little functional information. Before filtering, in each signal profile it interpolates values missing from the signal tracks. In multiple-mark mode, the signal is considered to be the concatenation of the signals from the various marks.

At each region ChromSMS first considers the signal value at the q^{th} percentile and filters out the region if this value falls below a threshold T_1 . It then considers the variance of the signal and filters out the region if this variance falls below a threshold T_2 . The parameter values are found empirically and depend on the mode and type of signal track (e.g., fold-change or p -value based). For example, for p -value based tracks with eight marks in multiple-mark mode we find that $q = 95$, $T_1 = 1.5$, and $T_2 = 0.1$ appear to work well.

Finally, in varied-length mode ChromSMS standardizes each signal to zero mean and unit variance before continuing.

4.3 Clustering Signal Profiles

At the core of ChromSMS are its techniques for clustering regions based on their signal profiles. We primarily work with linkage-based hierarchical clustering, and explain why in Section 4.3.2.3, although ChromSMS supports other methods. We begin by defining a dissimilarity measure between signal profiles.

4.3.1 Dissimilarity measurements between two profiles

We define separate measures depending on whether ChromSMS is run in varied-length or multiple-mark mode. In each we desire an *ad hoc* measurement designed for the explicit purpose of the mode.

4.3.1.1 Varied-length mode

In varied-length mode, our primary goal is to capture spatial patterns in profiles. To this end, the signals are standardized. When comparing two regions, we wish for ChromSMS to be agnostic to the length scale of regions and also to relative phase offsets in patterns between the two.

Even if two signals were to have the same length and thus the Euclidean distance could be applied, this would be a poor metric for the problem because it linearly compares values between two signals along the spatial axis; small scaling and phase offsets along the axis could produce large measured differences. This is exacerbated in the context of this problem because of considerable noise within small regions, potential difficulty in aligning signals, and the biological possibility that similar spatial patterns are present over different length scales.

A very similar problem appears in speech recognition, where one must compare spoken waveforms recorded along the time axis to determine whether they represent the same phrase. A popular solution in the field has been the Dynamic Time Warping (DTW) technique [1, 25]. DTW is a dynamic programming algorithm that warps two signals non-linearly along the time axis to construct a mapping between them, and does so in a way that minimizes a measure of the cost of the warping. Interestingly,

it is very closely related to techniques in computational biology for global sequence alignment [33]. ChromSMS applies the DTW algorithm to our problem of interest.

Suppose we wish to compare two signals x and y of dimension m and n respectively and assume WLOG that $m \geq n$. We begin by filling an $(m + 1) \times (n + 1)$ warping matrix D . We initialize $D(i, j) = \infty \forall 0 \leq i \leq m, 0 \leq j \leq n$ except $D(0, 0)$, which we initialize to $D(0, 0) = 0$. Then we apply the recursion

$$D(i, j) = \{x(i - 1) - y(j - 1)\}^2 + \min \left\{ \begin{array}{l} D(i - 1, j) \\ D(i, j - 1) \\ D(i - 1, j - 1) \end{array} \right\} \quad (4.1)$$

for all $1 \leq i \leq m$ and $1 \leq j \leq n$. The left summand is the cost we define between two signal values and makes the squared Euclidean distance a special, constrained case of this recurrence. The value $D(m, n)$ is the minimum cost of warping to match the two signals and is often referred to as the DTW measurement.

Performing this computation takes $O(mn)$ time. Given our use-case, in which we must compute it for every pair of input regions, this is too slow. Rather than filling the entire matrix, we can constrain our computations to a narrow “band” of constant width w centered around the diagonal of the matrix. In particular, for each row $1 \leq i \leq m$, we let the middle column on the diagonal be $r = \lfloor \frac{i \cdot n}{m+1} \rfloor + 1$ and only compute $D(i, j)$ for $\max(1, r - w) \leq j \leq \min(n, r + w)$. This has the potential to miss the global optimal warping, but improves the running time to $O(m + n)$. This is the approach followed by ChromSMS. We find empirically that $w = 20$ works well in balancing runtime and quality of approximation for the types of data processed.

Finally, note that the measurement $D(m, n)$ — or $\sqrt{D(m, n)}$ to make a parallel to Euclidean distance — has the consequence that longer signals are further apart than shorter ones because there are more dimensions over which warping cost is added. In the context of agglomerative clustering, this measurement could conceivably lead to a small number of large clusters containing short regions and a large number of small clusters containing large regions; this would fail to capture the patterns we seek. Hence, we wish to normalize out the length of the discovered warping path.

ChromSMS fills another matrix to determine the length of this path; let the length be l . Ultimately, ChromSMS reports as the dissimilarity between x and y the value $\sqrt{\frac{D(m,n)}{l}}$.

We note that, despite the added complexities, a centroid-based clustering approach could still be applied with this dissimilarity measurement. We could compute a centroid for a cluster by computing the median region length in the cluster, warping (e.g., by resampling) all the signals for the regions in the cluster to this median length, and then computing the mean across these resampled signals. Although ChromSMS supports a centroid approach for single-mark fixed-length regions, we experimented little with this approach for varied-length regions. Computing these dissimilarities involves relatively high constant factors; in a centroid-based approach we cannot precompute the required dissimilarities, so they must be computed on-the-fly, which adds considerable computational burden. In this context there are no obvious benefits over linkage-based approach.

4.3.1.2 Multiple-mark mode

In multiple-mark mode, our goal is for ChromSMS to compare regions based on the signals of multiple and possibly diverse chromatin marks, considering both magnitude and spatial patterns. A naïve approach — simply taking a Euclidean distance across a concatenation of signals for the different marks — is ill-suited for this goal. Among other reasons, it could capture more strongly differences in magnitude than in shape.

In a given region some chromatin marks hold more information content than others. For example, based on their known associations, in active enhancers H3K4me1 is usually more enriched than H3K27me3 and, therefore, will usually contain more functional signal [42]. In comparing two enhancers, differences in the signal for H3K27me3 are more likely to be the result of noise, and thus less interesting, than differences in H3K4me1. It is therefore important that the measure not weigh equally the differences in signal profiles among distinct marks; that is, it should account for the information content we observe.

Furthermore, in considering the large amount of information across multiple marks,

there is a stronger likelihood of encountering noise and of encountering signals that are not well-aligned relative to each other; the measure should not be adversely affected by these possibilities. Hence, we desire a dissimilarity measure that is robust to noise, agnostic to small signal offsets, accounts for the information content of different modifications, and does not give disproportionate weight to either magnitude or shape.

The central idea behind our dissimilarity measure is to divide a signal into some number of bins and consider each independently at first. There are two reasons for this. First, it allows us to separate magnitude and shape when comparing regions. We can use a bin’s mean magnitude to compare magnitude, and can use a Euclidean distance to effectively compare shape between two bins. Second, we can incorporate a bin’s magnitude into a measure of its information content, and thereby weigh measurements between two bins in a measurement calculation between two regions. The size of the bin, as well as the sliding length between bins, are parameters that we can empirically tune to yield robust results.

Notation. We first introduce notation that will help in defining the dissimilarity measure. ChromSMS accepts as input regions of fixed length $L = 2000$ bp. Let each region r contain signals $\{s_1^r, s_2^r, \dots, s_m^r\}$ where m is the number of incorporated histone marks. ChromSMS creates k (possibly overlapping) bins, of size w , for each signal; let s_{ij}^r be the length- w signal inside bin j of s_i^r . Let $\text{mag}(s_{ij}^r)$ be the mean magnitude of s_{ij}^r . For two signals s_{ij}^x and s_{ij}^y in corresponding bins of regions x and y , let the fold-change difference in magnitude between them be

$$d_{\text{mag}}(s_{ij}^x, s_{ij}^y) = \frac{\max[\text{mag}(s_{ij}^x), \text{mag}(s_{ij}^y)] + \epsilon}{\min[\text{mag}(s_{ij}^x), \text{mag}(s_{ij}^y)] + \epsilon} - 1 \quad (4.2)$$

where $\epsilon = 0.001$ is a pseudo-count. Let \widetilde{s}_{ij}^r be the signal s_{ij}^r standardized to zero mean and unit variance. Then we define $d_{\text{shape}}(s_{ij}^x, s_{ij}^y)$ to be the Euclidean distance between \widetilde{s}_{ij}^x and \widetilde{s}_{ij}^y .

Information content of bins. For each mark i and bin j , we estimate a probability density function $f_{ij}(\alpha)$, taken across all input regions, that approximates the likelihood of obtaining a signal in the bin where $\text{mag}(s_{ij}^r) = \alpha$. It estimates this empirically by computing a histogram over the values $\text{mag}(s_{ij}^r)$ for all input regions r in which it places the values into fixed-width histogram bins. Let $\text{Pr}_{ij}(\alpha)$ be $f_{ij}(\alpha)$ integrated over the width of the histogram bin into which α falls. Then we can measure the information content of the signal describing mark i in bin j of region r as

$$I_{ij}^r = -\log \text{Pr}_{ij}(\text{mag}(s_{ij}^r)) \quad (4.3)$$

It is also possible to work with a measure of the information content of mark i , obtained by summing I_{ij}^r over the bins j in mark i . However, working with a separate measure of information content for each individual bin allows for the more fine-grained possibility that different portions of a mark hold different degrees of information content.

Measuring dissimilarity between two bins. Before constructing a dissimilarity measure between two regions, we first construct one between two signals s_{ij}^x and s_{ij}^y . Recall that we seek to assign similar weight to differences in magnitude and in shape. Working directly with $d_{\text{mag}}(s_{ij}^x, s_{ij}^y)$ and $d_{\text{shape}}(s_{ij}^x, s_{ij}^y)$ is difficult because these values may take on different distributions across the pairs of signals. Thus, in each bin j of each mark i we standardize the values $d_{\text{mag}}(s_{ij}^x, s_{ij}^y)$ and $d_{\text{shape}}(s_{ij}^x, s_{ij}^y)$. We do this by considering these values across all pairs of regions. For each bin j of each mark i let $M_{ij} = \{d_{\text{mag}}(s_{ij}^x, s_{ij}^y)\}$ over all pairs of regions x and y and, similarly, let $S_{ij} = \{d_{\text{shape}}(s_{ij}^x, s_{ij}^y)\}$. We can standardize these dissimilarities as $z_{\text{mag}}(s_{ij}^x, s_{ij}^y) = \frac{d_{\text{mag}}(s_{ij}^x, s_{ij}^y) - \text{E}[M_{ij}]}{\sigma(M_{ij})}$ and $z_{\text{shape}}(s_{ij}^x, s_{ij}^y) = \frac{d_{\text{shape}}(s_{ij}^x, s_{ij}^y) - \text{E}[S_{ij}]}{\sigma(S_{ij})}$. With these standardized dissimilarities computed, we can sum them to compute a dissimilarity between the two signals s_{ij}^x and s_{ij}^y that gives similar relative weight to magnitude and shape.

Measuring dissimilarity between two regions. Suppose we wish to compute a dissimilarity between two regions x and y . For each bin j in each mark i , we assign a relative weight of $I_{ij}^x + I_{ij}^y$ to the dissimilarity between the two bins:

$$[I_{ij}^x + I_{ij}^y] \cdot [z_{\text{mag}}(s_{ij}^x, s_{ij}^y) + z_{\text{shape}}(s_{ij}^x, s_{ij}^y)] \quad (4.4)$$

Summing over the bins and normalizing for the total weight, we measure the dissimilarity $d(x, y)$ between the two regions as

$$d(x, y) = \frac{\sum_{i=1}^m \sum_{j=1}^k [I_{ij}^x + I_{ij}^y] \cdot [z_{\text{mag}}(s_{ij}^x, s_{ij}^y) + z_{\text{shape}}(s_{ij}^x, s_{ij}^y)]}{\sum_{i=1}^m \sum_{j=1}^k [I_{ij}^x + I_{ij}^y]} \quad (4.5)$$

We find empirically that using a bin size of $w = 500$ bp and a sliding length of 500 bp between bins produces robust results. Although this may appear to be a large bin size — and we did indeed experiment with smaller sizes — it reaches a strong balance between the two extremes of having many small bins and of having one large bin. With $L = 2000$ bp, this means there are four bins describing each chromatin mark and they are non-overlapping.

ChromSMS performs these computations and returns $d(x, y)$ as the dissimilarity between regions x and y .

4.3.2 Implementing agglomerative hierarchical clustering

ChromSMS implements agglomerative hierarchical clustering as its primary method for clustering. Because of the unique need to “flip” regions (see Section 4.3.2.1) while clustering, which is unsupported by existing tools, we implement the clustering algorithm from scratch. Recall that ChromSMS first computes dissimilarities between every pair of regions; we call $d(x, y)$ the dissimilarity between regions x and y .

ChromSMS supports the commonly used average-, single-, and complete-linkage distance measures. The average-linkage distance between two clusters X and Y is

defined to be

$$D(X, Y) = \frac{\sum_{x \in X, y \in Y} d(x, y)}{|X||Y|} \quad (4.6)$$

When reporting results we use the average-linkage distance. The single-linkage distance is $\min_{x \in X, y \in Y} d(x, y)$ and the complete-linkage distance is $\max_{x \in X, y \in Y} d(x, y)$.

Computing the distance between two clusters using these measures takes $O(|X||Y|)$ time. However, if we track pairwise distances between clusters, it is possible to perform updates while merging in just $O(1)$ time. Indeed, suppose we have three clusters X , Y , and Z . If we merge X and Y to form $X \cup Y$, then the average-linkage distance from this merged cluster to Z is

$$D(X \cup Y, Z) = \frac{\sum_{x \in X \cup Y, z \in Z} d(x, z)}{|X + Y||Z|} \quad (4.7)$$

$$= \frac{\sum_{x \in X, z \in Z} d(x, z) + \sum_{y \in Y, z \in Z} d(y, z)}{(|X| + |Y|)|Z|} \quad (4.8)$$

Because we have already computed and stored $D(X, Z)$ and $D(Y, Z)$, we can rewrite this as

$$D(X \cup Y, Z) = \frac{|X||Z| \cdot D(X, Z) + |Y||Z| \cdot D(Y, Z)}{(|X| + |Y|)|Z|} \quad (4.9)$$

$$= \frac{|X| \cdot D(X, Z) + |Y| \cdot D(Y, Z)}{|X| + |Y|} \quad (4.10)$$

This calculation can be performed in just $O(1)$ time, and ChromSMS exploits this while merging. ChromSMS performs similar calculations for single- and complete-linkage.

ChromSMS stores a matrix of distances between pairs of clusters and tracks the minimum of each row, which is effectively the nearest neighbor of each cluster. Thus, it requires $\Theta(n^2)$ space. Because there are $O(n)$ rows, this allows it to find the closest pair of clusters for merging in $O(n)$ time. It updates row minimums when they become invalidated due to a merge. Because each update takes $O(n)$ time and there may be

$O(n)$ on each merge in the worst-case, a merge may require $O(n^2)$ time. Given that there are $O(n)$ merges, this technique takes $O(n^3)$ time in the worst-case, which is no better than a brute-force technique for hierarchical clustering. As it turns out, however, this approach has been studied in the literature, where it has been called the “nearest neighbor heuristic” [14]. It performs strongly in practice when compared against other algorithms. Indeed, when n , the number of regions, is on the order of tens of thousands, ChromSMS’s implementation of this approach runs in a reasonable time. (In a test we clustered 30,232 enhancer regions, counted after region selection and profile filtering, by merging down to a single cluster. The clustering took 33.7 minutes.)

4.3.2.1 Flipping regions to maintain orientation

In general the genomic regions that are input to ChromSMS do not have an explicit directionality. (One that would, for example, is a promoter of a known annotated gene.) However, the signal profiles for a region may reveal an implicit directionality that must be considered when clustering. Failing to account for this could reveal symmetric patterns when in fact asymmetry is present. Figure 4-3 displays this phenomenon.

This authors of CAGT recognized this and devised a modification to hierarchical clustering that solves it; ChromSMS uses the same idea. In particular, for every pair of regions $\{x, y\}$, ChromSMS computes the dissimilarity $d(x, y)$ between x and y as well as the dissimilarity $d(x, y')$, where y' represents the region y with its signal profiles flipped about its center. Furthermore, rather than simply storing an $n \times n$ distance matrix between clusters during the hierarchical clustering process (where n is the number of regions), ChromSMS stores an $n \times n \times 2$ matrix. If F is this matrix, F_{ij0} stores the distance $D(X, Y)$ between clusters X and Y , and F_{ij1} stores the distance $D(X, Y')$ where $Y' = \{y' \mid y \in Y\}$ is the “flipped” version of Y . In this way, the global minimum of F yields the closest pair of clusters while accounting for the phenomenon of flipping.

Consider the merge of two clusters X and Y . If $D(X, Y') < D(X, Y)$, ChromSMS

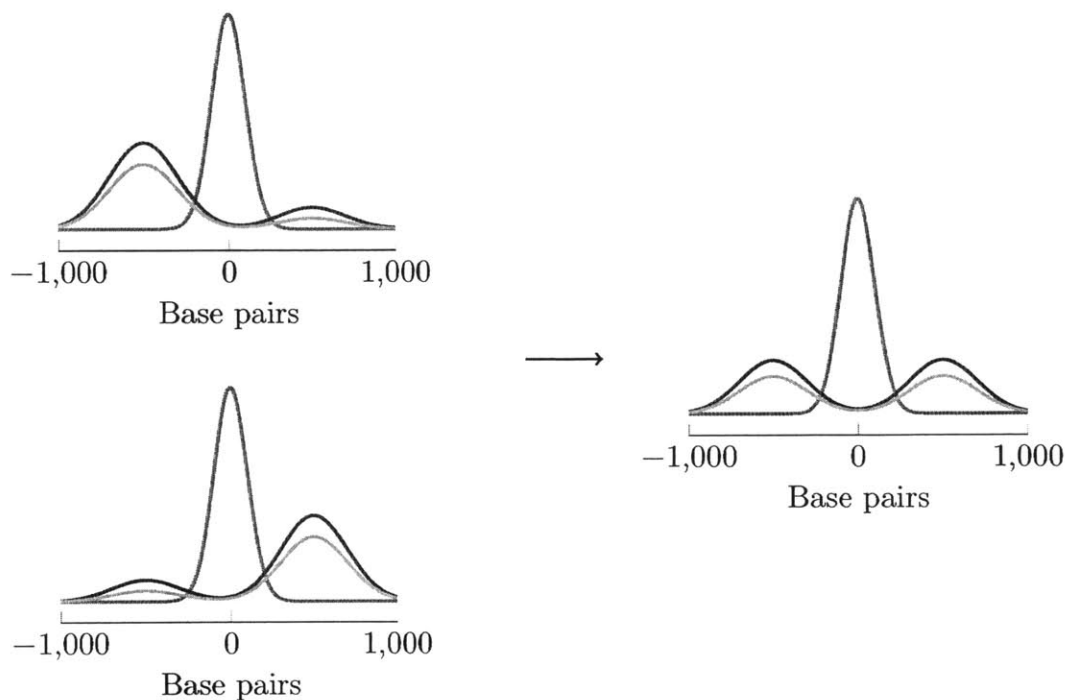


Figure 4-3: An example of two signal profiles with three marks each (left) that appear different, but should actually be considered identical. Because their regions may not have an explicit directionality, they may demonstrate identical patterns. When combined and averaged naïvely, they yield a symmetric profile (right).

in fact merges X with Y' . That is, the appropriate elements F_{ij0} store the distances from each cluster to $X \cup Y'$ and the elements F_{ij1} store the distances to $X \cup Y$. This ensures that flipped signal profiles are carried through the merging process and that orientation within clusters is maintained.

4.3.2.2 Terminating the merging process

One important problem in agglomerative hierarchical clustering is that of determining when to stop merging clusters. ChromSMS supports two simple methods — a parameter defining the maximum distance between clusters permitted for merging, as well as a parameter that defines the number of clusters — and can stop merging when either threshold is crossed.

However, acceptable values for these parameters are highly dependent on the dissimilarity measurements, which in turn depend on the number of regions, the

number of chromatin marks being considered, the signal track type, and other aspects of the input. Thus, we experimented with techniques that are more robust to changes in these features.

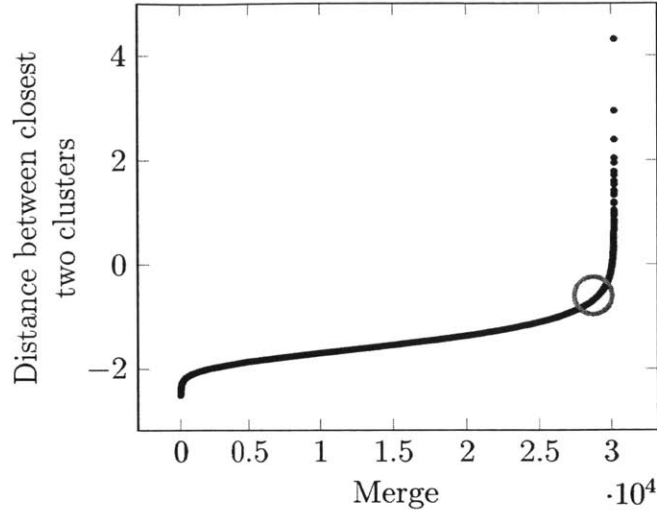


Figure 4-4: We clustered a set of 30,232 enhancer regions (counted after region selection and profile filtering) in multiple-mark mode by merging down to a single cluster. This points shows, at successive merges, the distance (or average pairwise dissimilarity) between the two closest clusters. The red circle indicates a region at which it might be reasonable to terminate merging. Although each run produces a curve over different scales, the general shape — and notably the pronounced “elbow” — is typical regardless of the input.

Figure 4-4 shows a typical plot of the distance between the closest two clusters at each successive merge. The curve shows a clear “elbow”, at which point the clusters no longer display the variance in the data, and instead merging is performed between clusters that are too dissimilar. One approach to finding this point is to estimate a second derivative and locate where it reaches its maximum. However, while mathematically valid, we find empirically that this is not as robust as we desire; this may be due in part to irregularities, particularly on small numbers of input regions.

ChromSMS uses a heuristic that works as follows: We approximate first derivatives with a central difference and second derivatives with a 2^{nd} order central difference. Let n be the number of input regions. We first compute H , the first derivative of the middle 25% of the curve. That is, H is the first derivative of the portion of the curve after $0.375n$ merges and before $0.625n$ merges. (We only consider the possibility of

terminating the merging process after $0.625n$ merges have been performed.) At each merge we consider the ratio of the first derivative of the plot at that merge to the value H . We terminate the merging process if this ratio exceeds a parameter α_1 . The particular value still has a dependency on the mode and signal track type; as an example, we find that, in multiple-mark mode with eight chromatin marks from p -value signal tracks, using $\alpha_1 = 7.0$ works well across very different input sizes.

(We also require that the minimum distance at that merge exceeds a parameter α_2 and that the second derivative at that merge exceeds a parameter α_3 ; however, these two thresholds are aimed at highly irregular curves, usually resulting from small inputs, and do not typically prolong the merging process beyond the crossing of the parameter α_1 . We find that $\alpha_2 = -2.0$ and $\alpha_3 = 0.0$ work well in multiple-mark mode.)

4.3.2.3 Other relevant clustering algorithms

Hierarchical clustering is appropriate in this context for a few reasons. First, the dissimilarity measures are not metrics (i.e., they do not satisfy the triangle inequality in general or, in the case of multiple-mark mode, non-negativity); unlike some other approaches, hierarchical clustering is suitable for measures that are not metrics. Furthermore, the combinatorial complexity introduced by multiple-mark mode may mean that the optimal clustering has very many clusters, making it difficult to pick this number in advance. It may also be helpful for analyses to visualize a hierarchy of clusters. Although hierarchical clustering is thought to be slower than some other algorithms, its runtime is not poor in this context because we can precompute pairwise dissimilarities in a distributed setting and the “nearest neighbor heuristic” applied by ChromSMS performs well.

However, we do not claim that it performs better, in terms of the general “quality” of its clusters, than other clustering approaches. In this thesis we spent a large fraction of time on designing and testing appropriate dissimilarity measures, as detailed in Section 4.3.1, compared to experimenting with other clustering algorithms. As one popular machine learning textbook observes, “Specifying an appropriate dissimilarity

measure is far more important in obtaining success with clustering than choice of clustering algorithm” [22, p. 506].

The k -means algorithm is not suitable here because it is based on a Euclidean distance [10, p. 427-428]. However, a k -medioids algorithm may achieve success with the dissimilarity measures used by ChromSMS. A “soft” clustering approach, based on an EM for Gaussian mixtures algorithm, may also be worth experimenting with. These would, of course, have to be modified in a way that accounts for the phenomenon of “flipping.”

4.3.3 Reporting results and creating visualizations

After generating clusters, ChromSMS computes and reports measures to validate the clustering. These include measurements of intra-cluster variance, as well as a statistic that can indicate whether each pair of clusters represents, at a significant level, a different pattern.

We refer to each of the clusters output by ChromSMS as a **shape**. This term holds even though such a cluster may incorporate more than just spatial patterns (e.g., magnitude as well). We use this term to ensure a distinction between ChromSMS’s clusters and activity clusters, the latter of which serve as input to ChromSMS, as discussed in Sections 2.4 and 5.1. Furthermore, we refer to the regions that pass the signal profile filtering step, described in Section 4.2.3, as forming the *Passed-Set*. We refer to the shape with the most regions as *Shape 0*, the shape with the second most regions as *Shape 1*, and so on.

The portion of this step that requires the most resources, and arguably yields the most helpful output, generates plots to visualize the patterns in the shapes. These include standard aggregate plots that account for the flipping phenomenon. In varied-length mode, prior to plotting, the signals in each shape are resampled to the median length of regions in the shape. The plot for each shape shows, averaged across the regions in the shape, the signal profiles that were used in clustering (i.e., signals that were generated by averaging across the specified input cell types). But ChromSMS also generates, as a control, plots for each shape in which the signals for the regions

in the shape are averaged across the tracks from all cell types that were not specified as input. Furthermore, for each shape and chromatin mark, ChromSMS generates a heatmap that displays that chromatin mark's signal, averaged across the regions in shape, for each individual cell type for which a signal track is available. Examples of these plots are shown in Chapter 5 and Figure B-1.

Due to the inclusion of the control aggregate plots and the heatmaps, generating the signals to display in these figures requires substantial computational resources because it requires reading from every available signal track for the input chromatin marks. It takes a time strictly slower than the signal extraction step outlined in Section 4.2.2. Fortunately, like that step, this is simple to parallelize; ChromSMS can compute the signals to display more quickly in a distributed setting.

4.4 Implementation Details and Performance

A goal from the outset of development was for ChromSMS to be accessible to users, both in terms of using and modifying it. Consequently, the tool is written in the Python programming language. It has, in total, about 13,000 lines of code.

Writing the tool in Python has presented a number of performance challenges. To this end, we make extensive use of the NumPy package for scientific computing; its support for large matrices helps tremendously in efficiently processing signals. We also have extensively profiled the tool's performance and rewritten key portions to take advantage of caching and vectorization. We also rewrote performance bottlenecks in Cython, a language that allows the invocation of a C routine, which has helped. These portions include the dissimilarity measure that ChromSMS uses in varied-length mode.

Beyond performance engineering the serial execution of ChromSMS, we also make extensive use of distributed computing on a computing cluster for expensive operations that fit the MapReduce model. These operations include reading signals and computing pairwise dissimilarity measurements. ChromSMS was written to take advantage of a distributed environment, and such operations would be prohibitively

expensive without access to a cluster. In particular, ChromSMS is written to use the Platform Load Sharing Facility (LSF) environment; however, its interface should make it easy to adapt the tool to other environments.

In evaluating ChromSMS, we are, for the most part, not concerned with optimizing its time or space requirements, as long as these are “reasonable” for analyses, assuming access to a computing cluster. And we believe they both are: An 8-mark run provided with on the order of 10,000 input regions can be performed in under six hours with 100 available nodes in which one coordinating node requires roughly 10 GB of RAM and the others require roughly 1 GB. Our primary focus in evaluating ChromSMS is a consideration of its ability to characterize and better understand the epigenome at a high-resolution. We describe our findings in Chapters 5 and 6.

Chapter 5

A Review of Discovered Shapes and Analyses Within These Shapes

In this Chapter we apply ChromSMS to data produced by the National Institutes of Health's Roadmap Epigenomics project [7]. We study directly its output by visualizing the aggregate plots it presents for each shape. (Recall from Section 4.3.3 that a **shape** is a cluster output by ChromSMS.) We also conduct various analyses of these results. These include a look at functional enrichment in regions within shapes and studies of the relation between discovered patterns and genomic sequence.

We begin, in Section 5.1, by describing in more detail the input data used for conducting our analyses. In Section 5.2 we present results from ChromSMS run in varied-length mode. We present an overview of results from multiple-mark mode in Section 5.3 and for the remainder of the Chapter we use multiple-mark mode.

We then focus on studies of these results that involve the regions within individual shapes. In Section 5.4 we present an analysis of directionality in promoters and enhancers. Then, in Section 5.5 we outline results on functional enrichment of regions in particular shapes. In Section 5.6 we describe observed correlations between signal profile patterns and GC-content. We conclude these analyses in Section 5.7 by describing how we search for enrichment of motifs in individual shapes. Finally, we discuss what we find in Section 5.8.

In Chapter 6 we perform studies that involve regions from multiple shapes at once.

5.1 Data Provided as Input to ChromSMS

The regulatory regions used as input are delineated and annotated as part of the Roadmap Epigenomics project, and are provided to us [7]. They are regions that are found to be regulatory in *both* a 5-mark 15-state model from observed data and a 12-mark 25-state model from imputed data, trained across 127 epigenomes, and whose DNase I hypersensitivity peak has a p -value satisfying $-\log_{10} p \geq 10$. Of these regions, there are 33,086 promoters (1% of the genome) and 474,004 enhancers (5.92% of the genome). The regions in each category are grouped into activity clusters, as explained in Section 2.4. For each activity cluster we are given the fraction of its regions that are active in each cell type. We say that an activity cluster is **active in a cell type** if this fraction exceeds 0.5.

We are given 22 promoter activity clusters and, of these, we identify 17 that are active in at least one cell type. Similarly, we are given 47 enhancer activity clusters and identify 41 that are active in at least one cell type. Thus, when producing and presenting results we use these 17 promoter activity clusters and 41 enhancer activity clusters.

The signal tracks with which we work are p -value based, as explained in Section 2.3, and are also provided by the Roadmap project. They are **imputed** tracks, and thus offer data for combinations of cell types and chromatin marks for which **observed** data is not available. We make use of signal tracks from 127 cell types for the marks DNase, H3K27ac, H3K27me3, H3K36me3, H3K4me1, H3K4me3, H3K9ac, and H3K9me3.

For a particular activity cluster, we provide as input to ChromSMS the cell types in which the activity cluster is active. Recall from Section 4.2.2 that ChromSMS computes an average signal track for each chromatin mark, taken across the input cell types. One benefit of this process, which is crucial in producing results, is that it helps noticeably in reducing noise. A signal profile of a region, taken from one cell type, may contain ample noise; a new signal profile, created from many cell types, can better display information in the region.

5.2 Review of Output from Varied-length Mode

To demonstrate the output of ChromSMS in varied-length mode, we display and discuss the results from a particular promoter activity cluster. We choose a promoter cluster because of the wide distribution of lengths they take on. Indeed, the lengths of the regions across all promoter activity clusters has a standard deviation of 901 bp, but just 270 bp in enhancers. The regions in the promoter activity cluster that we use as input have, as their lengths, a median of 341 bp and a standard deviation of 555 bp, and they span the range [69, 7541] bp. The activity cluster is active in 32 cell types that encompass, for the most part, embryonic stem cells, induced pluripotent stem cells, and various cells derived from embryonic stem cells. There are 1,022 regions in the activity cluster; ChromSMS selects 839 of them and filters out none based on their signal profiles.

Figure 5-1 shows aggregate plots output by ChromSMS. We display five shapes (as well as the *Passed-Set*) that together account for 66.74% of all the selected regions. We immediately observe that ChromSMS is able to capture signal profile patterns that appear concrete and distinct, although we quantify their differences more rigorously in Section 6.1.

We find that many of the larger shapes have little complexity or interesting information. This is likely a consequence of the small windows we are providing as input; for example, this phenomenon is prominent in *Shape 1* of Figure 5-1, which has a median region length of just 240 bp.

To offer ChromSMS more information about each region, we can try the following: we use the same procedure described in Section 4.2.2 for extracting and aligning signals, but additionally extend each region by 500 bp on each side. Our output with this input are shown in Figure 5-2. It is clear that expanding the window around each region yields more diverse patterns and a higher degree of complexity. Furthermore, we notice immediately the same asymmetric patterns that were identified by CAGT and are characteristic of promoters [28]. The results support the choice to use 2,000 bp as a fixed-size window in multiple-mark mode.

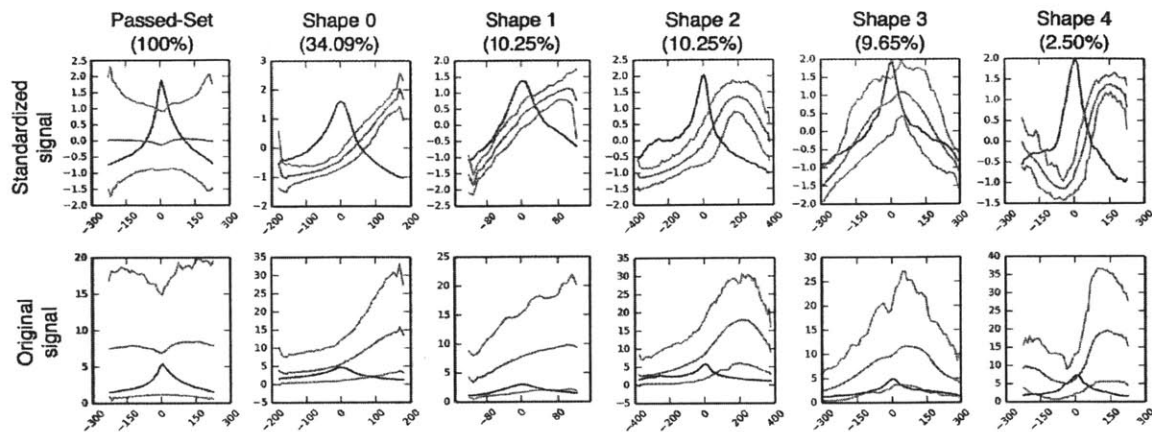


Figure 5-1: Aggregate plots of five shapes, as well as the *Passed-Set*, output by ChromSMS in varied-length mode on a promoter activity cluster. The y -axis displays signal magnitude and the x -axis displays base pairs offset from the anchor point. The red signal in a shape shows the mean signal for the mark H3K4me3, taken across all of the regions in the shape. The gray signals below and above the red signal show the 10th and 90th percentile signal values, respectively, for the regions in the shape. The black signal shows the mean signal for the mark DNase, taken across the regions in the shape. The particular signal that is averaged across the regions of each shape is the average signal across the cell types in which this activity cluster is active. Before the aggregate signals are computed, the signal profiles for each region in a shape are resampled to the median length of the regions in the shape. ChromSMS discovered these shapes by clustering *only* on the standardized signal profiles for the mark H3K4me3. The signals in the plots on the top row are standardized to zero mean and unit variance, and those on the bottom row are not standardized. The percentage under each shape name indicates the percentage of all selected regions in the activity cluster that belong to the shape.

5.3 Review of Output from Multiple-mark Mode

To survey the output of multiple-mark mode, we begin by displaying the output of ChromSMS when given, as input, the same promoter activity cluster used above. Figure 5-3 shows two of the output shapes (along with the *Passed-Set*) that together hold 62.57% of the selected regions. First, note the relatively high enrichment of H3K27me3 in *Passed-Set* and *Shape 0*. This is unusual and unexpected because, as noted in Section 2.2, the mark is typically associated with low expression, yet the displayed signals are averaged across cell types in which we deem this activity cluster to be active. Thus, it is possible that the mark holds important information about the regions in this activity cluster. At the same time, the spatial patterns and

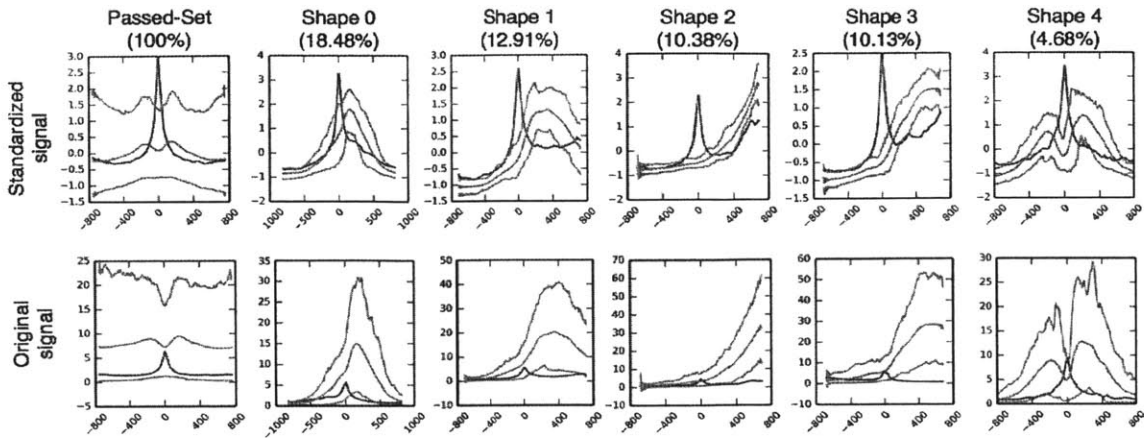


Figure 5-2: See the caption of Figure 5-1 for information about what these plots display; they are of the same form as the ones in that figure. Here, however, we enlarge the window of each region prior to running ChromSMS, and display the aggregate signals over these enlarged windows.

magnitude of the signal of H3K27me3 is evidently different than that of H3K4me3, the sole mark considered by ChromSMS in varied-length mode. This is evidence for a benefit provided by multiple-mark mode: varied-length mode can only capture the information contained by the signals of H3K4me3 or H3K27me3 in separate runs.

Furthermore, we observe that the signal profile patterns of H3K4me3 found by ChromSMS in multiple-mark mode resemble those found in varied-length mode. From a visual inspection, the aggregate H3K4me3 signal in *Shape 0* of Figure 5-3 resembles that in *Shape 0* and *Shape 1* of Figure 5-2. Also, the H3K4me3 signal in *Shape 1* of Figure 5-3 resembles that in *Shape 2* and *Shape 3* of Figure 5-2. We find similar results across many of the other activity clusters.

These observations, as well as others, make a strong justification for the use of multiple-mark mode in general analyses. Although it does not perform the same warping and shifting that varied-length mode explicitly implements, it still appears able to capture spatial patterns for individual marks that are similar to those discovered in varied-length mode. It does this while accounting for the information in magnitude and across more than one mark. Furthermore, the multiple-mark mode inherently makes it simpler to perform comprehensive analyses than varied-length mode because we do not have to consider, for each activity cluster, separate collec-

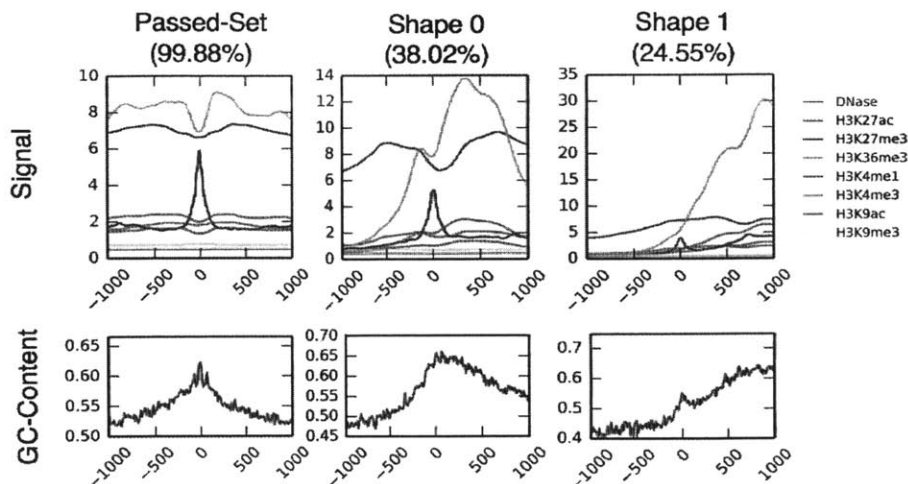


Figure 5-3: Aggregate plots of two shapes, as well as the *Passed-Set*, output by ChromSMS in multiple-mark mode on the same promoter activity cluster used in Section 5.2. We show signals for eight chromatin marks (DNase and seven histone modifications), all of which are considered by ChromSMS when clustering. These plots are not standardized because ChromSMS does not standardize signal profiles in multiple-mark mode. Below we show a signal representing GC-content averaged over the regions in each shape; these are referenced and described later in Section 5.6.

tions of shapes that are the result of running ChromSMS with different input marks. Hence, for the remainder of this thesis, we apply ChromSMS in multiple-mark mode.

Along with promoter activity clusters, we also apply ChromSMS in multiple-mark mode to 41 enhancer activity clusters. Here we display the output from one enhancer activity cluster in Figure 5-4. This enhancer activity cluster is active in 19 cell types; these 19 are a subset of the 32 active cell types for the promoter cluster studied in this Section and Section 5.2. It has 12,797 regions, of which ChromSMS selects 9,645. We see high enrichment of the marks H3K4me1 and H3K27ac, which is expected in active enhancer regions, as well as diversity across shapes in both spatial patterns and magnitude. While asymmetry is not a strong characteristic of these shapes, we do observe large shapes in other activity clusters that have more pronounced asymmetry.

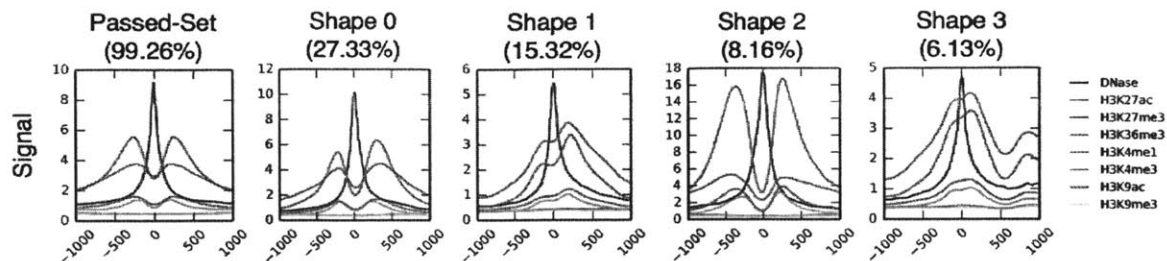


Figure 5-4: Aggregate plots of four shapes, as well as the *Passed-Set*, output by ChromSMS in multiple-mark mode on an enhancer activity cluster.

5.4 Associations Between Directionality and Positions of Nearby Genes

As we have noted, we observe from the output of ChromSMS that the signal profiles in promoter shapes take on strongly asymmetric patterns and that the enhancer shapes also display asymmetry, albeit with less prevalence. In this Section we aim to describe the association between this asymmetry and nearby, annotated genes. We do this for both promoter and enhancer shapes, though the results differ between the two.

5.4.1 Directionality in promoters

We consider promoter shapes that take on, based on the appearance of their aggregate plots, asymmetric patterns. For each region in such a shape, we use a binary search over a set of human genes to find the annotated **transcription start site** (TSS) — or position where RNA polymerase begins transcribing the gene — that is closest to the region¹. True promoters are located at TSSs. However, because the promoters with which we work are “putative,” we find many that are “far” (> 1000 bp) from their nearest TSS. In fact, in cell type-specific promoter activity clusters (as opposed to ones that are constitutively active across most or all cell types), *most* putative promoter regions are far from their nearest TSS.

For each shape, we determine an orientation of the shape by evaluating the **shape’s** most enriched mark (e.g., H3K4me3 in promoters) and finding the side of the **anchor**

¹In particular, we search over 57,820 annotated genes provided by the GENCODE project [21]. (GENCODE v19.) The distances are measured from the midpoint of a region.

point on which its peak lies. In its output, ChromSMS gives the orientation of each region relative to its shape (i.e., whether it is “flipped” while clustering). Then, knowing a shape’s orientation and the orientation of each region relative to its shape, we can determine an orientation of each region. That is, we determine whether the peak of a region’s signal profile lies upstream or downstream of its anchor point. In a region, we call this orientation the region’s “directionality.” This concept is only valid for regions in shapes that have clear asymmetry; in symmetric shapes, along with the *Passed-Set*, it is not well-defined because the shape’s orientation is not well-defined.

Among promoter regions in asymmetric shapes that are “near” (< 1000 bp) a TSS, we find that there is exists a strong correlation between their directionality and the side on which its nearest TSS lies, as expected. Namely, a region’s direction tends to be pointed toward the nearest TSS.

As an example, we look at two asymmetric shapes in the same promoter activity cluster that we have used in past sections of this Chapter. The two shapes, *Shape 0* and *Shape 1*, are shown in Figure 5-3. Because this activity cluster is cell type-specific, most of its putative regions are far from a known, annotated TSS. In Figure 5-5a we show, for increasing distances from the nearest TSS, the fraction of regions in each shape whose direction is oriented toward its nearest TSS.

Among regions in *Shape 0* that are within 1,000 bp of their nearest TSS, 73 of 105 (69.5%) have a direction oriented toward the TSS. (Recall that the window of the region extends 1,000 bp from the anchor point, and we are measuring distance from this anchor point.) By modeling the number of aligned orientations with a binomial distribution, we find that this result is significant at a level of $p < 10^{-4}$.

Note that in *Shape 0* the fraction showed (solid red in Figure 5-5a) stays roughly flat between 1,000 bp and 2,000 bp, indicating that this phenomenon is maintained even in putative promoters whose nearest TSS is located beyond the window we consider. This confirms a result found using CAGT [28].

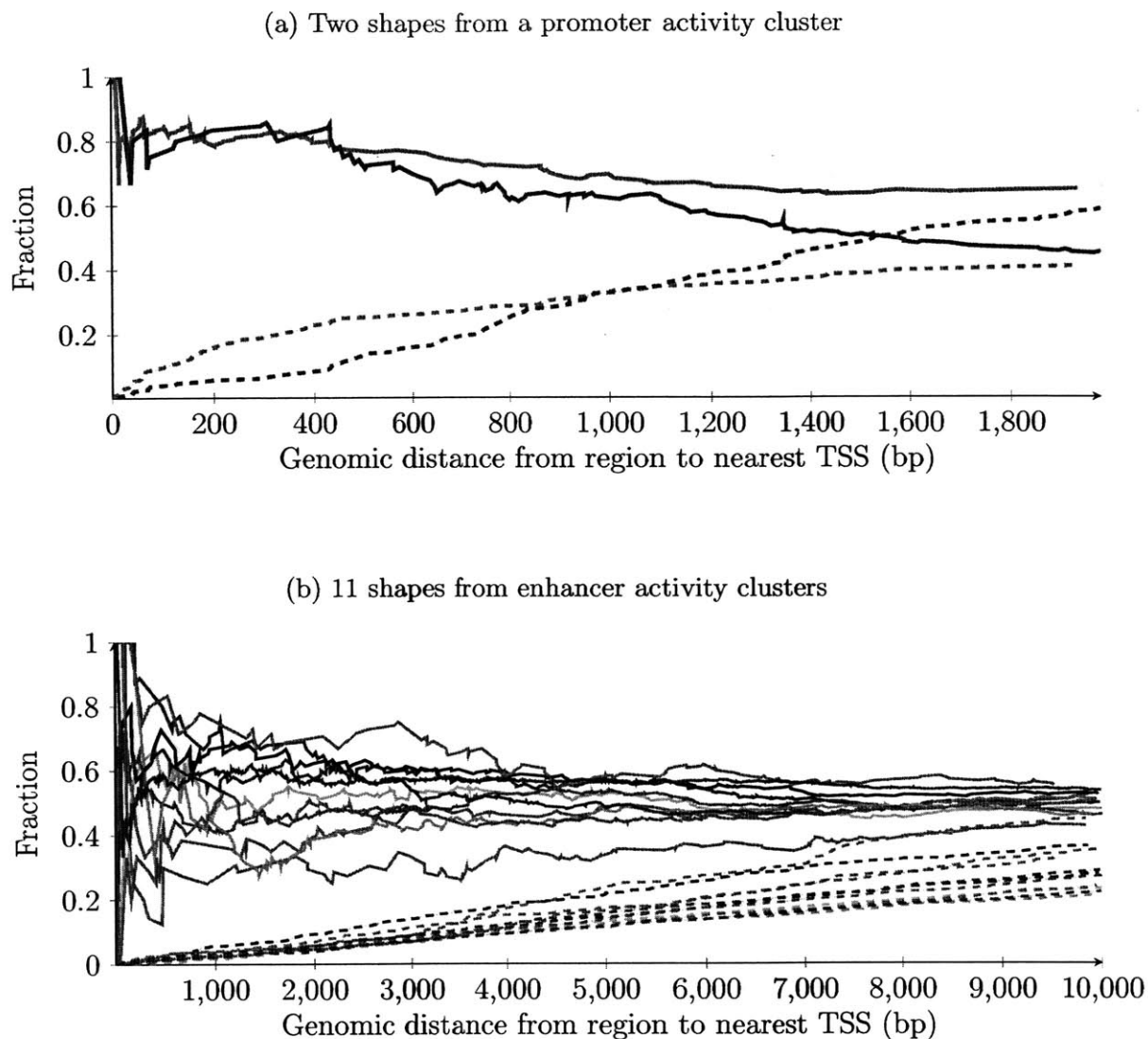


Figure 5-5: We show information about directionality for different shapes, by color. The x -axis shows the genomic distance from the center of a regulatory region to its nearest TSS. At distance x , the solid line shows the fraction of regions that are oriented “toward” their closest TSS, computed among the regions in the shape whose nearest TSS is within a distance x of the region. The dashed line shows the fraction of all regions in the shape whose nearest TSS is within a distance x . **(a)** Two shapes from a promoter activity cluster: *Shape 0* (red) and *Shape 1* (blue). These two shapes are shown in Figure 5-3. We only show up to 2,000 bp, although many regions are further. **(b)** 11 shapes from enhancer activity clusters. They were selected visually from the set of *Shape 0* and *Shape 1* shapes across the 41 enhancer activity clusters based on the asymmetry of the peak of the signal for either H3K4me1 or H3K27ac. The information for each shape is in a unique color. We only show up to 10,000 bp.

5.4.2 Directionality in enhancers

As one might expect, the picture of directionality in enhancer regions is less concrete than it is in promoter regions. To demonstrate this, we consider the set of all *Shape 0* and *Shape 1* shapes across 41 enhancer activity clusters. We manually select ones that display a clear asymmetry in which the peak of the signal for H3K4me1 or H3K27ac lies on one side of the anchor; we choose 11. These look much like *Shape 1* in Figure 5-4, and indeed this shape is one of the 11.

In Figure 5-5b we display, for each of these 11 enhancer shapes, the same information we showed for promoters. We notice that there are some shapes in which the regions within a few thousand base pairs of a TSS tend to have a direction pointed toward the TSS. There are other shapes whose direction appears to point *away* from its nearest TSS. However, even in the most extreme shapes (solid cyan and solid red in Figure 5-5b), the findings are not significant at a level of $p < 0.05$ for the sets of regions within 2,000 or within 4,000 bp of a TSS.

Our findings on enhancer directionality are inconclusive. Although many shapes in enhancer activity clusters appear markedly asymmetric, we are unable to find correlations between this asymmetry and the locations of annotated genes. This is in part due to the distal nature of enhancers, in which few are near a TSS. We also considered other data on annotated genes, including expression, that provide more information than location; however, we are unable to find significant associations between these data and the direction of enhancer regions.

5.5 Functional Enrichment of Discovered Shapes

In this Section we aim to assign functional meaning to the shapes discovered by ChromSMS. Because most of the regulatory regions that make up these shapes are not annotated, we instead look toward the annotations of each region's nearby genes. Thus, we can characterize a shape by analyzing the functions of the genes that are near the regions in the shape. We use the tool GREAT to perform these analyses [32].

In particular, we use GREAT to find statistical enrichments of the regions in each

shape, relative to some background set, in order to assign to each shape annotation terms that are likely not due to chance. There are two options for computing a background set in order to determine enrichments in some shape s :

Option 1. We consider the collection of regions in the activity cluster to which s belongs. Such a background serves as a control for enrichment that may be correlated with particular activity patterns present in s .

Option 2. We consider the entire collection of regions of the same regulatory type as the regions in s . For example, this may consist of the union of the regions in all the enhancer activity clusters with which we work. Unlike Option 1, this type of background clearly does not control for enrichment correlated with activity.

One appeal of Option 2 is that it provides a common background across all tests. However, in this Section we are concerned only with determining functional enrichment within individual shapes, rather than with comparing functional enrichment across shapes, so there is less of a need to have a common background set across all the shapes we test. Therefore, we choose Option 1 due to its ability to account for the bias in activity of each shape.

We work with the 41 enhancer activity clusters to which we applied ChromSMS, and for each we find the shapes that contain at least 1% of the regions in their activity cluster. There are 578 such shapes in total. We apply GREAT to the collection of regions in each of these shapes to obtain a set of statistically significant annotation terms for each shape.

For the majority of shapes (439 of 578) we observe no significant annotation terms in any category of associations. The terms that are associated with shapes come from a variety of ontologies and, among those shapes in which there is enrichment, typically a small number of terms are associated with the shape. Table 5.1 summarizes the degree of enrichment, based on the number of assigned terms, across seven ontologies.

There are a limited number of shapes that show extensive enrichment across multiple ontologies. For example, one (activity cluster 36, *Shape 10*) shows enrichment

Ontology	Number of terms tested	Number of shapes assigned at least one term	Median / mean number of terms assigned to shape
GO molecular function	3401	14	1 / 1.71
GO biological process	8761	22	3 / 13.64
GO cellular component	1125	18	1 / 4.11
Mouse phenotype	7293	21	5 / 20.43
Human phenotype	5982	14	14.5 / 21.71
Disease ontology	2233	17	2 / 12.94
Cancer neighborhood	427	21	2 / 4.29

Table 5.1: We use GREAT to test annotation terms for each separate shape. The terms are grouped into ontologies, seven of which we list here in the first column. The second column gives the number of terms tested by GREAT in each of the seven ontologies. The third column gives the number of shapes assigned at least one significant functional enrichment term. The fourth column gives the median and mean number of terms that GREAT associates with each shape, where we only count shapes assigned at least one term in the ontology. The first three ontologies listed are from the Gene Ontology vocabulary that describes the molecular functions, biological processes, and cellular components of gene products [4]. The mouse phenotype ontology describes associations in mouse that have been mapped to human genes [11]. The human phenotype ontology describes human genotype—phenotype associations [36]. The disease ontology describes associations with human diseases [34]. The cancer neighborhood ontology lists sets of genes that are correlated in expression with cancer-associated genes [38].

for 97 GO biological process associations, 190 mouse phenotype associations, and 111 disease associations. Although this shape contains just 2.17% (128) of regions in its activity cluster, it is distinguished from the other shapes in the activity cluster by its consistently high levels of *both* H3K4me1 and H3K27ac. The shape is shown in Figure 5-6a. Figure B-2 displays the top ten enriched terms for each of these three ontologies in this shape.

Similarly, there are a small number of other shapes with a high degree of enrichment in one ontology, and little to no enrichment in the others. One of these (activity cluster 4, *Shape 2*) shows enrichment for 42 disease ontology terms and one placenta disorder term, but is assigned no terms from any other ontologies. The shape, which contains 5.71% (275) of regions in its activity cluster, stands out due to its strong

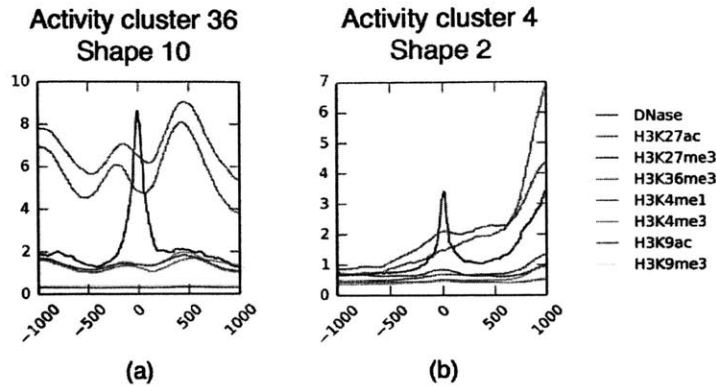


Figure 5-6: These two shapes are examples of ones that display a high degree of functional enrichment. It is difficult to evaluate them alone, and they should be compared to the other shapes in their activity cluster, which are shown in Figures B-1b and B-1a.

asymmetry and apparent proximity to a separate DNase I hypersensitivity site. It is shown in Figure 5-6b.

5.6 Correlations Between Signal Profiles and GC-content

Next we seek to characterize the relation between GC-content and the signal profile patterns we discover. **GC-content** is a measure of the fraction of bases in a window, or across a whole genome, that are either G or C. We work with a signal that encodes GC-content over 5-base windows of the genome². We observe from the plots produced by ChromSMS that there appears to be a correlation between this signal and those of particular histone modifications in many of the shapes, especially in those that take on asymmetric patterns. This is especially pronounced in promoter regions, as shown in Figure 5-3. Indeed, it is known that peaks in nucleosome positioning are often centered at regions with high GC-content [28]. However, in the aggregate plots we observe a less pronounced correlation in enhancer shapes. In this section we quantify

²The track we use — the GC Percent (gc5Base) Track in hg19 — is provided by the UCSC Genome Bioinformatics website.

this correlation using the output of ChromSMS.

We again work with shapes from 41 enhancer activity clusters. For a particular shape, we compute the Pearson correlation coefficient between a histone modification's signal and the signal of GC-content, both of which are taken as averages across the regions in the shape. We evaluate significance by modeling the null hypothesis, that the two have zero correlation, with a Student's t -distribution, and verify with a permutation test. We apply a Bonferroni correction once across all shapes and marks for which we compute p -values.

We first consider the *Passed-Sets* in the 41 activity clusters. For each of seven histone modifications, Figure 5-7a displays the number of shapes, among these 41, that meet a certain threshold of significance for the correlation of its signal pattern with GC-content. We observe that a significant positive correlation is most prevalent for H3K9me3, H3K27me3, and H3K36me3. Indeed, for each of these three modifications, 31 of the 41 clusters show significant positive correlation at the level $p < 0.01$ and 21 show significant positive correlation at the level $p < 10^{-10}$. However, other marks normally associated with active enhancers, including H3K27ac, do not show such abundant correlation in the activity clusters' *Passed-Set*. H3K27ac shows a positive correlation with GC-content in just 21 of these 41 shapes, and negative correlation in the others; the positive correlation is significant in just 3 of the 21 shapes at the level $p < 0.05$.

We also look at correlations between GC-content and the signals for histone modifications inside shapes discovered by ChromSMS. Across the 41 enhancer activity clusters, we consider all shapes that contain at least 5% of the regions in their activity cluster. There are 142 of these in total. Figure 5-7b displays the number of shapes, among these 142, that meet thresholds for significance in correlation with GC-content. We observe, across all histone modifications, more prevalence for significant positive correlations with GC-content than we did with the signals in the *Passed-Sets*. We also observe, across all histone modifications, even stronger significance for correlation. For all seven histone modifications, the signals in 114 of the 142 shapes show significant positive correlation with $p < 0.01$ and 77 show significant

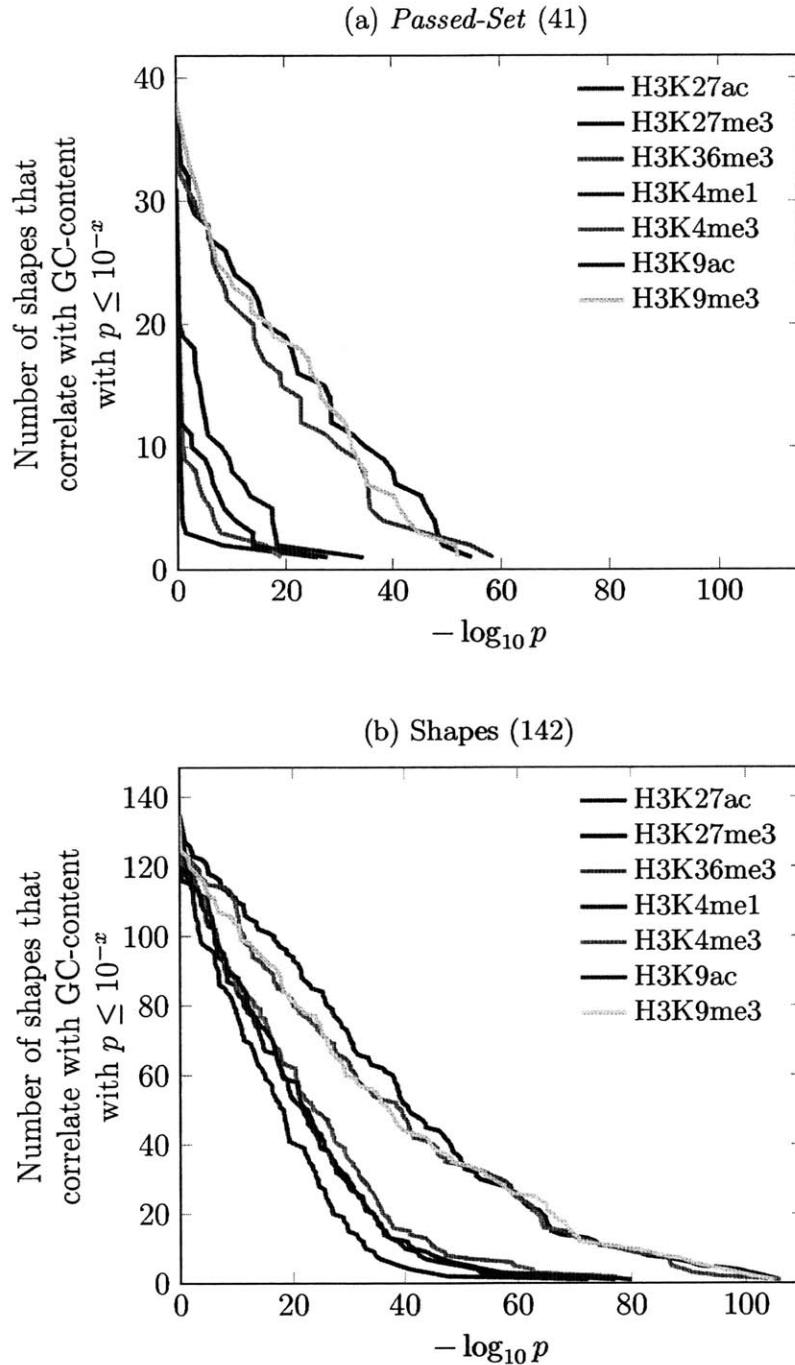


Figure 5-7: These plots show, for some collection of shapes, the number of shapes whose signal for a particular histone modification meets a threshold of positive correlation with GC-content. The x -axis displays p -values on a logarithmic scale. For a particular histone modification, the plotted value at x gives the number of shapes in the collection whose p -value, for the positive correlation with GC-content, is at most 10^{-x} . (a) The collection of shapes is the *Passed-Sets* from 41 enhancer activity clusters. There is just one from each activity cluster, so there are 41. (b) The collection of shapes consists of all shapes from 41 enhancer activity clusters that contain at least 5% of the regions in their activity cluster. There are 142 of these. Note that the value at $x = 0$ ($p \leq 1$) gives the number of shapes that show positive correlation with GC-content. In (b) this is almost all shapes, but in (a) there are modifications for which it is not. For example, in (a) there are only 19 out of 41 shapes for which H3K27ac shows a positive correlation ($p \leq 1$) with GC-content, even if it is non-significant; there are only 5 for which the correlation has significance $p < 1$, and only 3 for which it is significant ($p < 0.05$).

positive correlation with $p < 10^{-10}$. (In the *Passed-Sets*, we could only make a similar claim for three of the seven modifications rather than all.) We still do find that the significance for positive correlation with GC-content is more pronounced for the histone modifications H3K9me3, H3K27me3, and H3K36me3 than it is for the other four we consider.

5.7 Motif Enrichment Within Discovered Shapes

While we find a clear correlation between GC-content and signal profile patterns, it is natural to wonder whether there are additional mechanisms that may drive the diverse set of common signal profiles discovered by ChromSMS. One possible factor are sequence motifs. In the context of this thesis, a **sequence motif** is a frequently occurring sequence of bases in the genome to which particular transcription factors bind. At each site to which a transcription factor binds to a motif, the motif takes on a particular form; the specific sequence at such a site makes up a **motif instance**. Our first task in analyzing the importance of motifs is to discover, at a statistically significant level, motif instances that are associated with the shapes discovered by ChromSMS.

In this Section we describe our method for discovering enriched motif instances within individual shapes, we explain how we analyze the similarity of the motifs associated with each shape, and we summarize our findings. We use these enriched motifs extensively in Section 6.2.

5.7.1 Method for determining enrichment

We work with a total of 1,646 motif instances representing 651 motifs³. For each shape and motif instance, we count the number of regions belonging to the shape whose sequence, on either strand, contains the motif instance within its middle 200

³We are provided with these motif instances by Pouya Kheradpour, a member of Manolis Kellis's lab at MIT's CSAIL. We are also provided with the position weight matrices (PWMs) that describe these motif instances, as well as a tool for measuring the sequence similarity between them, which is used in Section 6.2.1.

bp. That is, motifs are only counted within a relatively small region centered at the peak of DNase I hypersensitivity. Let this count, for motif instance m in shape s , be s_m .

Studies of enrichment are computed relative to background data, and thus we must compute an expected level of motif occurrences. As in Section 5.5, we have two options for a background set used to determine enrichment in a shape s : (**Option 1**) all the regions in the activity cluster to which s belongs, or (**Option 2**) the set of all regulatory regions of the type with which we are working (e.g., all enhancers). There are certainly scenarios in which Option 1 might be preferable. This could include, for example, a look at how a signal representing coverage of enriched motifs correlates with the signals of histone modifications within individual shapes. It could also include studies that compare motif enrichment across the shapes in a particular activity cluster. However, we will be primarily concerned with a single study comparing motif enrichment across all the shapes whose regions are of a particular regulatory type; indeed, this is the focus of Section 6.2. Motif enrichment used in such an analysis should be relative to the same background set. Hence, in presenting these results we work with the background set described in Option 2.

We compute, for each motif instance, the number of regions in the background set whose sequence contains the motif instance. Then, we use a hypergeometric test to determine enrichment. In particular, our null hypothesis is that there is no significant enrichment of motif m in s . Under this null hypothesis, a hypergeometric distribution gives the probability that s has s_m regions containing m , drawn from the background set. We perform multiple testing correction on each shape using the Benjamini-Hochberg procedure to control for FDR at a level of 0.1 [6]. We report the rejected null hypotheses as positive enrichments.

5.7.2 Analyzing intra-shape motif similarity

We also examine how, within each shape, the enriched motif instances relate to each other in their sequence. We use a provided tool called `motif-util` to evaluate the Pearson correlation coefficient between each pair of motif instances based on their

position weight matrices (PWMs). This provides a measure of similarity between pairs of motif instances. Then, for each shape we compute the mean pairwise similarity among its motif instances to measure the motif instance sequence similarity within the shape.

We use bootstrapping to evaluate the significance of each shape’s motif instance similarity. We evaluate the probability that an enrichment is of motif instance m by computing the fraction of all enrichments in the background set that are enrichments of m . For each shape s , we use these probabilities to construct a random set of enrichments equal in size to the number of enrichments in s , and we compute the sequence similarity of this random set. We perform this a large number of times, and let the p -value — representing the probability that the sequence similarity we observe in s is due to chance — be the fraction of sampled sets whose similarity is at least as great as what we observe in s .

5.7.3 Findings

We look for motif enrichment in the 41 enhancer activity clusters that we have been using. We consider all of their shapes that contain at least 1% of regions in the activity cluster, as well as the *Passed-Sets*; there are 619 of these in total. (For convenience, here we will also refer to the *Passed-Set* as a shape.)

We find substantial variation in the number of enriched motifs across activity clusters and across shapes within each activity cluster. However, there is a general trend that shapes with more regions demonstrate a greater quantity of motif enrichment. Indeed, across the 41 activity clusters, the median number of enriched motifs in *Shape 0* (the shape with the most regions) is 56, in *Shape 1* is 15, and in *Shape 2* is 7. Just 283 of the 619 (45.72%) total shapes contain at least one enriched motif. Figure 5-8 shows the number of enriched motifs discovered in three shapes, and the *Passed-Set*, for each of the four largest activity clusters.

We analyze motif sequence similarity within shapes that contain at least two enriched motif instances. 254 (41.03%) of the shapes contain at least two enriched motif instances. We correct for multiple testing across these shapes by controlling

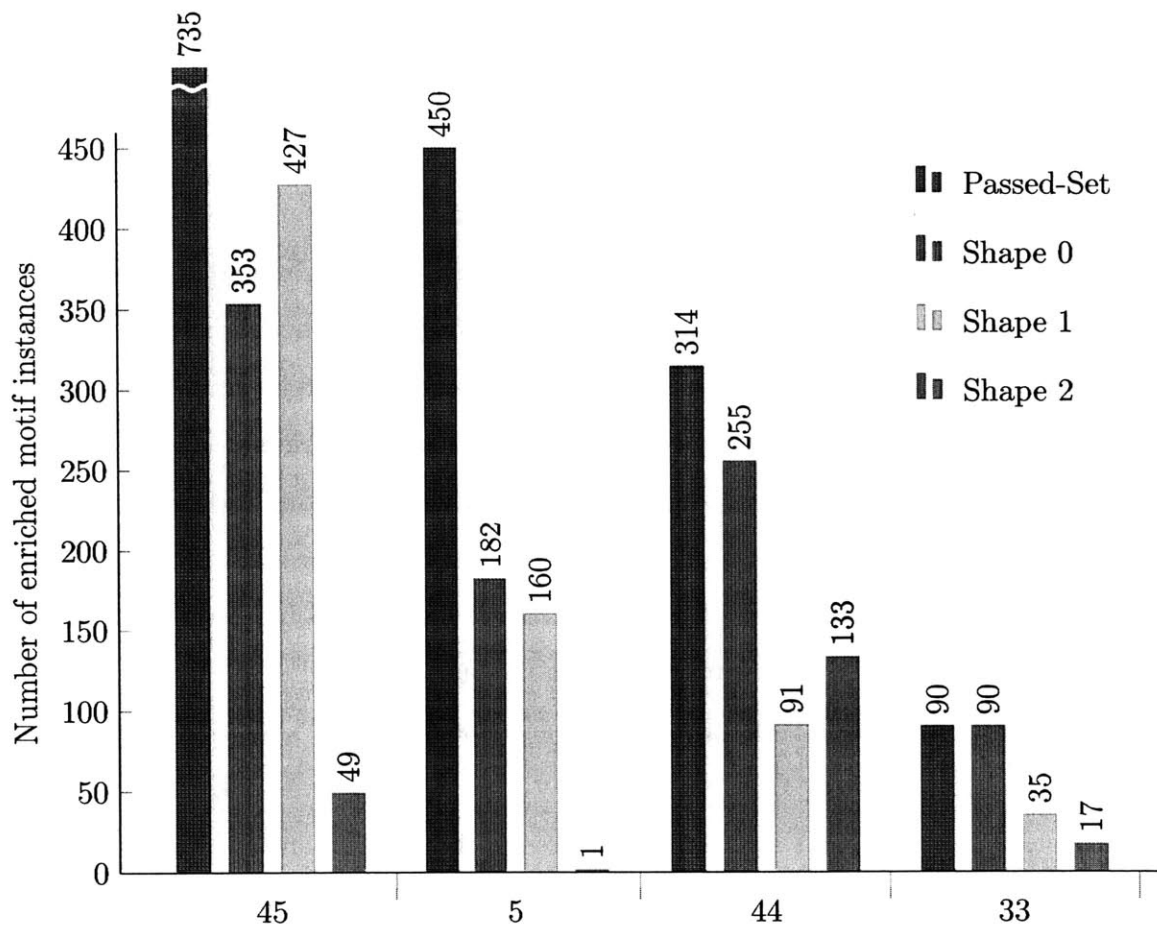


Figure 5-8: This plot displays the four largest enhancer activity clusters: 45, 5, 44, and 33. For each activity cluster, it shows the number of enriched motif instances discovered in the *Passed-Set*, as well as the three largest shapes.

for FDR, and we report adjusted p -values. We find significant similarity within most shapes. Among the 41 *Passed-Sets*, 31 show significant sequence similarity in their enriched motif instances at a level of $p < 0.05$. Of the 213 shapes discovered by ChromSMS that have at least two enriched motif instances, the instances in 190 of them display significant intra-shape sequence similarity at a level of $p < 0.05$.

5.8 Discussion

We had three primary goals in performing the analyses of this Chapter. In this Section we state these goals, and discuss what we find as it pertains to each of them.

First and foremost, we aimed to evaluate the effectiveness of ChromSMS in discovering meaningful biological patterns. We did this in part by visualizing aggregate plots of the shapes it produces and we find that, even with the complexity of warping signals and of incorporating information from multiple marks, ChromSMS is still able to produce shapes that appear distinct from one another. And we extended this largely by conducting analyses of the relations between the shapes we uncover and additional biological information. In promoters we notice significant correlations between the directionality of their signal profiles and the position of TSSs. We also observe strong correlations between GC-content and the signal profiles of histone modifications. In enhancers, the ubiquity of these correlations with four of the considered histone modifications — including those that normally define enhancer regions — across activity clusters is only apparent when we consider individual shapes. Furthermore, we find extensive functional enrichment in a limited set of shapes. The shapes used in these analyses were discovered by ChromSMS purely from input signal tracks and putative regulatory regions; that is, they were discovered without information indicating TSS positions, GC-content, or the functions of genes near the input regions. That they were discovered without this additional information, yet show associations with it, is a compelling indicator that the patterns discovered by ChromSMS do indeed possess intrinsic biological meaning.

The second and third objectives of our analyses are more biologically motivated. They are to understand the role that these shapes have in the chromatin environment, and to learn about the mechanisms that drive the shapes we see. Our results on functional enrichment suggest that there may be a small number of shapes that carry functional significance and thereby inform gene expression. In particular, based on the two examples we explored, these could be shapes that are distinctive from others that display similar activity patterns, either because of their unusually high enrichment of some modification or the unusual spatial patterns they exhibit. Enrichment in disease ontology terms is especially evident. While our results do not point to a causative effect, one might speculate that the abnormal magnitude or spatial composition of these shapes influences the structure of chromatin in their regions in a way that more

easily alters the expression patterns of their nearby genes.

To study the mechanisms behind the patterns we discover, we looked toward GC-content and motif enrichment. The correlations we observe with GC-content strongly suggest that this plays a large role in governing the deposited patterns of histone modifications both across promoters and enhancers, although this is not surprising due to significant role of GC-content in controlling nucleosome positioning, particularly in promoters [17]. We also identify motifs enriched at the binding sites of the regions within individual shapes; however, in order to learn about their significance and the roles they may play, we must compare this enrichment across shapes. This is a major focus of Chapter 6.

Chapter 6

Analyses Across Discovered Shapes

In this Chapter, we continue presenting and discussing results from the application of ChromSMS. Whereas in Chapter 5 we considered analyses performed on regions within individual shapes, here we describe and present results from studies that explicitly consider multiple shapes at once.

We begin in Section 6.1 by describing a method for conducting a statistical comparison of shapes, and we display what we find by applying it to the shapes discovered by ChromSMS. Then, in Section 6.2 we extend the look at motif enrichment from Chapter 5. In particular, we present a comprehensive study that compares the enrichment of motifs across many shapes, with the goal of better understanding the role of motifs in driving the signal profile patterns we observe. In Section 6.3 we discuss these results.

6.1 A Statistical Test for the Comparison of Shapes

An important quantity for evaluating the effectiveness of ChromSMS and for conducting analyses with its output is a measure of the similarity of two shapes. These could be two shapes within the same activity cluster or across different activity clusters. We have defined dissimilarity measures between individual signal profiles in Section 4.3.1 and these could be used to compare shapes as described in Section 4.3.2. However, we seek a measure that has a more intuitive meaning, that more directly accounts

for the average signal profile displayed by the shape, and that is independent of the methods used in clustering. Such a measure could then be used in a statistical test to evaluate the significance of the difference between two shapes.

6.1.1 Method: measuring similarity between two shapes

Suppose we wish to compare shape x and shape y . Consider the “signal” in a region for x and for y to be a concatenation of the signal across all used chromatin marks. For each position i in this signal, we wish to establish a confidence interval of the difference in the mean signal of x at i and the mean signal of y at i , where the mean is taken across the regions in the shape. There are two methods by which we can compute these confidence intervals:

Boostrapping. Let $|x|$ be the number of regions in x and likewise for $|y|$. We do the following a large number of times: We randomly sample $|x|$ regions from x with replacement and randomly sample $|y|$ regions from y with replacement. Let \tilde{x} be the mean signal across the randomly sampled regions from x and likewise for \tilde{y} . We flip \tilde{y} if doing so yields a smaller L_1 -distance to \tilde{x} . Then, we compute the difference vector $\tilde{x} - \tilde{y}$ and store this. By considering the difference vectors computed across all the repetitions, we can compute a 95% confidence interval on the difference in means at each position of the signal.

Using a t -score. Let \tilde{x} be the mean signal across all the regions in x and likewise for \tilde{y} . Similar to the bootstrapping approach, we flip all the regions in y if doing so yields a smaller L_1 -distance to \tilde{x} . Given the distribution of signals across the regions in x and in y , we compute the standard error of the difference in the mean signals at each position of the signal. We then use, as the critical value, the t -score with cumulative probability 0.975, and multiply this by the standard error to establish a margin of error at each position. Centered at $\tilde{x} - \tilde{y}$, this yields a 95% confidence interval on the difference in means at each position.

We check whether the confidence interval at each position overlaps zero. If it does, we determine that the mean signals of x and y are not different at that position at

a statistically significant level. We compute the fraction of confidence intervals that overlap zero, and this fraction serves as our similarity measure.

Note that this measure may represent a true difference despite being high. For example, with eight chromatin marks, consider two shapes with identical mean signals for seven of those marks, but significantly different mean signals for the other. It would be reasonable to state that these two shapes are different; however, the similarity measure we outline would yield of a value of 0.875.

While bootstrapping is free of assumptions, it is comparatively computationally expensive to compute. Hence, we use the latter method when we must compute the similarity measure many times, such as in the permutation tests described in Section 6.1.2. However, we use bootstrapping when it is computationally reasonable to do so; this is the case when comparing two collections of shapes, as described in Section 6.2.2.

6.1.2 Method: evaluating significance of difference between two shapes

It is difficult to evaluate the significance of the difference between two shapes directly using the absolute similarity measure described above. Instead, we use a permutation test on top of the similarity measure.

In particular, consider two shapes x and y whose similarity measure, as computed above, is $s(x, y)$. Let $|x|$ be the number of regions in x and likewise for $|y|$. We do the following some large number of times: We randomly permute the union of the regions in x and in y to form a shape x' of size $|x|$ and a shape y' of size $|y|$. We then compute the similarity of x' and y' , using a t -score, and store this.

The p -value for the difference between x and y is then the fraction of computed similarities from the permuted shapes, across all the repetitions, that are less than $s(x, y)$.

6.1.3 Findings

Within each of the 41 enhancer activity clusters, we evaluate the significance of the difference between every two shapes in the cluster and produce a p -value for each pair of shapes. For all pairs within an activity cluster, we use the Benjamini-Hochberg procedure to control for FDR and report adjusted p -values [6].

To restrict the number of shapes we consider, we only consider those that contain at least 2% of the regions in their activity cluster. In total, summing over the number of pairs of such shapes in each of the 41 activity clusters, we find 1,084 pairs. We find that 1,068 (98.52%) of these pairs contain shapes that are different at the level $p < 0.01$ and that 1,081 (99.72%) contain different shapes at the level $p < 0.05$. According to our statistical test, it is clear that, within activity clusters, ChromSMS produces shapes that are different at a statistically significant level.

6.2 Comparing Motif Enrichment Across Shapes

In this Section we describe our comparison of shapes on the basis of the enriched motifs discovered in Section 5.7. In particular, we are interested in the question of whether there are particular collections of motifs such that each collection is enriched in some collection of shapes. This could address whether the motif instances located at TF binding sites are correlated with the signal profile patterns we observe. We work with the same 41 enhancer activity clusters that we have used in previous sections.

We begin by clustering the 1,646 motif instances to reveal groups of motifs that are enriched in common collections of shapes. We represent each motif instance as a set consisting of the shapes, across all 41 activity clusters, in which the motif instance is enriched. We then use the Jaccard distance to measure dissimilarity between each pair of motif instances, and cluster these using average-linkage agglomerative clustering. Figure 6-1 shows the enrichment of motif instances, grouped by cluster, across 242 shapes. After removing clusters that have fewer than 4 motif instances and motif instances that are not enriched in any shapes, we are left with 19 motif clusters that together encompass 1,522 of the motif instances.

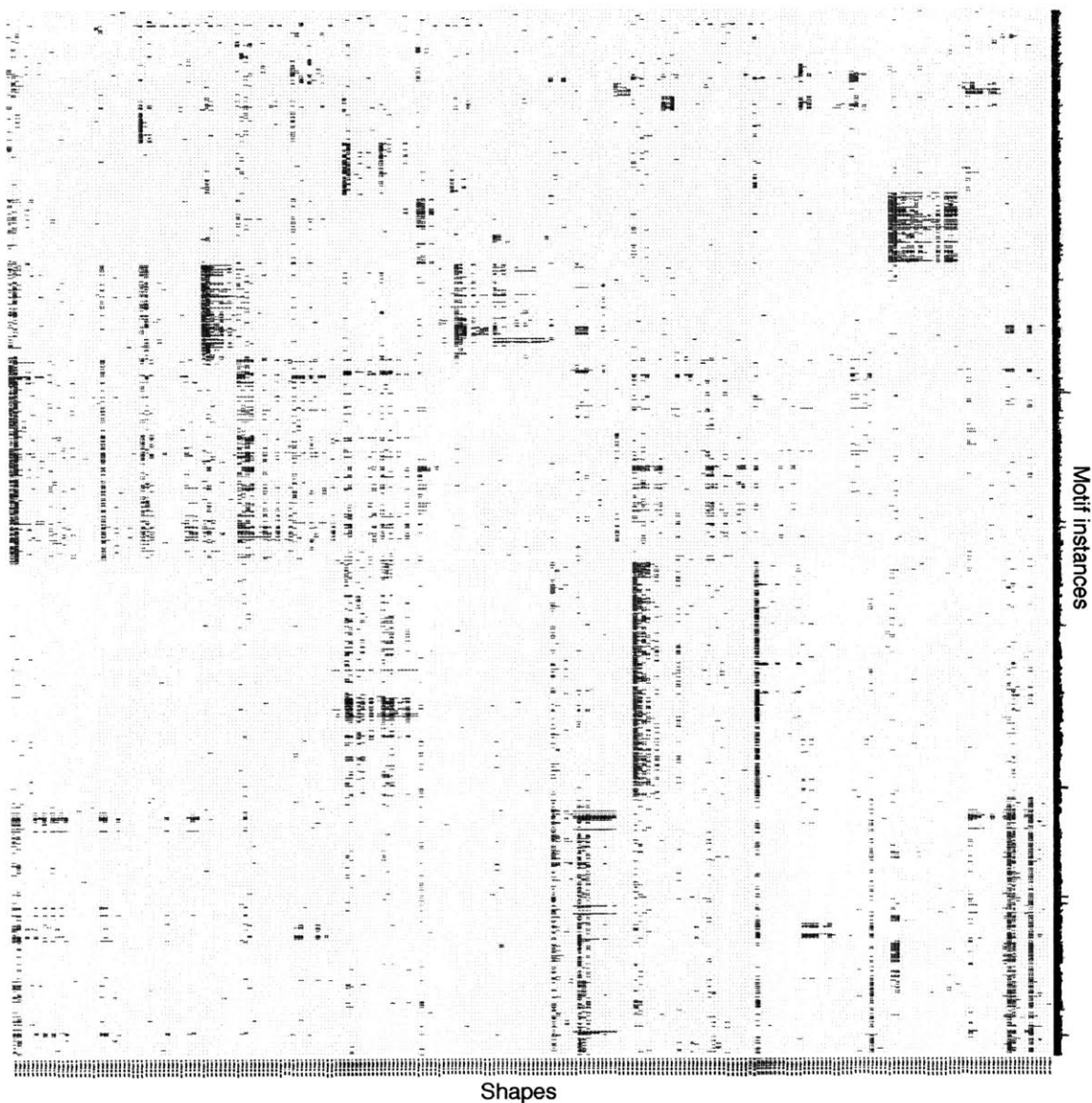


Figure 6-1: This binary matrix displays the enrichment of **motif** instances from shapes spanning 41 enhancer activity clusters. There are 1,544 **motif** instances, one on each row. (Only motif instances enriched in at least one shape **are** shown.) There are 242 shapes, one on each column. (Only shapes in which **there** is at least one enriched motif are shown.) The motif instances (rows) are grouped by their clusters.

We find that there are vivid and distinct collections of shapes in which the motif instances of each motif cluster are enriched. In many **cases** multiple shapes in a collection are from the same activity cluster, to a **higher** degree than one might expect; however, this is **likely** a consequence of the **chosen** background set and of

motif enrichment being correlated with differences in activity.

To better understand the relation between motif enrichment and signal profile patterns, we analyze these 19 clusters of motifs based on the sequences of their motif instances (Section 6.2.1), as well as on shapes that we can associate with each of these clusters (Section 6.2.2).

6.2.1 Comparison of motif clusters: sequence

After finding clusters of motif instances that are distinct in the shapes in which they show enrichment, we are interested in learning about the sequence of these motif instances. In particular, we are interested in comparing their sequences both within clusters and across clusters.

As in Section 5.7.2, we evaluate the Pearson correlation coefficient between each pair of motif instances. We take 1.0 minus this value to obtain a Pearson's distance for each pair. We use the average pairwise dissimilarity between clusters as a measure of cluster separation on the basis of their motifs' sequences. Similarly, we use an average pairwise dissimilarity among the motif instances within a cluster as a measure of cluster compactness.

Figure 6-2a displays the sequence-based similarities for each pair of motif clusters, and Figure 6-2c displays the significance of their difference. We find that the sequences in the vast majority of pairs of these motif clusters are different at a statistically significant level. Indeed, after adjustment for multiple testing, we find that, of the $19^2 - 19 = 342$ pairs of non-identical motif clusters, 292 display a separation in sequence that is not due to chance at a level of $p < 0.05$.

Furthermore, the values along the diagonal in Figure 6-2a demonstrate the compactness of each cluster's sequences. Some contain a large number of motif instances that are especially compact and well-represented by a single instance. For example, one (cluster 7) contains 21 instances, has an average intra-cluster pairwise similarity of 0.741, and is associated with 7 dense shapes; it is well-represented by a motif that binds the factor TCF3.

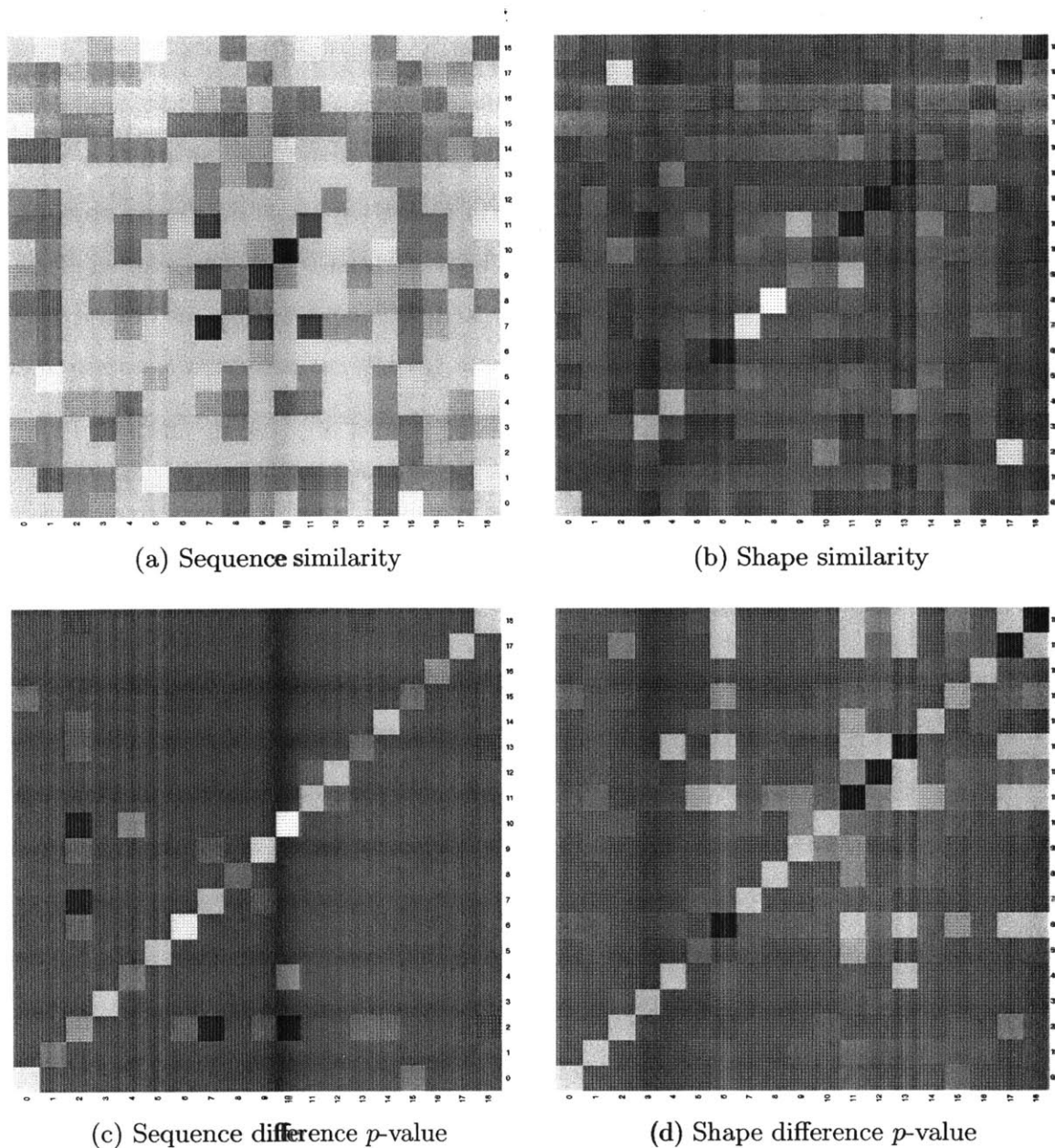


Figure 6-2: These symmetric heatmaps display the results of comparing the 19 discovered motif clusters. The diagonals serves as an indicator of each cluster's compactness. **(a) and (b)** These show a similarity measure between each pair of clusters based on the sequence of their motifs (a) and their dense shapes (b). In (a), to better display the range of values given by measuring sequence similarity, we scale the heatmap from 0.0 (blue) to 0.75 (red). In (b), we scale the heatmap from 0.0 (blue) to 0.5 (red). **(c) and (d)** These show a p -value for each pair of clusters representing the probability that the division of the pair is not an arbitrary permutation (i.e., the significance of the two being different) based on the sequence of their motifs (a) and their dense shapes (b). The p -value for each pair of motif clusters (both (c) and (d)) is computed via a permutation test. Namely, we randomly permute the items in the union of two clusters to produce two clusters that correspond in size to the true ones. We then evaluate the similarity (sequence or shape) between these permuted clusters. After doing this a large number of times, we let the p -value be the fraction of similarities that are at least as small as the measured similarity of the pair. We perform multiple testing correction, across all pairs in the matrix, to control for FDR and display adjusted p -values. Note that, due to the method by which we measure this p -value, the values on the diagonal should be roughly 0.5. There are indeed exceptions in (d) because some motif clusters have just one dense shape, and these p -values would be computed as 1.0. These heatmaps are scaled from 0.0 (blue) to 1.0 (red).

6.2.2 Comparison of motif clusters: dense shapes

For each cluster c of motifs and each shape s , we calculate a **density of s in c** , which is the fraction of motifs in c that are enriched in s . We say that the **dense shapes in c** are those shapes s where the density of s in c is at least 0.333. One can think of the dense shapes of a motif cluster as being the shapes in which, to some degree, the motifs in the cluster are enriched.

We find at least one dense shape for each of the 19 clusters of motifs. (The threshold was chosen empirically to ensure this.) The cluster of motifs with the greatest number of dense shapes has 10 dense shapes. The median number of dense shapes in a motif cluster is 3.

Following the approach of comparing the sequences of motifs within and across motif clusters, we also compare the dense shapes within and across motif clusters. To evaluate the similarity between two shapes, we use the measure described in Section 6.1. We take 1.0 minus this value to obtain a measure of dissimilarity between two shapes. Then, we use the average pairwise dissimilarity between clusters to evaluate cluster separation based on their dense shapes. We also use an average pairwise dissimilarity among the dense shapes within a cluster to evaluate cluster compactness. (As with comparing sequence across motif clusters, there are other measures for computing the separation between two collections of shapes. One that might be more suitable, especially if the dense shapes in a motif cluster are very different from each other, is a median Hausdorff distance. However, we wish to account for intra-cluster variance and to evaluate cluster compactness, which an average pairwise dissimilarity offers.)

In Figure 6-2b we display, for each pair of motif clusters, the similarity between these collections of shapes as computed in this manner; in Figure 6-2d we display the significance of their differences. We are unable to determine, at a statistically significant level, differences across the collections of shapes that are associated with each cluster of motifs. Based on our test, 236 of the 342 pairs of collections of shapes are different at a level of $p < 0.2$, after adjustment. This difficulty in establishing

differences is likely due in part to the methodology of our test, which considers shapes as a whole rather than pooling their regions. The small number of dense shapes per cluster (the most has 10) makes it difficult for this procedure to conclusively establish differences.

We find that, although there exist clusters of motifs that are well-separated in their sequences and the collections of shapes in which they are enriched, the shapes within these collections can be quite diverse. For example, one (cluster 9) has 10 dense shapes, but the average pairwise similarity among them is just 0.190; this is despite the fact that the cluster is compact in its sequence. And although we cannot conclusively establish differences, we can observe trends from an aggregate plot of the dense shapes. For example, some clusters (2, 4, 7, and 17) appear to be relatively highly enriched for H3K27ac compared to the dense shapes in other motif clusters. One of these (7) is highly compact in its sequence, which may suggest a group of similar motifs that are associated with high enrichment of H3K27ac.

Table A.1 displays more detail about each of the 19 motif clusters we find.

6.3 Discussion

In this Chapter we continued our analyses from Chapter 5 by comparing observations across the shapes discovered by ChromSMS. This is vital in establishing a more complete understanding of the role of and the mechanisms behind the signal profile patterns we observe.

We began with the important finding that the shapes uncovered by ChromSMS are, in the vast majority of pairs, distinct from one another at a significant level. Of course, this observation is based on the *ad hoc* statistical test we developed; other tests may not find such significance. Nevertheless, the framework established in developing this test proves useful in comparing shapes, which we need to perform in the next component of the Chapter.

We continued our analysis of motif enrichment by seeking connections between the enriched motifs at binding sites of regulatory regions and the shapes to which

these regions belong. An underlying goal in finding these connections is to determine whether sequence motifs or the transcription factors that bind to these motifs influence the patterns of histone modifications around these binding sites. We were able to find disjoint sets of motif instances such that the instances within these sets are enriched in distinct collections of shapes. We found these sets by only considering their patterns of enrichment across shapes, yet we observed that many of these sets are both compact and distinct from one another based on their sequence. This finding supports a possible association between sequence motifs and the signal profiles that encompass them. It is hindered by a comparison of the differences among these collections of shapes, which we do not find to be significant; however, this could simply be a consequence of the small number of shapes that we associate with each set of motifs.

These findings, though inconclusive, serve as motivation for further work in understanding the mechanisms behind the patterns we observe, and the particular role that motifs may have in this process.

Chapter 7

Further Applications of Ideas Behind ChromSMS

At its core, ChromSMS is about capturing information in signal profiles to find patterns in regulatory regions. But this basic idea can be applied to tools that perform functions other than clustering. We summarize a methodology and outline results for two such applications. The first is a Hidden Markov Model (HMM) that automatically learns hidden states behind regions of the epigenome based entirely on the *shapes* of chromatin mark patterns, without any information about their magnitude. The second is a use of locality-sensitive hashing, a recent development in computer science, to very efficiently discover regions in the genome that are “similar” to an input query region, where similarity is measured using signal profiles.

7.1 Detecting Regulatory Regions by Searching for Shapes¹

As explained in Section 2.3, tools like ChromHMM automatically learn hidden states of an epigenome by relying on a binarization of the level of enrichment of chromatin marks. However, such a binarization bypasses the complexity of the patterns in which

¹This section summarizes work submitted in a report for a final project in the MIT class *Computational Biology: Genomes, Networks, Evolution* (6.878).

chromatin marks are deposited, and it is possible that this complexity holds biological significance. Although we do not claim that this complexity is indeed significant, we can still develop a “proof-of-concept” to test the idea of uncovering hidden states based entirely on spatial patterns of chromatin marks, without information about large-scale relative degrees of enrichment. In particular, we develop a multivariate HMM, which is similar in many ways to ChromHMM, to accomplish this. The method constructs reference shapes using ChromSMS, searches an epigenome for these shapes, and then bins the epigenome so that each bin emits zero, one, or multiple shapes for each mark. Then, the HMM can learn the hidden state behind each bin.

7.1.1 Dataset we use and constructing reference shapes

We averaged 127 epigenomes for each of five well-studied histone modifications (H3K4me3, H3K27me3, H3K9me3, H3K36me3, and H3K4me1), resulting in a single epigenome (called AVG) with signals for each; this averaging should help reduce noise. (These signal tracks are computed by a fold-change enrichment, and are log-transformed to reduce the effect of outlier values.)

We begin in a supervised fashion, taking shapes output by runs of ChromSMS on a set of 103 promoter regions deemed active across all the epigenomes used in constructing AVG. Specifically, we run ChromSMS in varied-length mode, on each of the five modifications, on the set of regions that does not contain chromosome 1; this provides a reference set of shapes for various histone modifications. We then use only the regions from chromosome 1 for determining emissions and training an HMM, which ensures the parameters do not overfit the reference shapes. All clustered signal profiles are standardized to zero mean and unit variance so that magnitude is not accounted for.

For each reference shape we establish a threshold, based on the radius of the regions in its shape, to determine if a signal in a given window resembles the shape. If the distance between a signal and the reference shape’s signal falls below the threshold, we deem them similar.

7.1.2 Generating emissions for the HMM

We generate emissions over the first (roughly) 25 million bp of chromosome 1, by searching the AVG epigenome. We take the top 10 reference shapes for each mark (ranked on size) and, for each, slide through the epigenome at a rate of 10 bp with variable window lengths, where the lengths come from the distribution of lengths in the regions of the reference shape. We standardize the signal in the window and determine whether it is similar to the reference shape. Because one shape may lead into another, a portion of a region should be able to emit more than one shape. Hence, we divide the epigenome into bins of 200 bp each and each bin emits shapes it is similar to, for each mark. Namely, using 10 reference shapes for each of five marks, each emission is a binary vector $[e_{ij}]$ for $0 \leq i < 5$ and $0 \leq j < 10$ in which each bit indicates whether mark i takes on shape j in the bin.

7.1.3 Determining hidden states from the HMM

We use Baum-Welch training to learn the parameters of the HMM and posterior decoding to select the most likely hidden state behind each bin. We compute BIC scores for models with different numbers of states k , and the scores plateau at $k \geq 10$, so we select 10 states.

We detect regions based on contiguous annotations of the same hidden state. (Recall, based on our construction of reference shapes, that this is only intended to discover promoter regions.) Table 7.1 shows the number of regions annotated with each hidden state and the mean length of regions assigned the state.

Some hidden states annotate regions that are not of interest; for example, state 8, with the longest average length, may denote background, and shape 4, with the shortest average length, may denote shapes with little complexity. State 3 stands out. It is relatively small and, of the 103 promoter regions used in creating the reference shapes, 54.36% of them are annotated with state 3. 28.77% of the regions annotated with state 3 belong to these 103 promoter regions. This suggests that the state has some underlying significance relating to a set of “true” promoter regions.

State number	Number of regions	Average region length (bp)
0	495	1347.47
1	625	1363.84
2	1161	1236.69
3	358	2105.59
4	1113	694.88
5	3373	1037.36
6	312	1067.95
7	896	1381.70
8	4107	3072.61
9	2439	1125.13

Table 7.1: In a 10 state model, the number of discovered regions annotated by each state and their mean lengths.

7.1.4 Discussion

While the results in these promoter regions are strong, they could simply be a consequence of the positioning of nucleosomes in promoter regions, and may demonstrate little about the underlying information of the spatial patterns of histone modifications. Thus, it would be interesting to apply these techniques to enhancer regions. Furthermore, an ambitious extension would be to remove the need for supervised learning so that the tool could learn reference shapes *de novo* from the signals, perhaps via an iterative algorithm, without requiring annotated regions. Although there is certainly more to investigate, this work has promising implications in helping us discover regulatory regions of interest, and in helping us better understand relations between these regions.

7.2 Querying for Region Near Neighbors via Locality-Sensitive Hashing²

One potential problem involving regions of an epigenome is the following: Given a region q , quickly report a collection of regions in the epigenome that are “similar” to

²This section summarizes work submitted in a report for a final project in the MIT class *Geometric Computing* (6.850).

q , where two regions are defined as similar if the distance between them is within some parameter R . While a solution to this problem may be directly useful in certain analyses, it could also lay the framework for the problem of efficiently clustering a very large number of regions.

Even with a simple metric, this is a nontrivial problem to perform efficiently in practice because the number of regions n under consideration is often large and the number of signal values in any particular region, or the dimension d of their signals, is also large. A naïve solution would take $O(n)$ time, for fixed d , to examine each of the regions and report the similar ones. We seek a solution whose queries, after preprocessing, run in $o(n)$ time in theory and also run more efficiently in practice. This is an application of the *exact randomized R -near neighbor problem* [3]. The problem is the following: Given a set P of points in d -dimensional space, as well as parameters $R > 0$ and $\delta > 0$, perform preprocessing on P to develop a data structure that, given a query point q , returns each point $p \in P$ that is within a distance R of q with probability $1 - \delta$; we call p an R -near neighbor of q . There are solutions to this problem that use locality-sensitive hashing (LSH).

In fact, there is an existing tool, E²LSH, that solves the exact randomized R -near neighbor problem [2]. But it leaves out certain features that might be useful in this context, such as an implementation of the Hamming distance family and optimization over a particular parameter (w , defined later). Hence, we implement the approaches from scratch and evaluate their performance and the quality of their results. We begin by describing the approaches.

7.2.1 Comparing regions with Hamming distance

Because the signal profiles in any one region can contain substantial noise, it might be helpful to experiment with a binarization of the signal. Given a signal p , we construct a bitstring p' where $p'(i) = 1$ if and only if the value $p(i)$ exceeds an absolute threshold. Following the approach of [3], we use LSH functions $h_i(p) = p'(i)$ for $i \in \{1, \dots, d\}$. We must determine the parameters k and L . We vary k over a range and, for each k , set $L = \log_{1-P_1} \delta$, where $P_1 = 1 - R/d$. Each R -near neighbor for a query point is

guaranteed to be returned with probability $\geq 1 - \delta$.

The optimal value k depends on the query points [2]. We build the data structure for each value of k on the subset of points and perform a query on each of the points to measure the average query time. We then choose the pair (k, L) that minimizes the average query time. This is the approach of E²LSH.

7.2.2 Comparing regions with Euclidean distance

Another distance metric that incorporates all the information of a signal — and indeed one used by CAGT — is the L_2 distance. Again following the approach of [3], we pick a random projection vector $a \in R^d$ and a random value $b \in [0, w)$ and use the hash functions $h_{a,b}(u) = \lfloor \frac{a \cdot u + b}{w} \rfloor$ for a parameter w . Using the expression for the probability $p(u)$ of collision given in Section 4.1.2.1 of [2], and letting $u = R$, we find that

$$P_1 = \operatorname{erf}\left(\frac{w}{\sqrt{2}R}\right) + \frac{\sqrt{2}R(-1 + e^{-w^2/(2R^2)})}{\sqrt{\pi}w} \quad (7.1)$$

Like before, we can compute a sufficient value of L , namely $L = \log_{1-P_1^k} \delta$, given fixed values of k and w . Because k and w are degrees of freedom, we vary both over a range and choose the tuple (k, w, L) that minimizes average query time.

7.2.3 Results: query time

We evaluate the approaches on a set of $n = 2637$ enhancer regions. Each consists of 2000 bp, and we average across 10 bp, so that $d = 200$. We choose a subset (20%) of the input points for determining parameters, and use the remainder for testing. We only consider the signal corresponding to the histone modification H3K27me3. (Although this modification is not normally associated with active enhancer regions, it demonstrated substantial diversity in shapes found by ChromSMS across these regions.) We select a failure probability of $\delta = 0.1$.

Hamming distance. We select $R = 50$ as the Hamming distance at which to find neighbors; this is reasonable because it implies at most one-quarter of the bits can be

different for two regions to be considered neighbors. We determine average query time for $k \in \{1, \dots, 14\}$ and find that increasing k strictly worsens query time. This makes sense because, due to the way in which the bitstrings were constructed (a universal threshold), many of them are very similar; in fact, many are made up entirely of 0's or entirely of 1's. Thus, increasing k increases the number of times each query must be hashed without reducing the number of collisions. The minimum query time is at $k = 1$, and yields an average query time of 0.0543 seconds. This is in fact slower than finding neighbors naïvely via a linear scan, which takes an average of 0.0486 seconds per query.

Euclidean distance. We select $R = 30$ as the Euclidean distance at which to find neighbors and experiment with $k \in \{1, \dots, 11\}$ and $w \in \{100, 200, \dots, 500\}$. A minimum average query time of 0.0059 seconds is achieved at $k = 4$ and $w = 100$. Finding neighbors naïvely via a linear scan, while using the Euclidean distance, takes an average of 0.0289 seconds per query. (It is interesting that both are faster than the Hamming distance queries. This is likely due to constant factors; we compute Euclidean distance using NumPy and a Hamming distance with a loop in Python.)

7.2.4 Results: accuracy

To evaluate the quality of the output, we begin by using ChromSMS to cluster the input regions. We call a query region q a true neighbor of a reported neighbor p if q belongs to the same shape as p . Thus, the shapes serve as a “ground truth.”

We consider a subset of query points from each of these three largest shapes. We consider whether the reported neighbor is a true neighbor. Table 7.2 reports these results using Hamming distance and Table 7.3 reports them using Euclidean distance.

7.2.5 Discussion

The results using Hamming distance are poor. It is likely that the metric is only capturing coarse differences in magnitude. For example, about 50% of the reported

Query points from shape	Avg. number reported per query	Avg. fraction reported that are true neighbors	Shape size
0	1325	49.36%	31.66%
1	1178	28.21%	21.12%
2	613	18.53%	8.04%

Table 7.2: Results for query points (regions) using Hamming distance.

Query points from shape	Avg. number reported per query	Avg. fraction reported that are true neighbors	Shape size
0	239	66.93%	31.66%
1	135	39.74%	21.12%
2	1.43	46.00%	8.04%

Table 7.3: Results for query points (regions) using Euclidean distance.

neighbors for regions in *Shape 0* are also in *Shape 0*. Regions in *Shape 0* have low magnitude and their bitstrings are mostly 0-vectors. Thus, the reported neighbors are simply other regions with low magnitude, and this result suggests that about 50% of all regions with low magnitude fall into *Shape 0*.

The results indicate that Euclidean distance performs better. For example, about 67% of the reported neighbors for regions in *Shape 0* are also in *Shape 0*. Because this shape size is only 31.66% of all regions, this is substantially better than picking at random.

Indeed, the results are encouraging. With a Euclidean distance, reporting near neighbors is sped up almost fivefold against a linear scan, and this improvement is likely to be even more pronounced on a large dataset. These techniques could be useful in clustering very large sets of regions efficiently, beyond the number that ChromSMS can reasonably handle, if one accepts a trade-off in quality.

Chapter 8

Conclusions

The field of epigenomics has grown tremendously over the past decade, and there has been a similar swell in the amount of data available for conducting studies. Although at small scales much of this data displays an abundance of noise, it may also encode information that can prove useful in enhancing our understanding of biology. In this thesis we developed ChromSMS, a tool whose aim is to extract this information at small scales around regulatory regions by uncovering clusters, or “shapes,” of input regions that display distinct patterns in the signal profiles of a set of chromatin marks.

We designed ChromSMS with two primary use-cases in mind. One, the variable-length mode, draws on techniques in signal processing; it warps and shifts standardized signals to expose spatial patterns in the distribution of a chromatin mark around regulatory regions. The other, multiple-mark mode, incorporates information from multiple chromatin marks at once in order to demonstrate both magnitude and spatial differences in the signal profiles across regulatory regions. While it operates in a standard manner of clustering, many of its internal details — from its measurements of dissimilarity to the way in which it merges clusters — are designed *ad hoc* with an explicit goal of capturing the biological meaning we seek. With a wealth of information it draws on and a potentially large number of input regions, ChromSMS faced a considerable computational burden. Through a variety of techniques we were able to mitigate this burden, and we believe the tool runs with time and space requirements that are reasonable for most analyses.

After implementing ChromSMS, we evaluated its effectiveness and used it to conduct analyses that relate the patterns it discovers to additional biological information. We found that these patterns display significant correlations with information that was not explicitly provided to ChromSMS, including GC-content and the locations of transcription start sites. In these cases, we cannot establish the same degree of correlation without dividing the input regions into shapes. This suggests that the shapes uncovered by the tool hold biologically meaningful information. We also developed a statistical test and observed that the patterns uncovered by ChromSMS are different from each other at a statistically significant level.


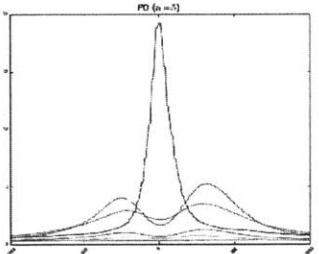

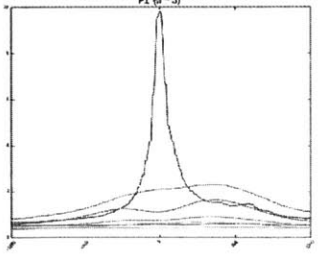
Furthermore, we performed analyses that attempted to reveal the importance of the diverse signal profile patterns we uncover, as well as the mechanisms that propel them. We detected the possibility that there are a limited number of shapes that demonstrate functional significance, notably in disease and phenotype abnormalities. In addition, we observed connections between sequence motifs and the profiles of the regions that surround them, namely that there exist distinct collections of motifs that are enriched in distinct collections of shapes. This lays the groundwork for further work in uncovering the role that particular motifs and transcription factors may play, alongside GC-content, in governing the diverse epigenomic landscape.

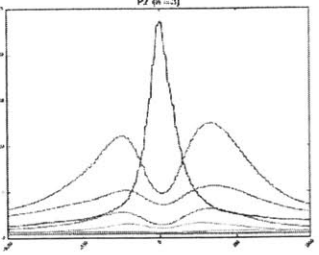
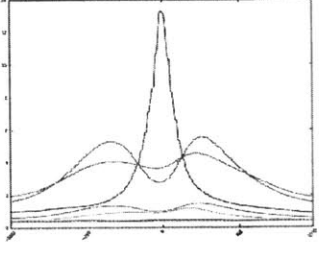
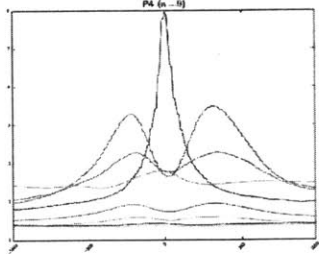
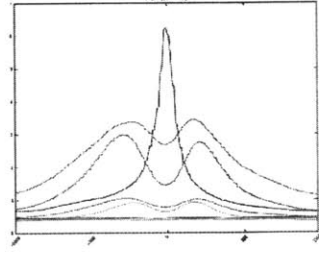
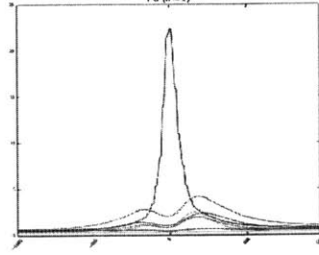
Based on our observations, we believe that ChromSMS is a powerful tool for discovering high-resolution epigenomic patterns within regulatory regions. As the field progresses, we suspect that ChromSMS could play an important role in aiding studies that further advance the field.


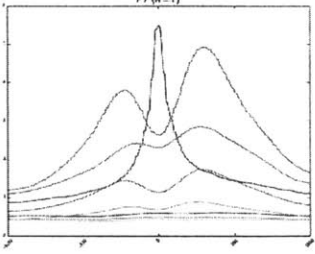

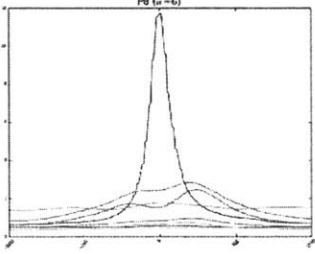

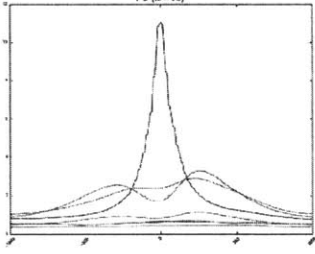

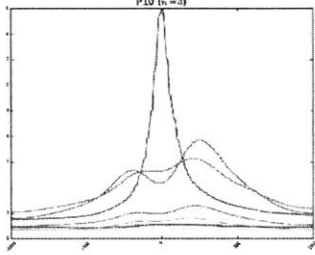

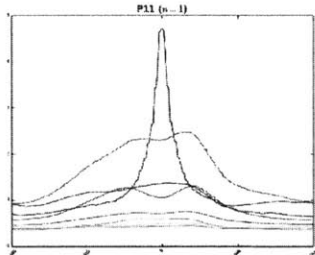
Appendix A


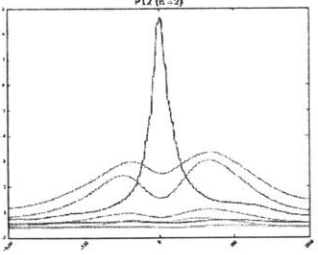

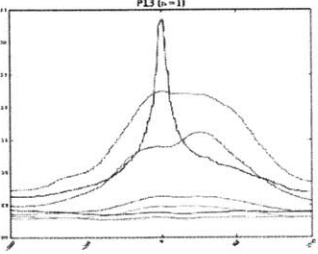

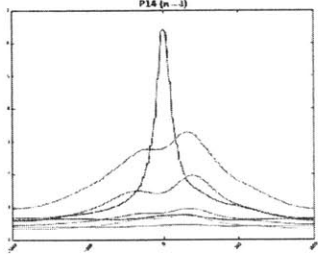

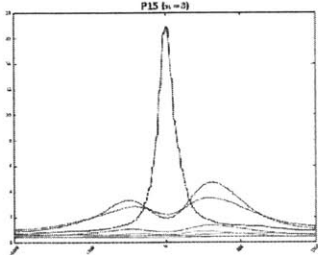

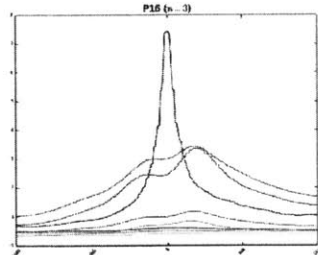
Tables


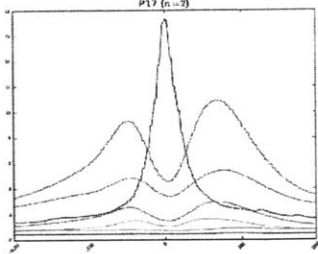

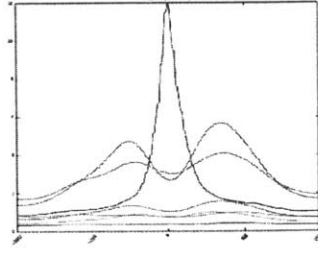
Table A.1: This table lists details about each of the 19 clusters of motif instances. The representative motif is the instance that has the great average similarity to the other motifs in the cluster. The aggregate shape is an average over the mean signal profiles for the dense shapes associated with the cluster, with each dense shape weighted equally in the average. The measures of sequence and shape similarity are based on the average pairwise similarities in the cluster. The number of motifs gives the number of motif instances in the cluster, and the number of shapes corresponds to the number of dense shapes associated with the cluster.

Cluster num	Representative motif (avg similarity to other motifs)	Aggregate shape	Sequence similarity / Shape similarity	Num motifs / Num shapes
0	 <p>TBP_3 (0.471)</p>		0.390 / 0.277	380 / 5
1	 <p>UNCX_3 (0.783)</p>		0.530 / 0.417	347 / 3

2	<p>TGACGTGG ATF6_1 (0.387)</p>		0.323 / 0.410	303 / 3
3	<p>ACCGGAAT ETV2_1 (0.589)</p>		0.461 / 0.224	140 / 7
4	<p>TGACCTTG NR2F1_1 (0.658)</p>		0.476 / 0.227	103 / 9
5	<p>AAAGGTTAA ZNF652_1 (0.492)</p>		0.428 / 0.389	77 / 3
6	<p>GCGGCC E2F1_21 (0.452)</p>		0.403 / 1.000	46 / 1

7	 <p>CAcCTG TCF3_7 (0.905)</p>		0.741 / 0.254	21 / 7
8	 <p>CCGC TFAP2C_2 (0.430)</p>		0.512 / 0.250	19 / 6
9	 <p>ACATCTG TAL1_5 (0.807)</p>		0.651 / 0.190	15 / 10
10	 <p>AAGGTC ESRR_A_3 (0.951)</p>		0.716 / 0.417	13 / 3
11	 <p>CACGTC HEY2_1 (0.908)</p>		0.611 / 1.000	10 / 1

12	 AR_5 (0.391)		0.438 / 0.532	9 / 2
13	 RXRG_2 (0.336)		0.411 / 1.000	9 / 1
14	 TFAP2A_11 (0.977)		0.602 / 0.321	8 / 4
15	 PAX3_3 (0.389)		0.521 / 0.368	7 / 3
16	 NR1H_3 (0.287)		0.460 / 0.430	6 / 3

17	 E4F1_1 (0.311)		0.435 / 0.500	5 / 2
18	 NFY_5 (0.345)		0.526 / 1.000	4 / 1

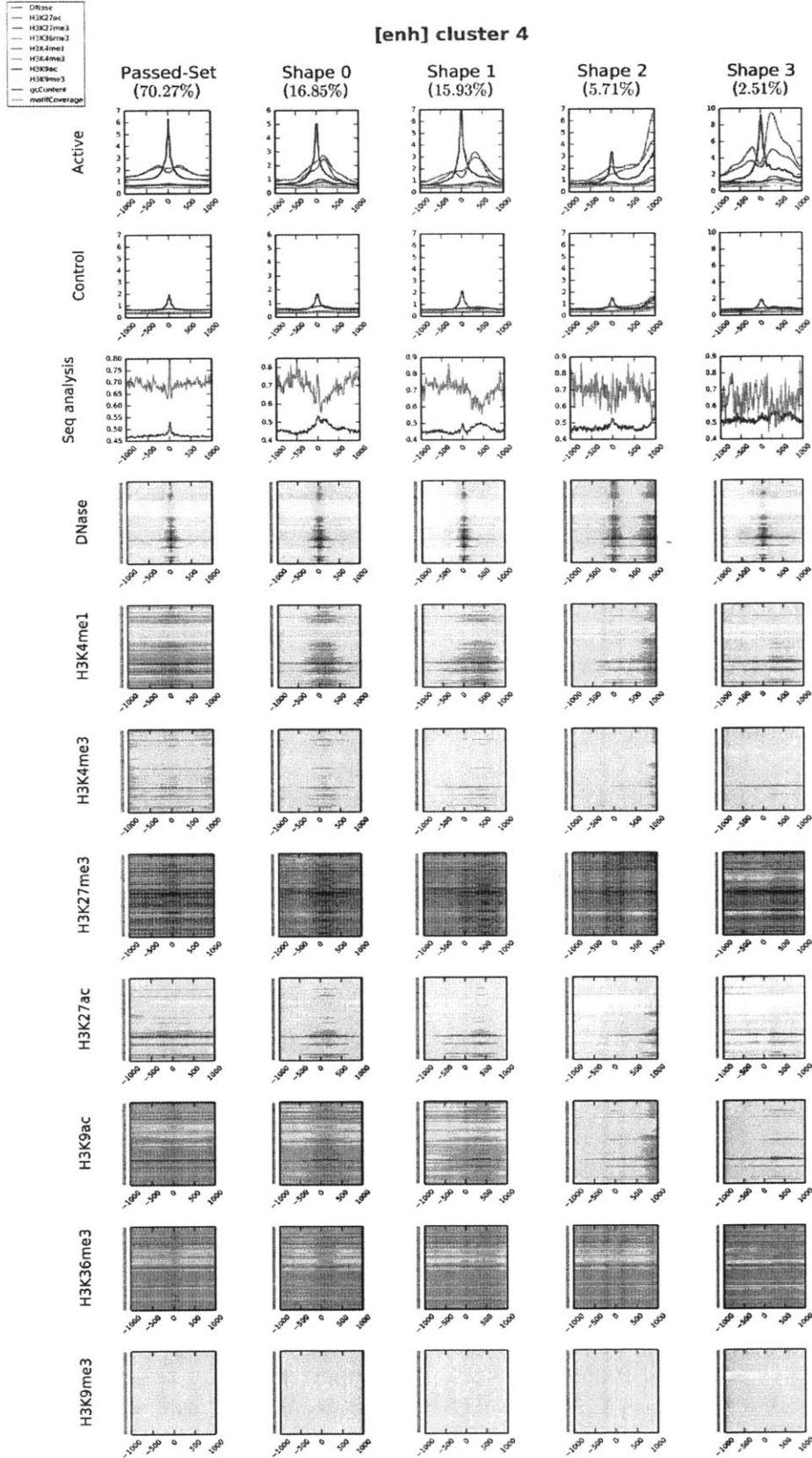
Legend for aggregate shape plots	<ul style="list-style-type: none"> — DNase — H3K27ac — H3K27me3 — H3K36me3 — H3K4me1 — H3K4me3 — H3K9ac — H3K9me3
----------------------------------	---

Appendix B

Figures

Figure B-1: (Begins on next page.) Each of these figures show detailed output from ChromSMS, run in multiple-mark mode, on an enhancer activity cluster. The activity cluster number is indicated in the title of the figure (e.g., “[enh] cluster 4”). The first column represents the *Passed-Set*, and the other columns represent each of the largest shapes discovered by ChromSMS, ordered from left to right in decreasing order of their size. The value directly under each column name gives the percentage of all regions selected by ChromSMS that are clustered into each of the shapes. (Or, in the case of *Passed-Set*, the percentage of all selected selection regions that pass the signal profile filtering stage.) In each case, the plots display signals averaged across the regions of the shape. The top row (“Active”) shows the aggregate plots, for each of eight chromatin marks, where the signal (before being averaged across regions of the shape) is taken as the average of the signals for the epigenomes that are **active** in the activity cluster. In the second row (“Control”), the signals are taken as the average of the signals in all the available epigenomes that are *not* active in the activity cluster. The third row (“Seq analysis”) shows two signals: one displays sequence motif coverage averaged across the regions in the shape, and the other displays GC-content averaged across the regions of the shape. Each of the last eight rows corresponds to a chromatin mark used by ChromSMS. They each display a heatmap for each shape in which each row of the heatmap corresponds to a cell type. (The rows of the heatmap are ordered in the same order as the rows of Figure 2-3.) The row for a cell type displays enrichment, averaged across the regions in a shape, solely in the signal representing that cell type’s epigenome. The cell types are along the row labels of these heatmaps, and the ones in red are active in the activity cluster.

(a)



(b)

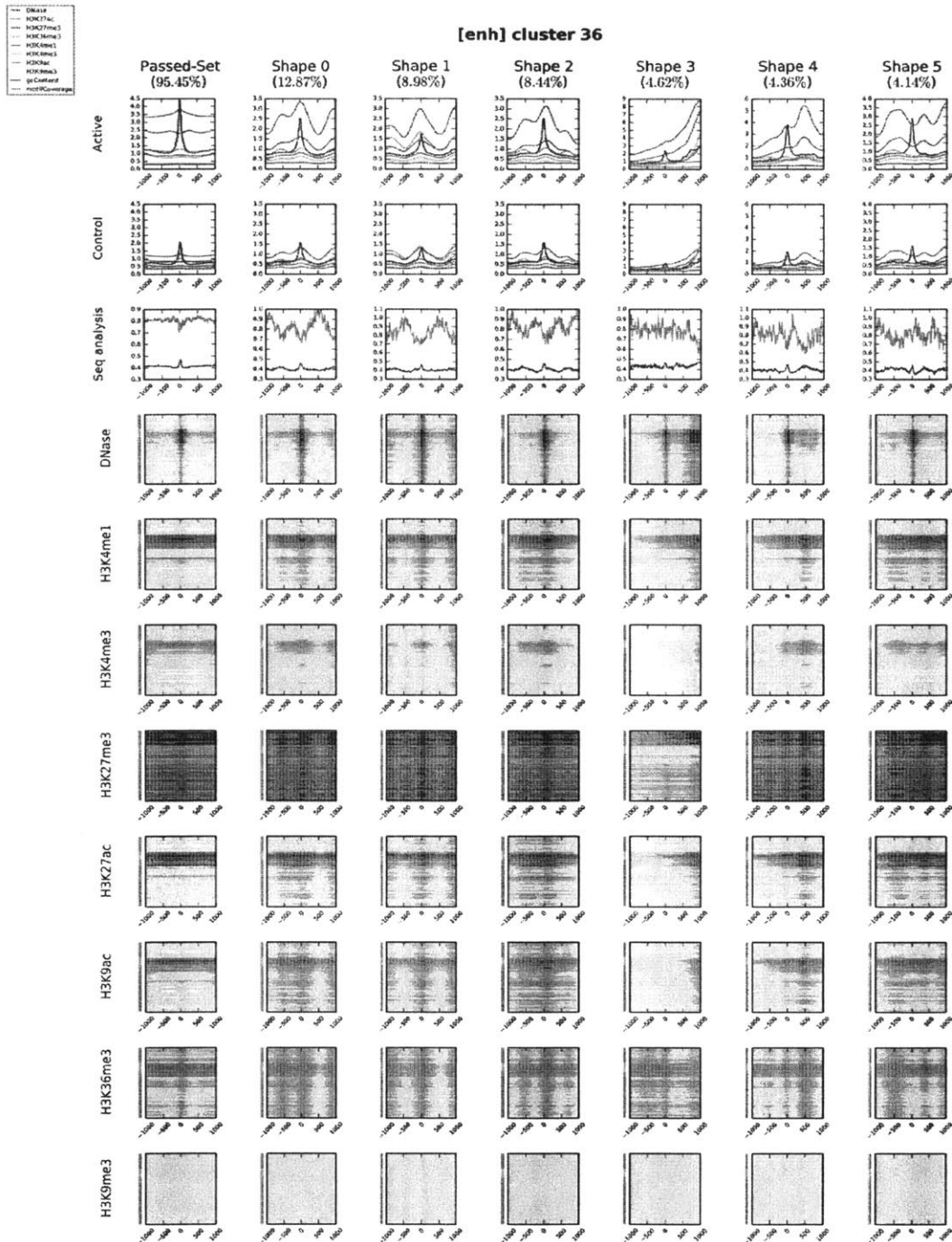
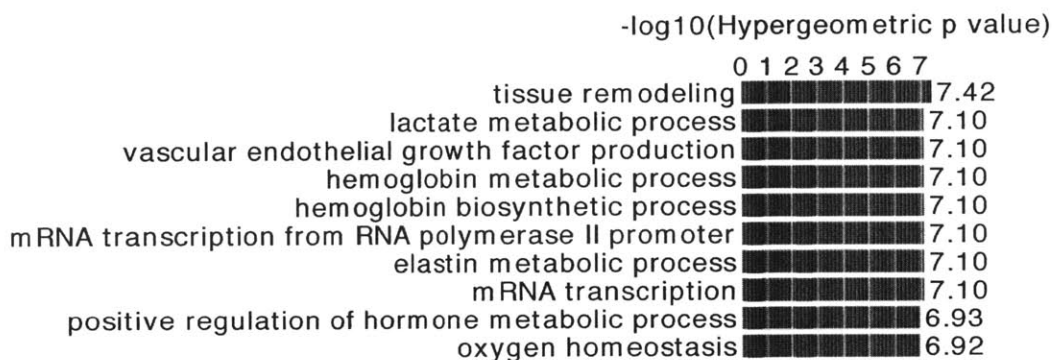
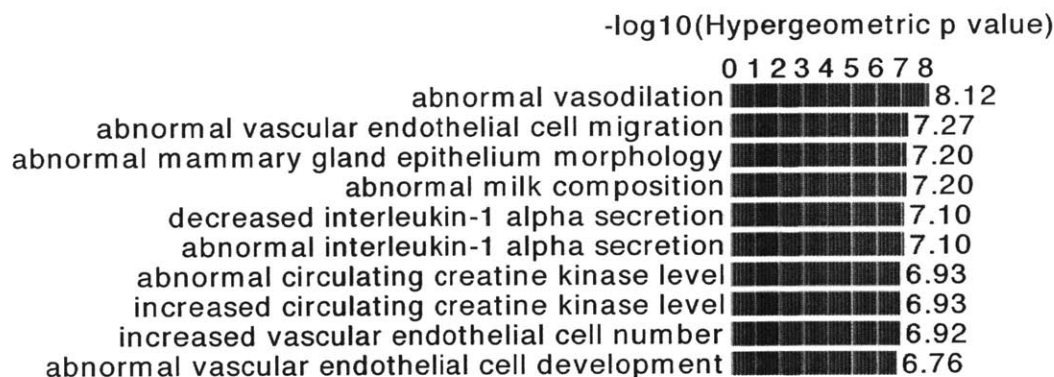


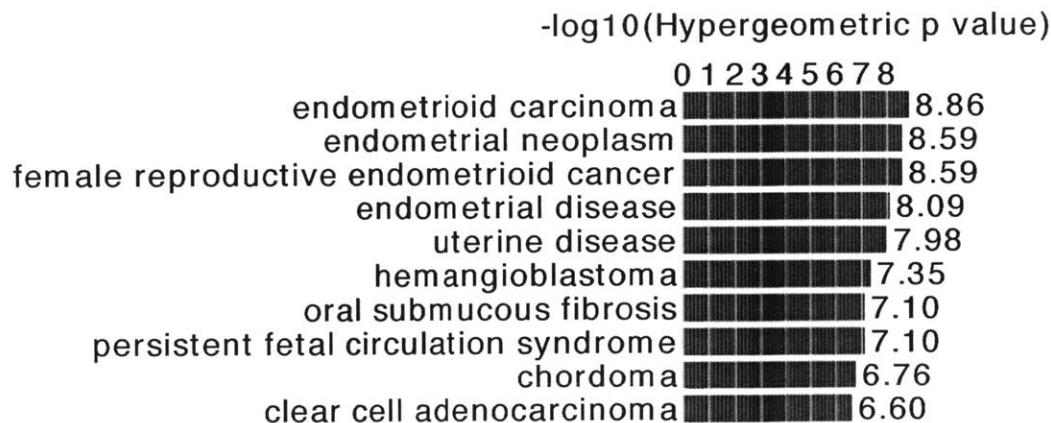
Figure B-2: These figures display the highest ten ranked terms associated with activity cluster 36, *Shape 10* for three ontologies. This shape has enrichment for 96 GO biological process terms, 190 mouse phenotype terms, and 111 disease terms. Each term is listed alongside its $-\log_{10}$ p-value. These figures were generated by the tool GREAT.



(a) GO Biological Processes



(b) Mouse Phenotype



(c) Disease Ontology

Bibliography

- [1] Ghazi Al-Naymat, Chawla Sanjay, and Javid Taheri. SparseDTW: a novel approach to speed up dynamic time warping. In ACM, editor, *The 2009 Australasian Data Mining*, volume 101, pages 117–127, 2009.
- [2] Alexandr Andoni et al. *Nearest Neighbor Methods in Learning and Vision: Theory and Practice*, chapter Locality-sensitive hashing scheme based on p-Stable distributions (ed. T. Darrell et al.). MIT Press, 2006.
- [3] Alexandr Andoni and Piotr Indyk. Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions. *Communications of the ACM*, 51(1):117–122, January 2008.
- [4] Michael Ashburner et al. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nature Genetics*, 25(1):25–29, 2000.
- [5] Teresa K. Barth and Axel Imhof. Fast signals and slow marks: the dynamics of histone modifications. *Trends in Biochemical Sciences*, 35(11):618–626, November 2010.
- [6] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society*, 57(1):289–300, 1995.
- [7] Bradley E. Bernstein et al. The NIH roadmap epigenomics mapping consortium. *Nature Biotechnology*, 28(10):1045–1048, October 2010.
- [8] Bradley E. Bernstein, Alexander Meissner, and Eric S. Lander. The mammalian epigenome. *Cell*, 128(4):669–681, February 2007.
- [9] Adrian Bird. The essentials of DNA methylation. *Cell*, 70(1):5–8, July 1992.
- [10] Christopher M. Bishop. *Pattern recognition and machine learning*. Springer, 2006.
- [11] Judith Blake et al. The mouse genome database genotypes::phenotypes. *Nucleic Acids Research*, 37:712–719, January 2009.
- [12] Jason D. Buenrostro et al. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome positioning. *Nature Methods*, 10(12):1213–1218, December 2013.

- [13] Editorial. Making a mark. *Nature Biotechnology*, 28(10):1031, October 2010.
- [14] David Eppstein. Fast hierarchical clustering and other applications of dynamic closest pairs. In *SODA '98: Proceedings of the ninth annual ACM-SIAM symposium on discrete algorithms*, pages 619–628, 1998.
- [15] Jason Ernst and Manolis Kellis. Discovery and characterization of chromatin states for systematic annotation of the human genome. *Nature Biotechnology*, 28(8):817–825, July 2010.
- [16] Jason Ernst and Manolis Kellis. ChromHMM: automating chromatin-state discovery and characterization. *Nature Methods*, 9(3):215–216, February 2012.
- [17] Romain Fenouil et al. CpG islands and GC content dictate nucleosome depletion in a transcription-independent manner at mammalian promoters. *Genome Research*, 12:2399–2408, December 2012.
- [18] Daniel J. Gaffney et al. Controls of nucleosome positioning in the human genome. *PLOS Genetics*, 8(11), November 2012.
- [19] Aaron D. Goldberg, C. David Allis, and Emily Bernstein. Epigenetics: a landscape takes shape. *Cell*, 128(4):635 – 638, February 2007.
- [20] Diana C. Hargreaves and Gerald R Crabtree. Atp-dependent chromatin remodeling: genetics, genomics, and mechanisms. *Cell Research*, 21(3):396–420, March 2011.
- [21] Jennifer Harrow et al. GENCODE: the reference human genome annotation for The ENCODE Project. *Genome Research*, 22(9):1760–1774, September 2012.
- [22] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer, 2nd edition, 2009.
- [23] Thomas Jenuwein and C. David Allis. Translating the histone code. *Science*, 293(5532):1074–1080, August 2001.
- [24] P.A. Jones et al. Moving AHEAD with an international human epigenome project. *Nature*, 454:711–715, August 2008.
- [25] Eamonn J. Keogh and Michael J. Pazzani. Scaling up dynamic time warping for datamining applications. In ACM, editor, *Proceedings of the sixth ACM SIGKDD international conference on knowledge discovery and data mining*, pages 285–289, 2000.
- [26] Peter V. Kharchenko, Michael Y. Tolstorukov, and Peter J. Park. Design and analysis of ChIP-seq experiments for DNA-binding proteins. *Nature Biotechnology*, 26(12):1351–1359, December 2008.

- [27] Tony Kouzarides. Chromatin modifications and their function. *Cell*, 128(4):693–705, February 2007.
- [28] Anshul Kundaje et al. Ubiquitous heterogeneity and asymmetry of the chromatin environment at regulatory elements. *Genome Research*, 22:1735–1747, September 2012.
- [29] Stephen G. Landt et al. ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia. *Genome Research*, 22(9):1813–1831, September 2012.
- [30] Ben Langmead, Cole Trapnell, Mihai Pop, and Steven L. Salzberg. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biology*, 10(3), March 2009.
- [31] Bing Li, Michael Carey, and Jerry L. Workman. The role of chromatin during transcription. *Cell*, 128(4):707–719, February 2007.
- [32] Cory Y. McLean et al. GREAT improves functional interpretation of cis-regulatory regions. *Nature Biotechnology*, 28(5):495–501, 2010.
- [33] Saul B. Needleman and Christian D. Wunsch. A general method application to the search for similarities in the amino acid sequence of two proteins. *Journal of Molecular Biology*, 48(3):443–453, March 1970.
- [34] John David Osborne et al. Annotating the human genome with Disease Ontology. *BMC Genomics*, 10, July 2009.
- [35] G. Robertson et al. Genome-wide profiles of STAT1 DNA association using chromatin immunoprecipitation and massively parallel sequencing. *Nature Methods*, 4(8):651–657, August 2007.
- [36] Peter Robinson and Stefan Mundlos. The human phenotype ontology. *Clinical Genetics*, 77(6):525–534, June 2010.
- [37] Richard I. Sherwood et al. Discovery of directional and nondirectional pioneer transcription factors by modeling DNase profile magnitude and shape. *Nature Biotechnology*, 32(2):171–179, February 2014.
- [38] Aravind Subramanian et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences*, 102(43):15545–15550, October 2005.
- [39] Jeff Vierstra et al. Coupling transcription factor occupancy to nucleosome architecture with DNase-FLASH. *Nature Methods*, 11(1):66–72, January 2014.
- [40] Conrad H. Waddington. The epigenotype. *Endeavour*, pages 18–20, 1942. Reprinted in *The International Journal of Epidemiology* 2012;41:10-13.

- [41] Jie Wang et al. Sequence features and chromatin structure around the genomic regions bound by 119 human transcription factors. *Genome Research*, 22(9):1798–1812, 2012.
- [42] Zhibin Wang et al. Combinatorial patterns of histone acetylations and methylations in the human genome. *Nature Genetics*, 40(7):897–903, July 2008.
- [43] Gabriel E. Zentner and Steven Henikoff. Regulation of nucleosome dynamics by histone modifications. *Nature Structural and Molecular Biology*, 20(3):259–266, March 2013.
- [44] Yong Zhang et al. Model-based analysis of ChIP-seq (MACS). *Genome Biology*, 9(9), September 2008.