

**Epigenetic determinants of cellular differentiation,
transcriptional reprogramming, and human disease**

by

Khoi Thien Nguyen

B.S., Bioengineering
University of Washington (2013)

Submitted to the Department of Biological Engineering
in partial fulfillment of the requirements for the degree of

Doctor of Philosophy in Biological Engineering

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

May 2020

© Massachusetts Institute of Technology 2020. All rights reserved.

Signature redacted

Author

Department of Biological Engineering

April 3, 2020

Signature redacted

Certified by

.....

Manolis Kellis

Professor

Signature redacted Thesis Supervisor

Accepted by

Katharina Ribbeck

Associate Professor

Chair, Biological Engineering Graduate Program

Epigenetic determinants of cellular differentiation, transcriptional reprogramming, and human disease

by

Khoi Thien Nguyen

Submitted to the Department of Biological Engineering
on April 3, 2020, in partial fulfillment of the
requirements for the degree of
Doctor of Philosophy in Biological Engineering

Abstract

Much of the diversity we observe in cellular and organismal phenotypes can be attributed to epigenetic and genetic variation. DNA provides the instructions for life, while epigenetic modifications regulate which parts of the genetic information contained in DNA can be read out in a given cell and how this information is interpreted. In recent years, epigenetic and genetic variation has been profiled on a large scale with sequencing-based assays, generating many datasets to be explored. In this thesis, I present three projects which apply computational techniques to identify and characterize epigenetic mechanisms that may contribute to the regulation of phenotypic variance. First, we mine a dataset characterizing the epigenomes of diverse cell types in order to discover signatures of adult stem cell differentiation. We identify a novel marker of the multipotent state, a chromatin state characterized by the histone marks H3K36me3 and H3K9me3, and describe biological processes that may be linked to the loss of this chromatin state in fully differentiated cell types. Next, I present what we learned from profiling the epigenetic state of cells before and after transplantation into *Xenopus* oocytes, a process that transcriptionally reprograms the cells. This analysis elucidates how the initial epigenetic state of a cell influences the success of cellular reprogramming and identifies transcription factors that help regulate this process. Finally, we integrate studies measuring the effects of genetic variants on disease with studies measuring the effects of genetic variants on transcriptional and epigenetic activity. This identifies specific mechanisms underlying disease processes, and demonstrates that transcriptional and epigenetic mechanisms may independently contribute to disease pathogenesis. Together, these projects demonstrate the biological insights that can be gained from epigenetic profiling, and expand our understanding of the potential effects of epigenetic modifications.

Thesis Supervisor: Manolis Kellis
Title: Professor

Acknowledgments

I am grateful to many people for their contributions to my thesis work and my experience as a PhD student.

First, I would like to thank my research advisor, Manolis Kellis, for his support and patience, and for giving me the opportunity to be a member of the MIT Computational Biology group. The breadth of people and ideas I've encountered here has helped me grow immensely as a scientist. I would also like to thank Doug Lauffenburger, Ernest Fraenkel, and Mary Gehring, for serving on my thesis committee. Your support, mentorship and advice over the years is appreciated.

Next, I would like to thank my collaborators, who have been integral to my achievements. Even if projects didn't always go the way we envisioned, I always learned a lot. I thank Pouya Kheradpour, Andreas Pfenning, Melina Claussnitzer, Je-Yoel Cho, Kei Miyamoto, Dinesh Kumar, Yongjin Park, Hailiang Huang, Lei Hou, and Carles Boix for working with me over the years and teaching me so much.

By extension, I would also like to thank everyone else who has been a member of the Kellis lab, including but not definitely not limited to Alvin, Bob, Kunal, Angela, Xinchun, Abhishek, Max, Yaping, Laura, Xushen, Patty, Jackie, Na, Suvi, Leandro, Remy, Eloi, Maria, Irwin, Peter, Nate, and Nicola. You have all helped me along in my scientific journey in one way or another and have spent many hours laughing and commiserating with me.

Many other people have also been a source of knowledge and comradeship, and have helped me get me through the inevitable rough patches in my PhD. Thank you to the students, faculty, and staff of the Biological Engineering department, the Muddy Charles, and all of my friends. Special thanks go to the BE class of 2013, the humans (and animals!) I've had the pleasure to live with, the inspirational women in my life, and those I can count on when I need to get my mind off of science. Adrienne, Deena, Natasha, Lynna, Claire, Andee, Kelly, Ashvin, Ben, Shawn, Santi, Griffin, Isaak, Naveen: thank you.

Finally, I would like to thank my parents, Thoa and Dung, and my siblings, Danh

and Qui. They have been beside me my entire life and I am forever indebted to them for the sacrifices they've made to give me the opportunity to be where I am today. Thank you for teaching me that material (and even academic) success is not the measure of a life well lived.

Contents

1	Introduction	17
1.1	Motivation	17
2	Chromatin state signatures of cellular differentiation	21
2.1	Introduction	21
2.2	Materials and methods	23
2.2.1	Epigenome and gene expression data	23
2.2.2	Model of differentiation effect	23
2.2.3	Transitions	24
2.2.4	Gene state classification	24
2.2.5	Biological pathways	24
2.3	Results and discussion	25
2.3.1	The loss of the "ZNF/Rpts" chromatin state is a signature of late-stage cellular differentiation	25
2.3.2	Chromatin state dynamics of "ZNF/Rpts"-associated regions and genes	29
2.3.3	Changes in gene expression linked to loss of "ZNF/Rpts"	35
2.3.4	Biological processes linked to genes that lose "ZNF/Rpts"	37
2.4	Conclusion	39
3	Chromatin accessibility impacts transcriptional reprogramming in oocytes	41
3.1	Introduction	42

3.2	Methods	44
3.2.1	Experimental methods	44
3.2.2	Statistical analysis	44
3.2.3	ATAC-seq data processing and peak calling	44
3.2.4	Genomic distribution of peaks	45
3.2.5	RNA-seq data	45
3.2.6	ATAC-seq read coverage	45
3.2.7	Correlation between samples	46
3.2.8	Pathway analysis	46
3.2.9	Motif analysis	46
3.3	Results	47
3.3.1	Optimization and Evaluation of ATAC-seq for Analyzing Cells Transplanted into Oocytes	47
3.3.2	ATAC-seq reveals changes in chromatin accessibility after NT of somatic cells to <i>Xenopus</i> oocytes	50
3.3.3	Open promoters of silent genes are permissive for transcriptional activation in oocytes	52
3.3.4	Open chromatin is often maintained near genes down-regulated after NT	53
3.3.5	Oocyte-Mediated Reprogramming Involves Opening of Closed Chromatin	58
3.3.6	Transcription Factors in Oocytes Are Involved in Transcriptional Reprogramming	58
3.3.7	Closed Chromatin Is Associated with Reprogramming-Resistant Genes	61
3.3.8	Chromatin Opening Facilitates Transcriptional Reprogramming	61
3.4	Conclusions	63
4	Multi-omic mediation analysis across multiple disease traits and cell types	67

4.1	Introduction	67
4.2	Materials and Methods	69
4.2.1	Datasets and processing	69
4.2.2	Mediation analysis	70
4.2.3	Functional analysis	73
4.3	Results and Discussion	73
4.3.1	Discovering molecular mediators of disease	73
4.3.2	Every QTL type contributes to heritability explained	80
4.3.3	Genes are linked to disease through varied mechanisms	86
4.3.4	Functional dissection of GWAS loci	91
4.4	Conclusion	95
5	Conclusion	97
5.1	Summary of contributions	97
5.2	Limitations	98
5.2.1	Chromatin state signatures of differentiation	99
5.2.2	Impact of chromatin accessibility on transcriptional reprogramming	100
5.2.3	Multi-omic mediation analysis	100
5.3	Looking to the future	101
A	Supplementary information related to Chapter 2	103
B	Supplementary information related to Chapter 4	105

List of Figures

2-1	Strategy to identify epigenomic signatures of cellular differentiation	25
2-2	Effect of differentiation stage on chromatin state prevalence.	26
2-3	Abundance of the "ZNF/Rpts" in cell types, grouped by differentiation stage and cell lineage	26
2-4	Effect of differentiation stage on other chromatin features.	27
2-5	Annotation enrichment of chromatin states from the 15-state ChromHMM model	28
2-6	Differentiation-associated change in chromatin state proportion at "ZNF/Rpts" regions	30
2-7	Genes clustered by chromatin state composition	31
2-8	Relative frequency distribution of genes across chromatin state clusters	32
2-9	Chromatin state composition of repressed and active genes with "ZNF/Rpts"	33
2-10	Cluster membership of genes that lose "ZNF/Rpts" during differentiation	33
2-11	Track image of a region that switches from the "ZNF/Rpts" to the "EnhG" chromatin state	34
2-12	Gene expression changes associated with "ZNF/Rpts" loss	36
2-13	Expression of genes in active or repressed epigenetic configuration, with or without "ZNF/Rpts"	37
2-14	Gene sets enriched in genes that lose "ZNF/Rpts" during differentiation	38
2-15	Overlap between genes that lose "ZNF/Rpts" during differentiation in each cell lineage	39
3-1	ATAC-seq enables the identification of open chromatin regions	48

3-2	Correlation between mapped ATAC-seq reads for samples with varying numbers of cells	49
3-3	Chromatin accessibility changes in oocyte reprogramming	51
3-4	Track image of MEF ATAC-seq compared to reference datasets	52
3-5	ATAC-seq reads in MEFs before and after NT were compared around TSSs in groups of genes	54
3-6	ATAC-seq, DNase-seq, H3K4me3 ChIP-seq, and Pol II ChIP-seq at TSSs of reprogramming genes	55
3-7	ATAC-seq reads for MEFs before and after NT around <i>Eif4e2</i> and <i>Ift46</i>	56
3-8	Transcription factor analysis of genes downregulated after NT	57
3-9	Track image of newly produced open chromatin during reprogramming	59
3-10	Gene categories associated with newly opened chromatin peaks after NT	59
3-11	Sequence motifs enriched in open chromatin regions after NT	59
3-12	RAR Influences Transcriptional Reprogramming in Oocytes	60
3-13	ATAC-seq reads near reprogramming-resistant genes in donor MEFs .	62
3-14	Efficient transcriptional reprogramming is associated with increased chromatin accessibility	64
3-15	ATAC-seq reads around the <i>Gapdh</i> and <i>Jun</i> gene loci for donor MEFs, MEFs after NT to control oocytes, and MEFs after NT to oocytes overexpressed with <i>Toca1/Fnbp1l</i>	65
3-16	Genomic distribution of ATAC-seq peaks detected in MEFs after NT to <i>Toca1/Fnbp1l</i> -overexpressed oocytes	65
4-1	GWAS and mediation signal for Crohn's disease	75
4-2	GWAS and mediation signal for ulcerative colitis	76
4-3	GWAS and mediation signal for coronary artery disease	77
4-4	GWAS and mediation signal for Alzheimer's disease	78
4-5	GWAS and mediation signal for systemic lupus erythematosus	79
4-6	Heritability explained by each QTL study, for each GWAS	82

4-7	Heritability explained by each QTL type, compared to number of features tested	83
4-8	Enrichment of heritability explained through mediation relative to number of features tested, for each QTL type	84
4-9	Heritability explained by combinations of QTL studies, for Crohn's disease	85
4-10	GWAS significance as a function of total mediation signal	87
4-11	UMAP, by QTL	88
4-12	UMAP, by cluster	89
4-13	Score of points belonging to each cluster	90
4-14	Crohn's disease risk locus PTGER4	92
4-15	Ulcerative colitis risk locus RBM6	93
4-16	Coronary artery disease risk locus LIPA	93
4-17	Alzheimer's disease risk locus BIN1	94
4-18	Lupus risk locus ITGAX,ITGAM	94

List of Tables

A.1	Epigenomes analyzed for differentiation signatures	104
B.1	Genome-wide association studies analyzed	105
B.2	Molecular mediators associated with Crohn’s disease (top 10 per QTL type, ranked by log odds)	106
B.3	Molecular mediators associated with ulcerative colitis (top 10 per QTL type, ranked by log odds)	107
B.4	Molecular mediators associated with coronary artery disease (top 10 per QTL type, ranked by log odds)	108
B.5	Molecular mediators associated with Alzheimer’s disease (top 10 per QTL type, ranked by log odds)	109
B.6	Molecular mediators associated with systemic lupus erythematosus (top 10 per QTL type, ranked by log odds)	110

Chapter 1

Introduction

1.1 Motivation

Humans exhibit a remarkable diversity of traits, and the cells that make up each individual also have a wide variety of functions and characteristics. This diversity emerges despite the fact that DNA, the molecule that contains instructions for the development, growth, and function of organisms, is 99.9% identical in sequence between human beings and is effectively identical in all of the cells in an individual. The recent increase in the availability of DNA sequencing, RNA sequencing, and other sequencing-based assays have enabled computational biologists to study these phenomena in much greater depth.

We now understand that certain chemical modifications (to DNA and the proteins it is attached to) regulate which parts of the genetic information contained in DNA can be read out in a given cell and how this information is interpreted, thereby providing a molecular basis for differences in function between cells. These DNA-regulating chemical modifications are known as epigenetic modifications or marks, and the collection of epigenetic marks in a cell is known as its epigenome. We have characterized how a subset of epigenetic marks correlate with gene expression and the activity of gene regulatory regions (such as enhancers and promoters), but the function of some epigenetic marks remains unclear. The varying landscape of epigenetic modification has also been an area of study. The epigenome changes naturally during the course

of development, as cells differentiate and take on more specialized roles in the body, but it can also change as a response to signals from the environment. Profiling the epigenetic state of a cell type can reveal the regulatory regions that control its gene expression program and influence how it responds to stimuli.

To understand how sequence variation can contribute to phenotypic variance, researchers have profiled the genomes of many different groups of people, identifying single-nucleotide differences that correlate strongly with ancestry and traits such as disease status and height, among others. Careful examination of disease-associated variants can provide insight into the cellular components and processes that cause disease. For example, genetic association studies linked variants regulating the expression of the complement component 4 (C4) genes to schizophrenia risk, guiding the execution of experiments demonstrating that C4 expression levels affect synaptic pruning in the brain. However, for many disease-associated variants, the molecular mechanisms by which the variant leads to disease risk are still uncertain.

The epigenome can influence the effect of genetic variants and vice versa. A genetic variant is more likely to have an effect in the cell types in which it is located in an active regulatory region. On the other hand, genetic variants may also cause changes in the epigenome, as suggested by studies linking genetic variants with varying levels of epigenetic modifications in gene regulatory regions.

Despite the wealth of data, there remain many open questions about the role of epigenetics in governing phenotypic variance. In this thesis, I present three projects computationally analyzing epigenetic and genetic data that aim to address specific gaps in our knowledge in this area. The first contribution is the discovery that a combination of epigenetic marks, H3K36me3 and H3K9me3, is a signature of adult stem cell differentiation. This provides new insight into the function of these epigenetic marks and suggests a mechanism by which stem cells maintain their unique cellular state. Next, I present what we learned from profiling the epigenetic state of cells before and after transplantation into *Xenopus* oocytes, a process that transcriptionally reprograms the cells. This analysis elucidates how the initial epigenetic state of a cell influences the success of cellular reprogramming and identifies transcription factors

that help regulate this process. Finally, we integrate studies measuring the effects of genetic variants on disease with studies measuring the effects of genetic variants on transcription and epigenetic activity, suggesting specific mechanisms underlying disease processes and generally broadening our understanding of how genetic variants can lead to disease.

Chapter 2

Chromatin state signatures of cellular differentiation

2.1 Introduction

In the nucleus of eukaryotic cells, DNA is packaged with histone proteins into a complex known as chromatin. These histone proteins can be post-translationally modified in many ways, including methylation, acetylation, phosphorylation, and ubiquitination[11]. These histone modifications (also known as histone marks) may modify the structure and accessibility of chromatin, e.g. by neutralizing the charge interaction between the DNA backbone and the histone tails, or recruit other enzymes and proteins to specific regions of the genome[11]. The specific amino acid on the histone protein that is modified may also influence its function. For example, the methylation of histone H3 lysine 36 (H3K36me) are associated with transcribed chromatin, while methylation of H3 lysine 9 (H3K9me), H3 lysine 27 (H3K27me), and H4 lysine 20 (H4K20me) are associated with repressed chromatin[11]. Further, histone marks can act in combination to produce additional types of functional chromatin domains[40].

Genome-wide maps of individual histone marks can be produced using ChIP-seq[78]. In turn, data from multiple histone marks can be combined, using methods such as ChromHMM, to learn functional combinations of epigenetic modifications, or

"chromatin states"[25].

The distribution of chromatin states across the genome is influenced by environmental signals and genetic factors, and in particular, varies widely across cell lineages and developmental stages. [77, 103, 11]. Several genome-wide epigenomic changes have been associated with cellular differentiation. Embryonic stem cells are characterized by bivalent chromatin, consisting of regions marked by methylation of H3K27 and H3K4, which is thought to silence developmental genes while keeping them poised for later activation[12, 63]. The chromatin in embryonic stem cells is also to contain fewer repressive epigenetic domains (e.g. H3K27me, H3K9me, DNAm) and be generally looser than the chromatin in cells that have committed to a cell fate [100, 13, 31, 62, 104]. In adult stem cells, differentiation-associated changes in DNA methylation in the binding sites of lineage-associated transcription factors and of chromatin state at gene promoters have been identified[13, 101]. However, studies of adult stem cells have not yet examined differentiation-associated changes across both multiple epigenetic marks and multiple cell lineages.

Recently, the number of epigenetic datasets available has increased sharply, enabling us to investigate adult stem cell differentiation across across several cell lineages. The Roadmap Epigenomics Project comprehensively characterized the epigenomes of 111 cell types and tissues, generating genome-wide maps of five core histone modification marks (H3K4me3, H3K4me1, H3K36me3, H3K27me3, and H3K9me3) and of chromatin state for each of these cell types[77]. In this study, using chromatin state segmentations for 45 epigenomes generated by the Roadmap Epigenomics Consortium, including samples from neural progenitors, fetal and adult brain, mesenchymal and hematopoietic adult stem cells, and mesenchymal and hematopoietic tissues, we identify signatures of adult stem cell differentiation and explore the biological function of these signatures.

2.2 Materials and methods

2.2.1 Epigenome and gene expression data

We used epigenome data from the Roadmap Epigenomics Consortium[77], which can be found online at https://egg2.wustl.edu/roadmap/web_portal/. The samples we used in our analysis are listed in Table A.1. The differentiation stage and cell lineage label for each sample are also listed in Table A.1. We assigned these labels based on prior knowledge.

For each sample, we calculated the genome-wide coverage of each chromatin state in the core 15-state model. Genome-wide coverage of a chromatin state is defined as the number of base pairs in the genome, excluding chromosome X, that are annotated as that state. We also generated epigenomic features following this procedure with the 25-state model.

Gene expression data (in the form of an RPKM expression matrix derived from RNA-seq experiments) for 12 out of 45 of the samples we selected were available from the Roadmap Epigenomics Consortium.

2.2.2 Model of differentiation effect

For each chromatin stage feature, we fit a linear mixed model, where differentiation stage is a fixed effect (0="Multipotent", 1="Terminal") and cell lineage group is a random effect on the intercept. Chromatin state features, differentiation stage, and cell lineage group are as defined in the previous section. To calculate the significance of the effect of differentiation stage on the value of the chromatin feature, we performed a likelihood ratio test comparing this model to a null model, where there are no fixed effects and cell lineage group is a random effect on the intercept. This was implemented using the lme4 package, version 1.1-21 [8].

2.2.3 Transitions

In a given cell lineage, we selected all regions that were in the "ZNF/Rpts" chromatin state in any of the cell types belonging to that lineage. We then calculated the fraction of these regions assigned to each chromatin state in the multipotent cell types, and in the terminally differentiated cell types. In the text, we report the mean difference of these fractions and the geometric mean log fold change of these fractions.

2.2.4 Gene state classification

For each protein-coding gene in the gene expression matrix ($n=19795$), we calculated the fraction of the gene body assigned to each chromatin state in each epigenome, creating a feature matrix with (19795×45) rows and 15 columns. These feature matrix was scaled before performing k-means clustering with the Lloyd algorithm. Gene annotations are derived from the GENCODE v10 human gene annotation set, which was used by the Roadmap Epigenomics Consortium for RNA-seq processing. The body of a given gene is defined as the region between the transcription start site (TSS) and transcription end site (TES) of the longest transcript produced by the gene.

2.2.5 Biological pathways

We performed pathway enrichment analysis using gene sets from the Gene Ontology (GO) consortium, including GO Biological Process gene sets, GO Molecular Function gene sets, and GO Cellular Component gene sets. Gene sets with over 500 members were excluded. Significance of enrichment was calculated using the hypergeometric test.

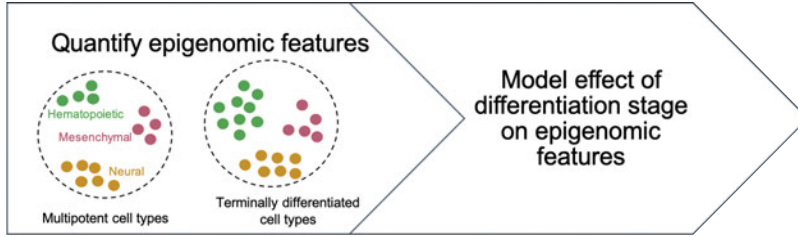


Figure 2-1: Strategy to identify epigenomic signatures of cellular differentiation.

2.3 Results and discussion

2.3.1 The loss of the "ZNF/Rpts" chromatin state is a signature of late-stage cellular differentiation

We selected epigenomes from the Roadmap Epigenomics Project that could be organized into a developmental lineage. This data set included 45 epigenomes, representing multipotent and fully differentiated samples from the neural, mesenchymal, and hematopoietic lineages. For each epigenome, we had ChIP-seq for five histone modifications (H3K4me1, H3K4me3, H3K27me3, H3K9me3, H3K36me), as well as a chromatin state segmentation learned by ChromHMM. We generate quantitative features that we can use to compare the epigenomic state of varied cell types by calculating the genome-wide coverage of each chromatin state in each epigenome. We then use a linear mixed model to find features whose value changes with differentiation stage across all cell type groups. This analysis suggested that differentiation stage has a significant effect on the abundance of the "ZNF/Rpts" state (FDR q -value $< 1e-3$), a chromatin state characterized by the presence of the histone marks H3K36me3 and H3K9me3 (Figure 2-2). The abundance of other chromatin states in the 15-state model do not significantly change with the differentiation of adult stem cells.

The effect of differentiation stage on "ZNF/Rpts" is negative, suggesting that this chromatin state becomes less common when cells differentiate (Figure 2-2). Indeed, we see that "ZNF/Rpts" abundance generally decreases from the multipotent to the terminally differentiated stage for all three cell lineages, although the effect is stronger

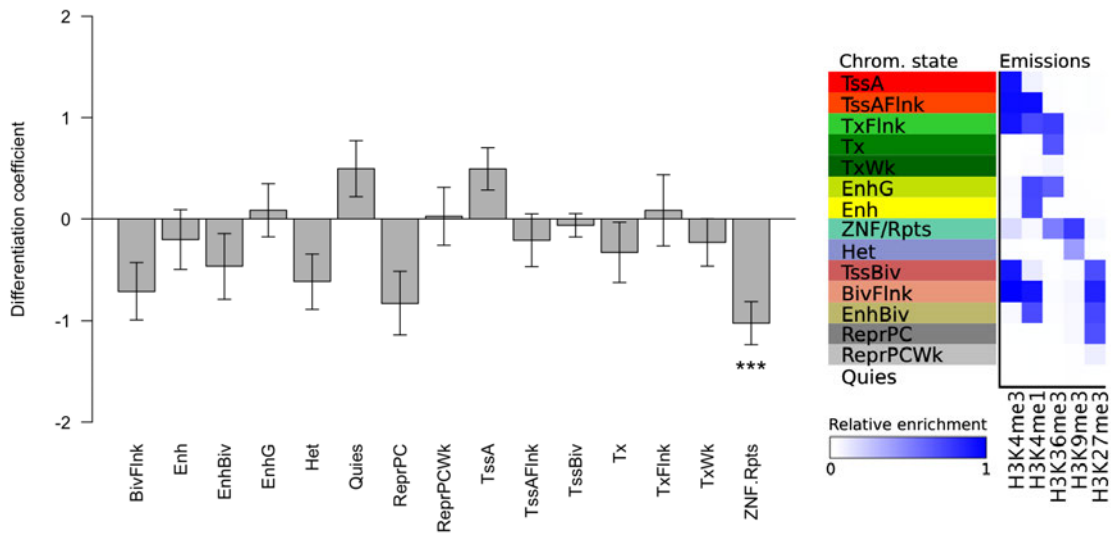


Figure 2-2: Effect of differentiation stage on the prevalence of each chromatin state, estimated using a linear mixed model. Error bars represent standard error of estimate. Stars indicate significance of differentiation effect (***: FDR q-value < 1e-3, **: q < 0.01, *: q < 0.05). Histone marks associated with each chromatin state are provided for reference on the right.

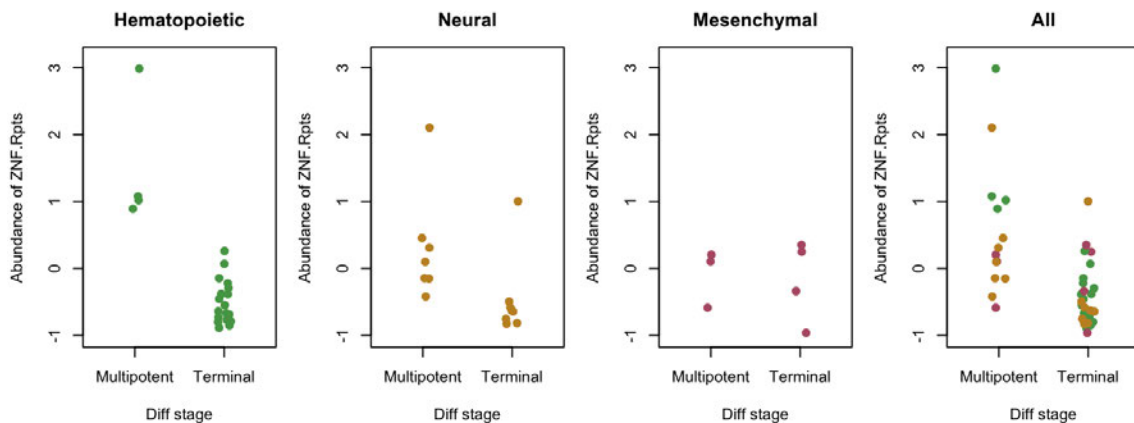


Figure 2-3: Abundance of the "ZNF/Rpts" in cell types, grouped by differentiation stage and cell lineage. Abundance y-axis is base pair coverage, scaled across epigenomes.

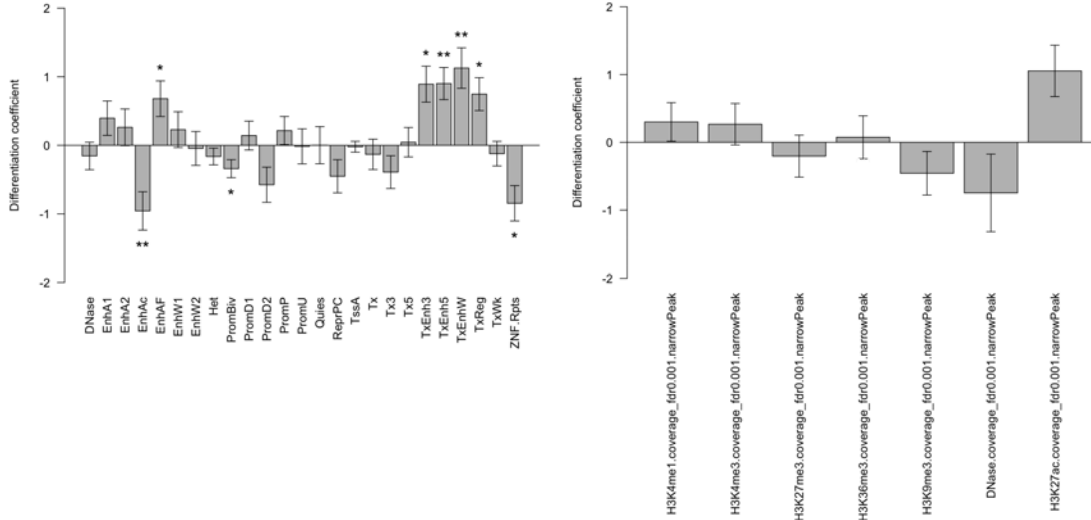


Figure 2-4: Effect of differentiation stage on other chromatin features: the prevalence of each chromatin state in the 25-state model (trained on imputed data) and on the genomic coverage of individual epigenetic marks. Effects are estimated using a linear mixed model. Error bars represent standard error of estimate. Stars indicate significance of differentiation effect (***: FDR q-value < 1e-3, **: q < 0.01, *: q < 0.05).

in the hematopoietic and neural than in the mesenchymal lineage (Figure 2-3).

The role of "ZNF/Rpts" in adult stem cell differentiation is corroborated by the results of this analysis on the chromatin states in the 25-state ChromHMM model, which is trained on a larger set of imputed epigenomic data. This analysis again showed a significant negative effect of differentiation on "ZNF/Rpts" abundance (Figure 2-4). It also suggests that the gain of transcribed enhancer states and other chromatin state changes may be associated with differentiation, though these effects are not significant in the core 15-state chromatin state model (Figure 2-2, 2-4). Further, we did not observe that differentiation is linked to changes in the genomic coverage of individual epigenetic marks (Figure 2-4). Taken together, the epigenomic signature of adult stem cell differentiation with the strongest evidence is the loss of the "ZNF/Rpts" chromatin state.

As defined by the Roadmap Epigenomics Consortium, "ZNF/Rpts" is the chromatin state characterized by H3K36me3, a histone mark associated with transcriptional elongation and transcribed regions of genes, and H3K9me3, a histone mark

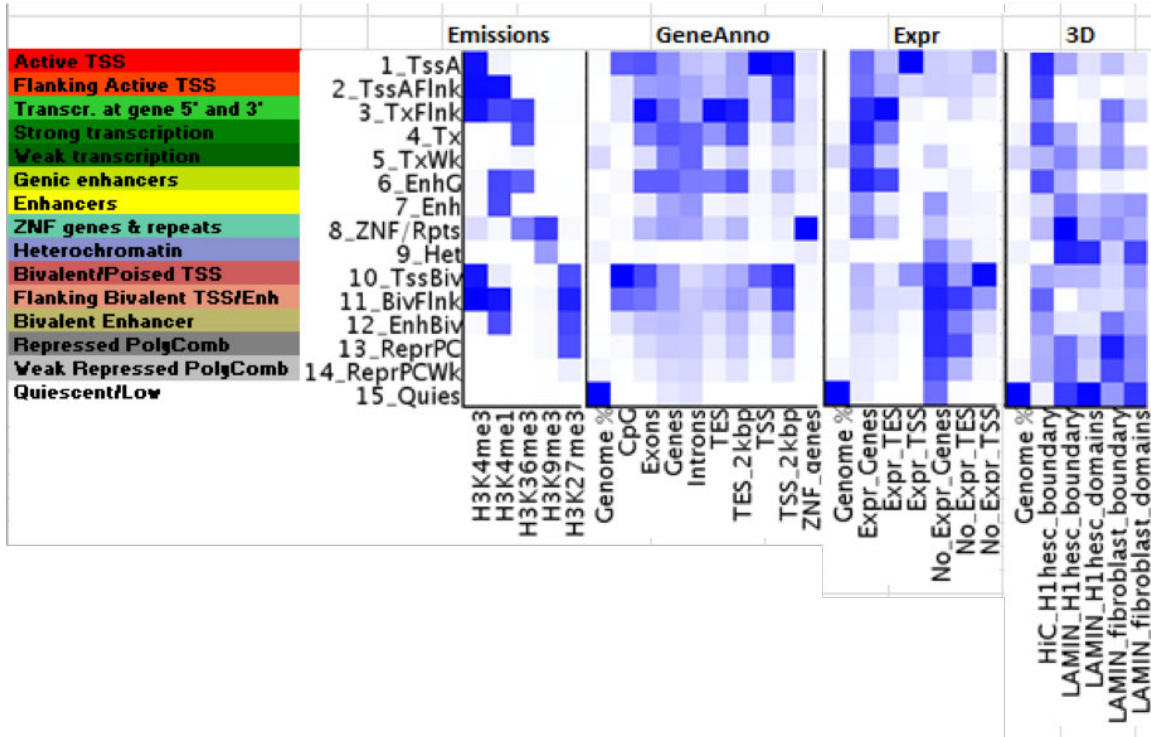


Figure 2-5: Annotation enrichment of chromatin states from the 15-state ChromHMM model, adapted from the Roadmap Epigenomics Consortium [77].

associated with heterochromatin domains [77]. It has not previously been described as a marker of differentiation. As its name suggests, the "ZNF/Rpts" chromatin state is often found near genes in the zinc finger family and/or near repetitive elements, including tandem and interspersed repeats (Figure 2-5) [77]. It is also associated weakly with both expressed and non-expressed genes (Figure 2-5).

Interestingly, this signature of adult stem cell differentiation parallels a signature of embryonic stem cell differentiation, namely the loss of the bivalent promoter state [12, 63]. Both "ZNF/Rpts" and the bivalent promoter state ("Tss-Biv" in the 15-state ChromHMM model) are characterized by both an active mark (H3K36me3 in "ZNF/Rpts", H3K4me3 in "TssBiv") and a repressive mark (H3K9me3 in "ZNF/Rpts", H3K27me3 in "TssBiv"). While we did not observe widespread loss of "TssBiv" in the context of adult stem cell differentiation (Figure 2-2), the loss of chromatin states with both active and repressive histone marks may be a general feature of cellular differentiation in both early and late stages.

2.3.2 Chromatin state dynamics of "ZNF/Rpts"-associated regions and genes

After observing that the abundance of "ZNF/Rpts" generally decreases, we explored what chromatin states are gained in its place during differentiation. In a given cell lineage, we selected all regions that had "ZNF/Rpts" in any of the cell types belonging to that lineage. We then compared the chromatin state composition of these regions in the multipotent cell types to the composition in the terminally differentiated cell types (Figure 2-6).

The abundance of "ZNF/Rpts" in these regions decreased substantially in terminally differentiated cell types, as expected. On average, approximately 15.2% of these regions lose "ZNF/Rpts" from the multipotent to terminal stage, representing a 1.97-fold decrease (Figure 2-6). In addition to the decrease in "ZNF/Rpts", the largest changes in chromatin state proportion we observed were the decrease in heterochromatin (6.8% of regions lose "Het"), the increase in quiescent regions (9.1% of regions gain "Quies"), and the increase in weakly transcribed regions (8.3% of regions gain "TxWk") (Figure 2-6, top). If we look at the fold change in chromatin state composition instead of the absolute difference, the largest changes we saw were the increase in enhancer regions ("Enh", 2.1-fold increase) and genic enhancer regions ("EnhG", 1.9-fold increase), as well as the increase in weakly repressed regions ("ReprPCWk", 1.9-fold increase) (Figure 2-6, bottom).

Next, we looked at chromatin state changes on a gene level, since the majority of "ZNF/Rpts" regions are within genes. We calculated the chromatin state composition for each gene, in each of the 45 epigenomes, and clustered these genes into 4 groups using k-means clustering. The distribution of chromatin state abundance in genes in each of these four clusters is shown in Figure 2-7. We generally see that genes are either in an active (Cluster 3 and 4) or repressed epigenetic configuration (Cluster 1 and 2). Cluster 1 genes are generally polycomb-repressed, while Cluster 2 genes are generally quiescent. Genes in both Cluster 3 and 4 are in an active epigenetic configuration, but Cluster 4 may represent shorter genes, where the promoter region

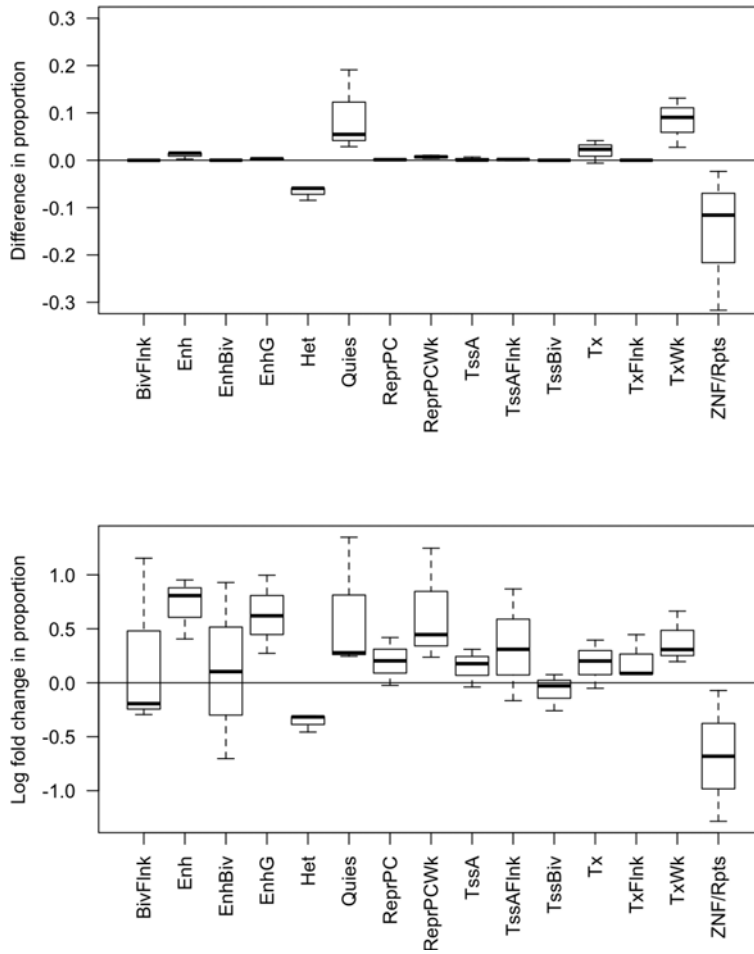


Figure 2-6: Differentiation-associated change in chromatin state proportion at "ZNF/Rpts" regions, for each cell lineage. Top: Difference between proportion of regions assigned to each chromatin state in terminally differentiated cell types, relative to multipotent cell types. The distribution of this difference across cell lineages is shown. Bottom: Log fold change in proportion of regions assigned to each chromatin state in terminally differentiated cell types, relative to multipotent cell types.

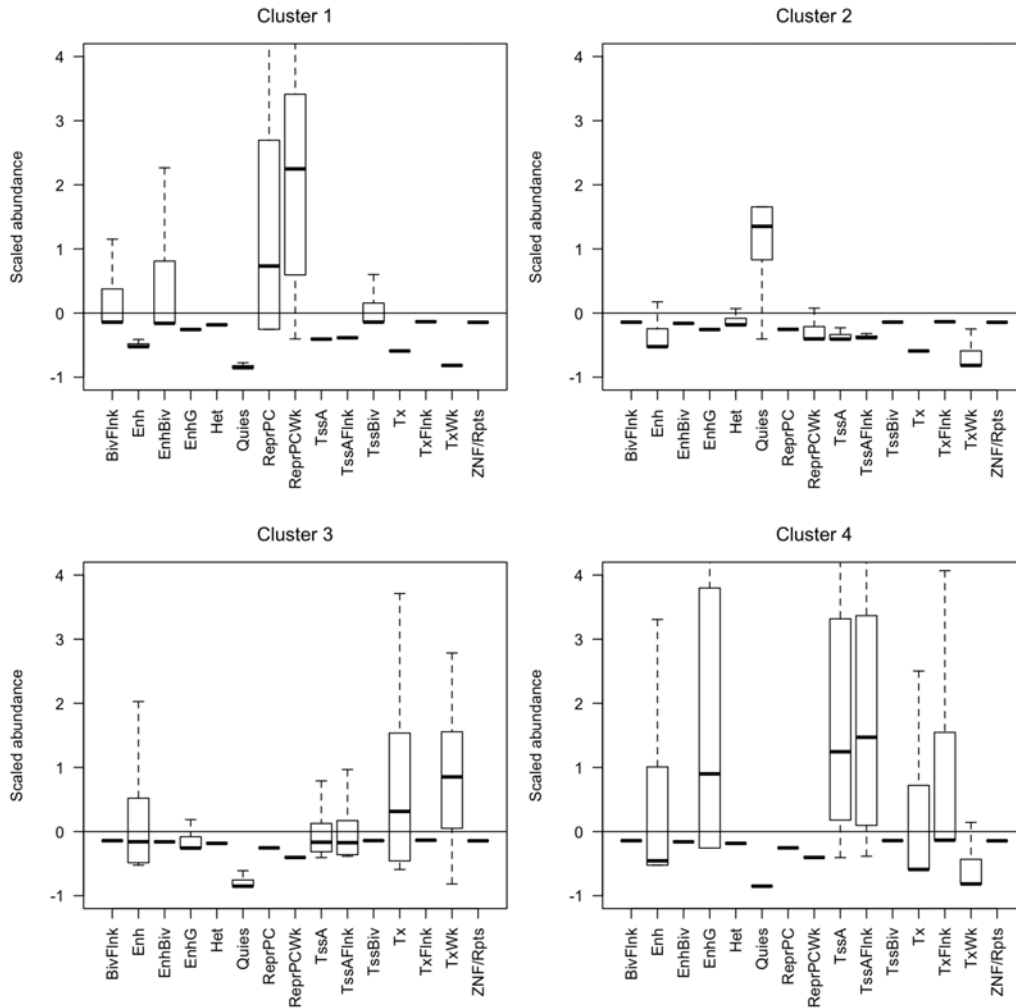


Figure 2-7: Genes clustered by chromatin state composition. For each cluster, we show the distribution of the abundance of each chromatin state in genes belonging to that cluster.

covers a larger fraction of the gene. Clusters 2 and 3 are the largest clusters, with 38.0% and 43.7% of genes in those clusters, respectively (Figure 2-8).

82.2% of genes with "ZNF/Rpts" are in Cluster 3, the main active gene cluster, which is a 1.9-fold enrichment compared to all genes (Figure 2-8). Most of the remaining genes with "ZNF/Rpts" are in Cluster 2, the quiescent gene cluster, while very few genes with "ZNF/Rpts" belong to Cluster 1 or 4.

Compared to other genes in Cluster 2, the genes with "ZNF/Rpts" in this cluster tend to have more chromatin in the "ZNF/Rpts", heterochromatin ("Het"), and weakly transcribed states ("TxWk"), and less chromatin in the quiescent ("Quies")

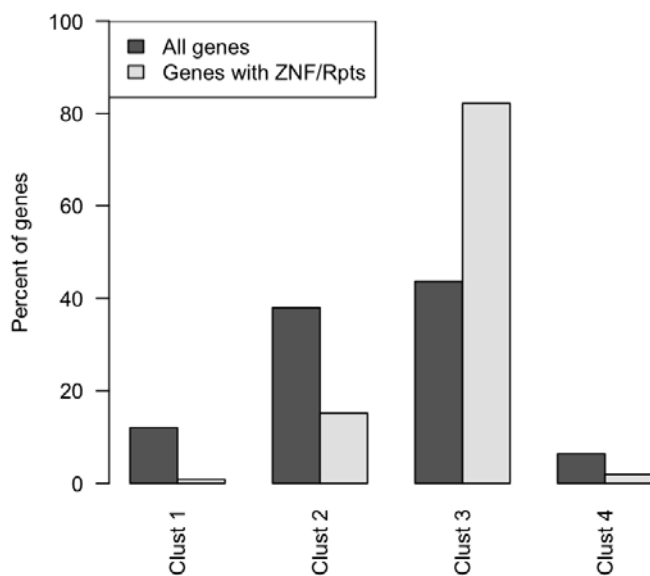


Figure 2-8: Relative frequency distribution of genes across chromatin state clusters, for all genes and for genes with "ZNF/Rpts".

state (Figure 2-9, 2-7). Compared to other genes in Cluster 3, the genes with "ZNF/Rpts" in this cluster tend to have more chromatin in the "ZNF/Rpts" state, and less chromatin in the weakly transcribed and enhancer states.

Next, we examine what happens to genes that lose "ZNF/Rpts" during differentiation in each cell lineage (switching genes). The subset of these switching genes that start in an active epigenetic configuration (Cluster 3) generally remain in active after differentiation, though about 19% of these genes become inactivated (Figure 2-10). In contrast, about 57% of the switching genes that start in a repressed epigenetic configuration (Cluster 2) remain repressed after differentiation, and about 57% become activated. An example of a region losing "ZNF/Rpts" in differentiation is shown in Figure 2-11. This genomic region on chromosome 6 is "ZNF/Rpts" in all 4 hematopoietic stem cell samples, but switches to the genic enhancer state in 18/19 blood cell types. The gene is in an active epigenetic configuration both before and after differentiation.

To sum up, these results indicate that most genes with "ZNF/Rpts" are in an

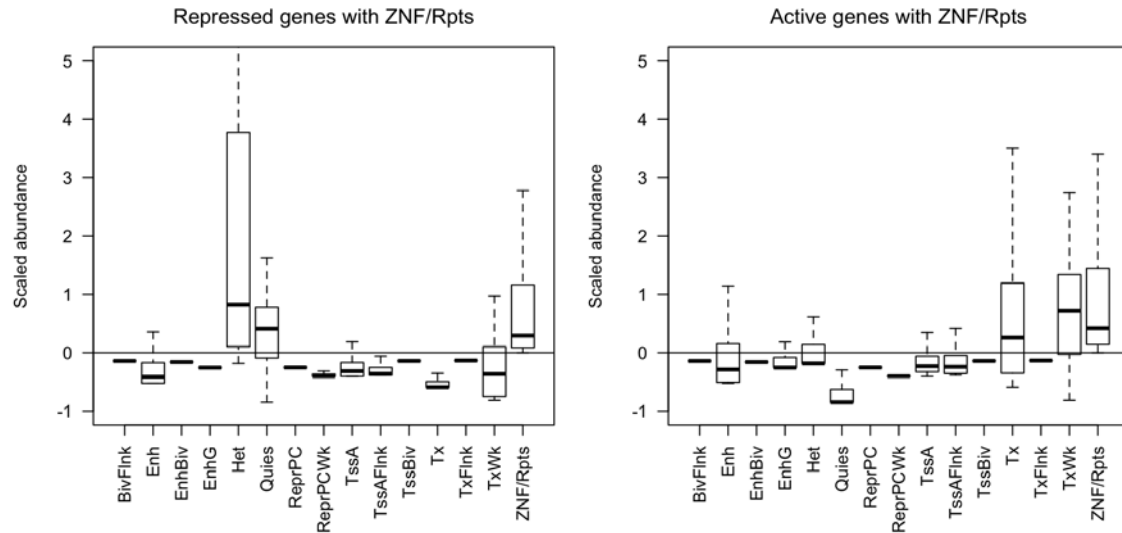


Figure 2-9: Left: chromatin state composition of genes with "ZNF/Rpts" in Cluster 2. Right: chromatin state composition of genes with "ZNF/Rpts" in Cluster 3.

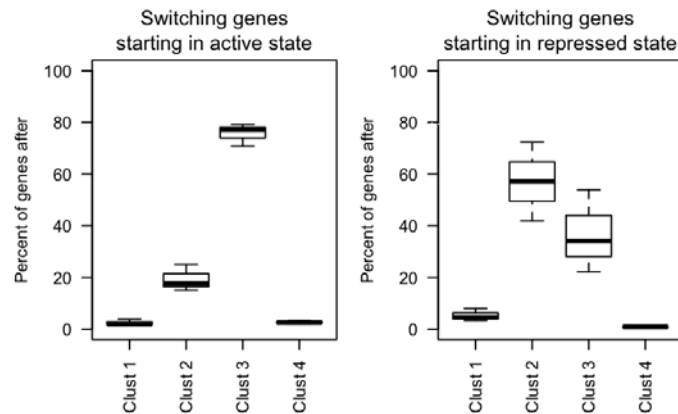


Figure 2-10: Cluster membership of genes that lose "ZNF/Rpts" during differentiation. Left: membership of genes in differentiated cell types that had "ZNF/Rpts" and were in cluster 2 in progenitor cells. Right: membership of genes in differentiated cell types that had "ZNF/Rpts" and were in cluster 3 in progenitor cells.

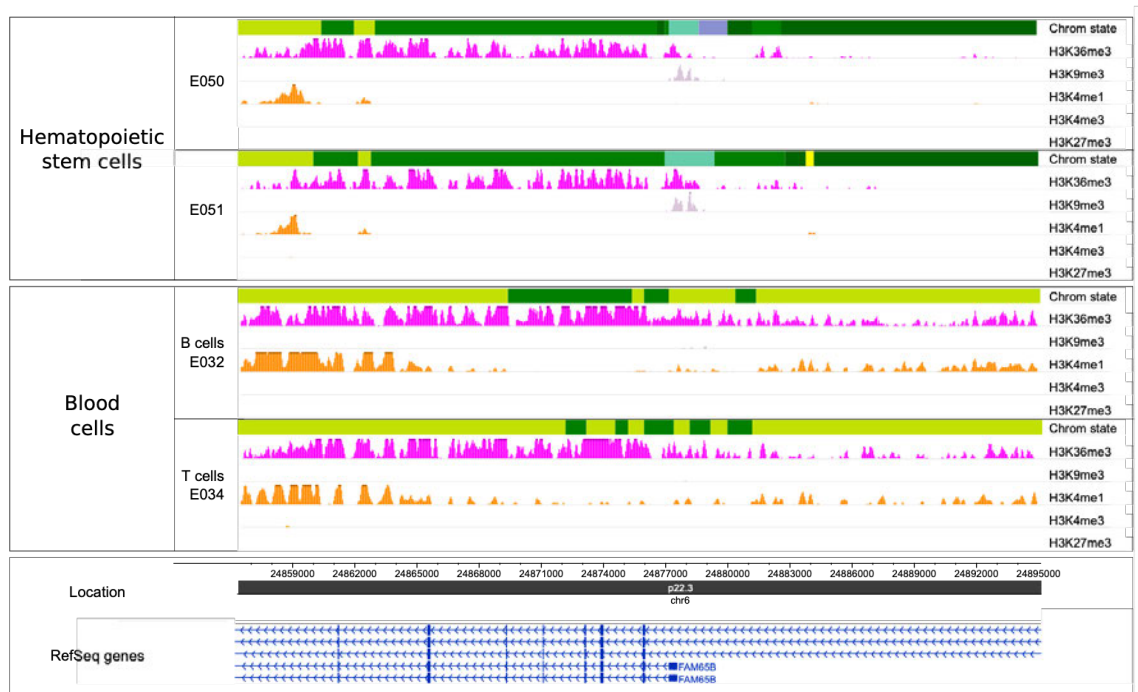


Figure 2-11: Track image of a region that switches from the "ZNF/Rpts" chromatin state (aqua color, H3K36me3 and H3K9me3) in hematopoietic stem cells (all 4 samples, 2 shown) to the "EnhG" (light green, H3K36me3 and H3K4me1) chromatin state in blood cell types (18/19 samples, 2 shown).

active epigenetic configuration. During differentiation from the multipotent to the terminal stage, many of these genes lose "ZNF/Rpts", but stay active, gaining weakly transcribed and enhancer chromatin states. However, some of these active genes do become inactive, losing "ZNF/Rpts" and gaining quiescent chromatin states.

2.3.3 Changes in gene expression linked to loss of "ZNF/Rpts"

After characterizing the epigenetic changes associated with late-stage cellular differentiation, we were interested in the functional consequences of these changes and how "ZNF/Rpts" may contribute to the maintenance of a multipotent cellular state. First, we asked whether the loss of "ZNF/Rpts" regions are linked to changes in gene expression. To do this, we categorized genes in each cell lineage based on the chromatin state dynamics we observed. First, we identified genes that have "ZNF/Rpts" in the multipotent cell types but do not have "ZNF/Rpts" in the differentiated cell types. Next, we determined the epigenetic status (active or repressed) of these genes in the multipotent cell types and in the differentiated cell types, based on the clustering described above. This defines four groups of genes that lose "ZNF/Rpts" in the context of differentiation: genes that start in the active state and remain there after losing "ZNF/Rpts", genes that start in the active state and become inactivated after losing "ZNF/Rpts", genes that start in the repressed state and become activated after losing "ZNF/Rpts", and genes that start in the repressed state and remain repressed after losing "ZNF/Rpts". For each gene, we determined the log fold change in expression in the differentiated cell types relative to the multipotent cell types (Figure 2-12). Gene expression data was available for samples in the hematopoietic and neural lineages.

On average, we found when a gene in the active state loses "ZNF/Rpts" but remains in the active state (the largest category of genes), its gene expression does not change significantly. In the hematopoietic lineage, the median log fold change in expression is approximately 0, while in the neural lineage, the median log fold change is -0.25. Genes that start in the inactive state and become activated after losing "ZNF/Rpts" also seem to not undergo drastic changes in expression, although the sample size is very small, in both the hematopoietic and neural lineages. In the

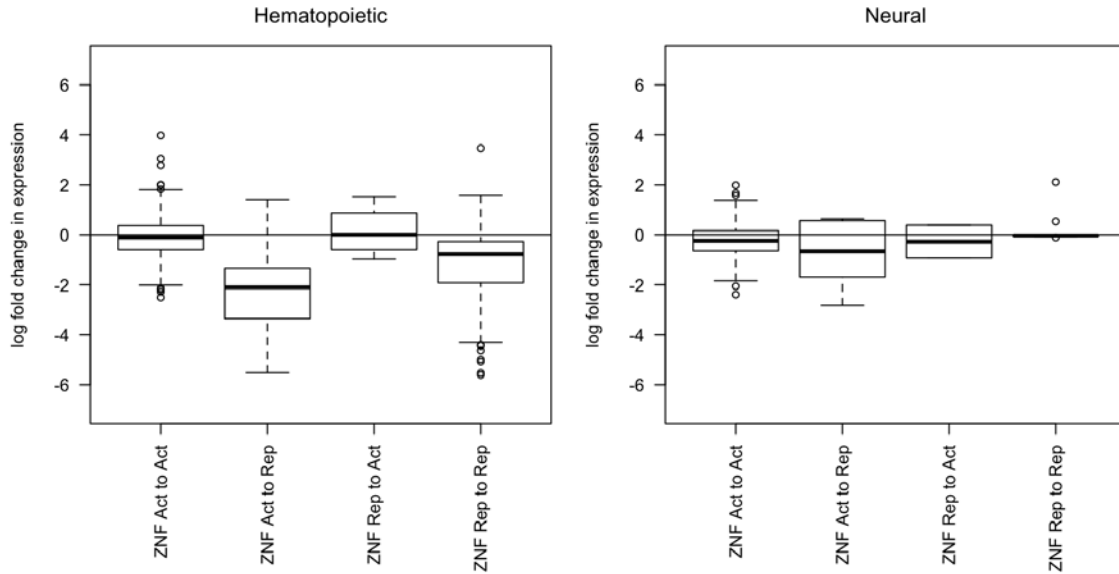


Figure 2-12: Gene expression changes associated with "ZNF/Rpts" loss. Genes that lose "ZNF/Rpts" from the multipotent to terminally differentiated stage are grouped into four categories, based on their overall epigenetic configuration before and after differentiation. Results are shown for the hematopoietic cell lineage (left) and the neural cell lineage (right). No expression data was available for terminally differentiated cell types in the mesenchymal cell lineage.

hematopoietic lineage, when a gene loses "ZNF/Rpts" and ends up in the inactive state, its expression tends to decrease, by 10.2-fold for genes that start in the active state, and by 3.3-fold for genes that start in the inactive state. In the neural lineage, there are few genes that lose "ZNF/Rpts" and end up in the inactive state, so the patterns in gene expression change are not statistically significant.

These results are consistent with gene expression data when we do not look specifically at genes that lose "ZNF/Rpts" in the context of differentiation. Genes with "ZNF/Rpts" in the active epigenetic configuration are expressed at similar levels to genes in the active epigenetic configuration overall (Figure 2-13). Genes with "ZNF/Rpts" in the repressed epigenetic configuration are transcribed at lower levels than genes in the active epigenetic configuration, but are transcribed at higher levels than other genes in the repressed epigenetic configuration (Figure 2-13).

In summary, most of the differentiation-associated loss in "ZNF/Rpts" regions are not linked to changes in gene expression. The majority of genes that lose "ZNF/Rpts"

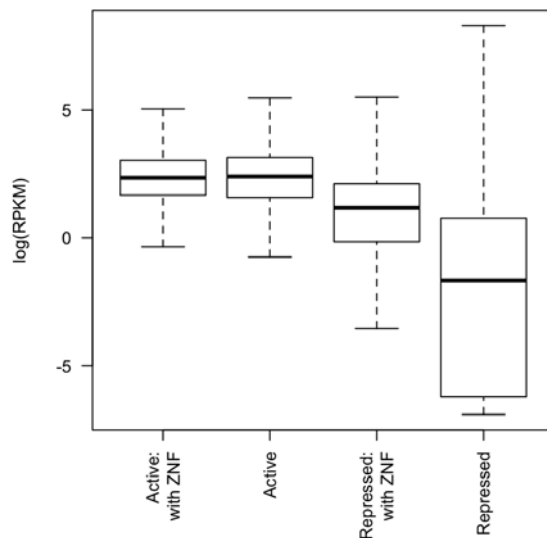


Figure 2-13: Expression of genes in active or repressed epigenetic configuration, with or without "ZNF/Rpts". All genes and epigenomes with expression data were included, not only genes that lose "ZNF/Rpts" in the context of cellular differentiation.

are in an active chromatin configuration before and after the loss of "ZNF/Rpts", and these transitions are not associated with broad changes in gene expression level.

2.3.4 Biological processes linked to genes that lose "ZNF/Rpts"

Next, we asked if the genes that lose "ZNF/Rpts" in the context of differentiation are connected to any specific biological process or function, performing Gene Ontology analysis on the genes that lose "ZNF/Rpts" in each of the three cell lineages. This analysis suggests that genes involved in ubiquitin-mediated protein degradation, cell projection assembly, organelle assembly and localization, Golgi vesicle transport, protein phosphorylation, DNA repair, and cell division may be regulated by "ZNF/Rpts" in multipotent stem cells.

The genes that lose "ZNF/Rpts" during differentiation do overlap between the three cell lineages, which may account for some of the shared pathway enrichment. Out of the 6808 genes that lose "ZNF/Rpts" in these lineages, 1357 (19.9%) lose "ZNF/Rpts" in all three cell lineages we examined (Figure 2-15). 2368 (34.8%) genes only switched from "ZNF/Rpts" in neural cells, 851 (12.5%) genes only switched

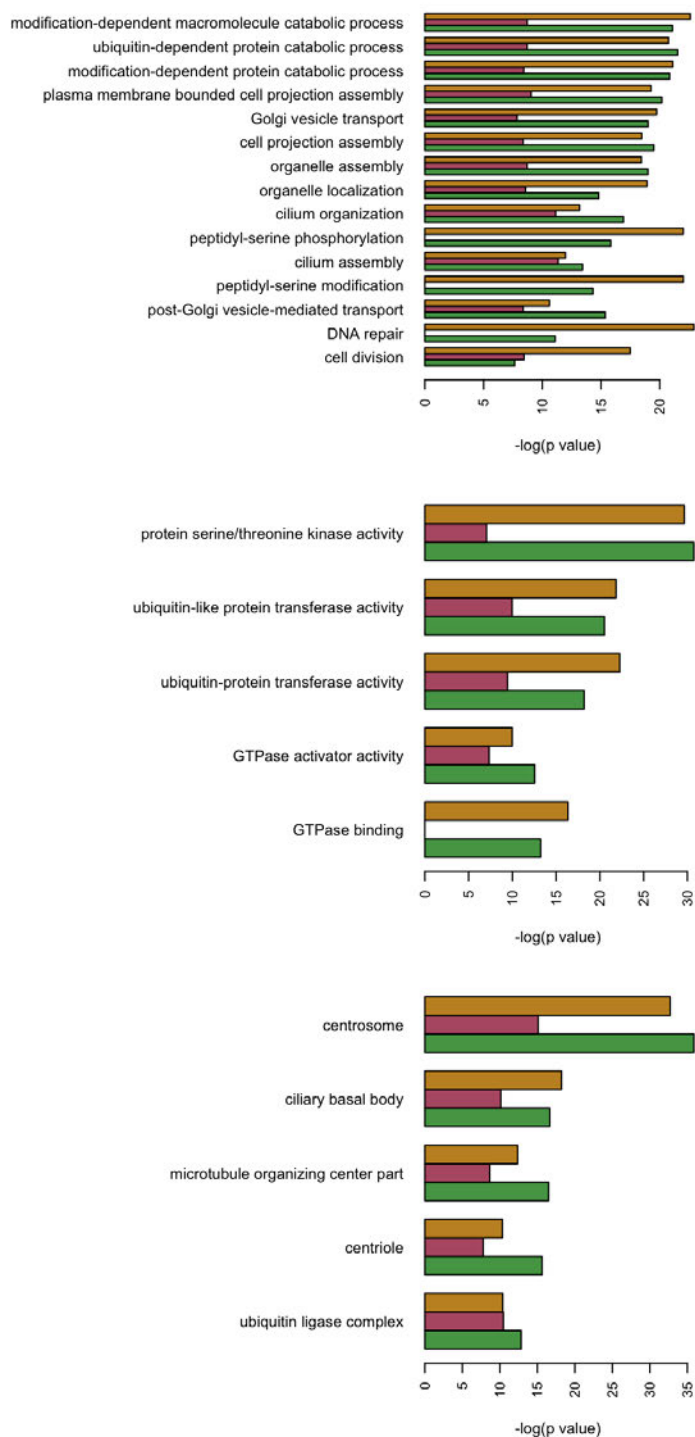


Figure 2-14: Gene sets enriched in genes that lose "ZNF/Rpts" during differentiation. For each pathway, the geometric mean of the enrichment p-values in each of the cell lineages was calculated. Top pathways with this mean p-value $< 1e-4$ are displayed, with the corresponding p-values in the cell lineages (neural: gold, mesenchymal: pink, hematopoietic: green). Top: GO Biological Process gene sets, middle: GO Molecular Function gene sets, bottom: GO Cellular Component gene sets. Gene sets with over 500 members were excluded.

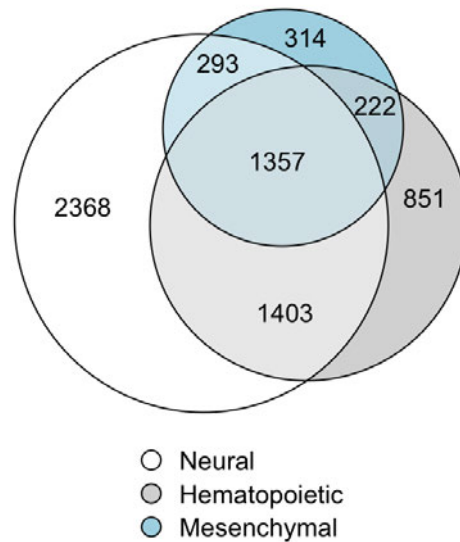


Figure 2-15: Overlap between genes that lose "ZNF/Rpts" during differentiation in each cell lineage.

from "ZNF/Rpts" in hematopoietic cells, and 314 (4.6%) genes only switched from "ZNF/Rpts" in mesenchymal cells, meaning 3533 (51.9%) genes only switch in one cell lineage.

Taken together, "ZNF/Rpts" seems to specifically mark certain groups of genes in multipotent cell types, related to cell functions such as ubiquitin-mediated protein degradation, organelle assembly and localization, and others. These cellular processes may play roles in stem cell proliferation and renewal.

2.4 Conclusion

Epigenetic modifications, including histone modifications and DNA methylation, help determine cell phenotype. Previous studies have identified changes in specific epigenetic modifications that occur during differentiation in some cell lineages. Here, we analyzed chromatin state segmentations generated from genome-wide maps of five histone modification marks in 45 epigenomes, including samples from neural progenitors, fetal and adult brain, mesenchymal and hematopoietic adult stem cells, and mesenchymal and hematopoietic tissues, and identified a chromatin state signature

of adult stem cell differentiation that applies to multiple cell lineages.

We found that the loss of the "ZNF/Rpts" chromatin state, which is characterized by the H3K9me3 and H3K36me3 histone modifications, is associated with adult stem cell differentiation. This is the first time that this chromatin state has been linked to differentiation. This chromatin state is primarily found on specific groups of actively transcribed genes in multipotent cell types, associated with biological processes including protein degradation and organelle assembly. We observed that most of these genes remain active after losing "ZNF/Rpts" during differentiation and do not significantly change in expression. The functional role of "ZNF/Rpts" in multipotent cells is still unclear, since it does not appear to regulate gene expression levels. Future work may elucidate how the "ZNF/Rpts" chromatin state contributes to the maintenance of a multipotent cellular state.

Chapter 3

Chromatin accessibility impacts transcriptional reprogramming in oocytes

Somatic cells are generally highly specialized, with stable gene expression programs. However, when somatic cells are transplanted into oocytes, they undergo dramatic changes in gene expression, turning off some genes and activating others. We used ATAC-seq to investigate the changes in chromatin accessibility that occur as a result of oocyte transplantation, and explored how chromatin accessibility and transcriptional reprogramming are linked¹. In somatic cells, we found that silent genes with accessible promoters are more likely to be transcriptionally activated after transplantation than silent genes without pre-existing open chromatin. While some regions of closed chromatin opened to enable transcription, closed chromatin seems to act as a barrier to complete transcriptional reprogramming. We also predicted transcription factors that may regulate this system by examining regions of accessible chromatin before and after transplantation. We demonstrate experimentally that one of our predicted TFs affects transcriptional activation, and that enhancing chromatin accessibility globally aids transcriptional reprogramming. In conclusion, chromatin

¹This work has been published in Miyamoto et al.[66]. The text of this chapter has been adapted from the published paper.

accessibility plays a crucial role in determining the outcome of oocyte-mediated transcriptional reprogramming.

3.1 Introduction

In the nucleus of eukaryotic cells, DNA is packaged with proteins into a complex known as chromatin. Some sections of the DNA are packaged in such a way that nuclear machinery, such as RNA polymerase and transcription factors, cannot access the DNA and initiate transcription and other processes[48]. The regions of the genome that are accessible change during cell state transitions. For example, chromatin remodeling is observed in mouse embryonic development, which is accompanied by altered transcriptional activities of the genes associated with the remodeled regions[99, 56]. Thus, the particular pattern of chromatin accessibility in a cell helps define its gene expression program and identity, and studying chromatin accessibility dynamics provides insight into cell fate changes.

Chromatin accessibility across the genome can be measured with several techniques, including DNase sequencing (DNase-seq) and the Assay for Transposase-Accessible Chromatin Sequencing (ATAC-seq). DNase-seq is a widely used method for assaying accessibility, which identifies genomic regions that can be cut by the DNaseI enzyme[15, 84, 90]. ATAC-seq is a more recently developed method, which identifies regions where a hyperactive Tn5 transposase is able to insert a short sequence into the DNA[16]. ATAC-seq produces maps of chromatin accessibility similar to DNase-seq, but requires far fewer cells as input[16].

One type of cell fate transition that is of particular interest is the transformation of a differentiated, somatic cell to a pluripotent cell. During the development of a multicellular organism, a pluripotent cell gradually differentiates into many types of somatic cells, each with a characteristic gene expression program. This establishment of somatic cell identity was generally considered permanent until various discoveries demonstrated that the process of differentiation could be reversed with experimental interventions[34]. Understanding how somatic cells can be reprogrammed to a

pluripotent state may enable multiple applications, including the development of individualized cell therapies[34, 88].

Several methods for reprogramming differentiated somatic cells to a pluripotent state have been established[34, 88, 43]. In induced pluripotent stem cell (iPSC) technology, specific transcription factors (TFs) are over-expressed in somatic cells. Over time and many cell divisions, the cells acquire a pluripotent state[88, 43]. In oocyte nuclear transfer, somatic cell nuclei are transplanted into the nucleus of a prophase amphibian oocyte. Proteins and other components present in the oocytes flow into the transplanted nuclei, interacting with the somatic chromatin and causing pluripotency genes to be transcriptionally activated in the somatic nuclei within a day, without cell division reprogramming[43, 42, 69]. Other methods include egg nuclear transfer, cell fusion, and exposing somatic cells or nuclei to cell extracts from pluripotent cells[34, 43].

The dynamics of chromatin accessibility during the generation of induced pluripotent stem cells (iPSCs) have been extensively studied. During this process, the accessibility of gene regulatory regions is altered, especially at the binding sites for the reprogramming TFs[52, 97]. The binding of these reprogramming TFs to chromatin is inhibited by a closed chromatin conformation marked by histone H3 lysine 9 trimethylation (H3K9me3), but enhancing chromatin accessibility during reprogramming by knocking down the histone chaperone CAF-1 accelerates TF binding and the transcriptional activation of pluripotency genes[82, 19].

However, the chromatin accessibility changes that occur during oocyte nuclear transfer have not yet been characterized. In this study, we use ATAC-seq to identify open chromatin regions in mouse somatic nuclei before and after nuclear transfer into *Xenopus laevis* oocytes, and explore how these changes in chromatin accessibility are linked to transcriptional reprogramming.

3.2 Methods

3.2.1 Experimental methods

For information on the experimental methods (animal protocols, ATAC-seq, cell culture, oocyte nuclear transfer, and oocyte transcriptional profiling), see the Experimental Procedures and Supplemental Experimental Procedures sections of our paper[66].

3.2.2 Statistical analysis

The number of biological replicates is shown as n. In transcriptional assays by qRT-PCR, the statistical difference was calculated by the two-tailed Student's t test. Error bars were represented as the standard error of the mean (SEM). The levels of significance were set as $*p < 0.05$. The statistical difference in mean ATAC-seq signal around groups of genes was calculated using the two-tailed Mann-Whitney U test. The statistical significance of the overlap between ATAC-seq peaks and other genomic datasets was calculated using the hypergeometric test. These statistical tests were conducted using the R software package (version 3.4)[76].

3.2.3 ATAC-seq data processing and peak calling

ATAC-seq data were processed as previously described[22]. In brief, we mapped the ATAC-seq reads to the mm9 genome assembly using Bowtie2 (version 2.2.1) with the parameter `-X 2000`[51]. Next, we filter out all reads meeting any of the following criteria: unmapped, mate unmapped, not primary alignment, multi-mapped, or duplicates marked by Picard MarkDuplicates (version 1.1)[2]. MACS2 (version 2.1) was used to call peaks on these reads, using the parameters `--nomodel --shift 75 --extsize 150 --keep-dup all`, and a minimum p-value threshold of 0.01[102]. To call higher confidence peaks, for a given cell type, we use only peaks called from the pooled data supported by both replicates. Next, we filtered out peaks landing in regions determined by the ENCODE Consortium to frequently have artificially high sequencing signal (blacklist regions)[21] Finally, we selected only those peaks with

an FDR value less than 0.01. An implementation of this protocol is available online (https://github.com/kundajelab/atac_dnase_pipelines).

3.2.4 Genomic distribution of peaks

The mm9 RefSeq Genes annotation available from the UCSC table browser was used to define all genomic features except enhancers[73, 46]. If a gene had multiple transcription start sites (TSSs), all were included in the analysis. Enhancer and super-enhancer regions for embryonic stem cells and myotubes were previously defined by mapping regions of key transcription factor binding [98]. MEF enhancers are from the mouse ENCODE project (<http://chromosome.sdsc.edu/mouse/download.html>)[80]. ATAC-seq peaks were also compared to myotube H3K4me3 ChIP-seq peaks[6]. Reported enrichment is $\frac{a/b}{c/d}$, where a is the number of peaks overlapping a given genomic feature, b is the number of total peaks, c is the number of regions corresponding to the feature, and d is the estimated number of discrete regions in the genome where the peaks and feature could overlap. Specifically, d is equal to $\frac{\text{genome size}}{\text{mean peak size} + \text{mean feature size}}$, following the implementation in the bedtools fisher software tool[75]. Significance of overlap between datasets was calculated using the hypergeometric test.

3.2.5 RNA-seq data

In Jullien et al. 2014, gene expression was measured in donor MEFs before and after reprogramming in *Xenopus* oocytes, and genes that are activated or silenced after NT were identified[42]. In Jullien et al. 2017, genes that were resistant to transcriptional reprogramming (reprogramming-resistant genes) were identified[44].

3.2.6 ATAC-seq read coverage

Only reads remaining after all filtering steps described above are used for visualization (track images and heatmaps), and to calculate read coverage around TSSs. Plotted reads of ATAC-seq, DNaseI-seq (GEO accession number: GSM1014172), and ChIP-seq (GEO accession numbers: GSM769029, GSM918761) around TSSs were visualized

as heatmaps using the Seqplots R package (version 1.10.2)[86].

3.2.7 Correlation between samples

For each sample, we counted the number of reads in a 1,000 bp window centered on each TSS. The correlation coefficient was calculated between pairs of samples (Pearson or Spearman, as noted in the text). For hierarchical clustering, we used average linkage and a distance metric $d = 1 - correlation$.

3.2.8 Pathway analysis

We used the Genomic Regions Enrichment of Annotations Tool (GREAT, version 3.0.0) to perform pathway enrichment analysis to lists of regulatory regions[60]. We used the default settings that are implemented on the web interface for GREAT (<http://bejerano.stanford.edu/great/public/html/>). Namely, to assign regions to genes for Gene Ontology enrichment analysis, each gene is assigned a basal regulatory domain of a 5kb upstream and 1kb downstream of the TSS. The gene regulatory domain is extended in both directions to the nearest gene's basal domain but no more than 1000kb in one direction. Curated regulatory domains were also included. The whole genome was used as a background for the enrichment.

3.2.9 Motif analysis

The HOMER tool (version 4.7, <http://homer.ucsd.edu/>) was used to calculate enrichment of known motifs[33]. The set of known vertebrate motifs included with HOMER (264 motifs total) were tested for enrichment in our data.

3.3 Results

3.3.1 Optimization and Evaluation of ATAC-seq for Analyzing Cells Transplanted into Oocytes

We developed a modified ATAC-seq protocol developed for our experiments (Figure 3-1A). We performed this protocol on mouse C2C12 myoblast samples with varying numbers of cells (50,000, 1,000, 100, 10, and 1). In samples where cell numbers were small (less than 100 cells), many reads were derived from mitochondria and the produced libraries showed characteristics of those with a low quality, such as duplicated reads and low library complexity. Hierarchical clustering of the mapped ATAC-seq reads around transcription start sites (TSSs) showed that the 1,000-cell and 50,000-cell samples clustered together while those created from 100, 10, and 1 cell did not cluster with the 1,000-cell and 50,000-cell samples (Figure 3-2). These results suggest that ATAC-seq reads that cover the whole genome can only be obtained from libraries prepared from 1,000 or more cells, at least using our reported method. Therefore, we decided to prepare ATAC-seq libraries with at least 1,000 cells in subsequent experiments.

Genomic locations of ATAC-seq peaks, representing open chromatin sites, were examined in C2C12 cells, and 52,697 unique sites were found. Open chromatin regions were over-represented within 1 kb of TSSs by 15.9-fold relative to the whole genome ($p < 1E-6$) and by 13.1-fold in myotube enhancers ($p < 1E-6$) (Figure 3-1B). Further, 40.0% of all TSSs and 40.3% of all myotube enhancers ($p < 1E-6$) contained ATAC-seq peaks. ATAC-seq peaks were 10.9-fold enriched in regions marked with histone H3 lysine 4 trimethylation (H3K4me3) on a genome-wide scale ($p < 1E-6$; Figure 3-1C, track image), in agreement with the previously noted association of H3K4me3 with open and active promoter elements[28, 32]. These data verify that our modified ATAC-seq protocol captures characteristic open chromatin features.

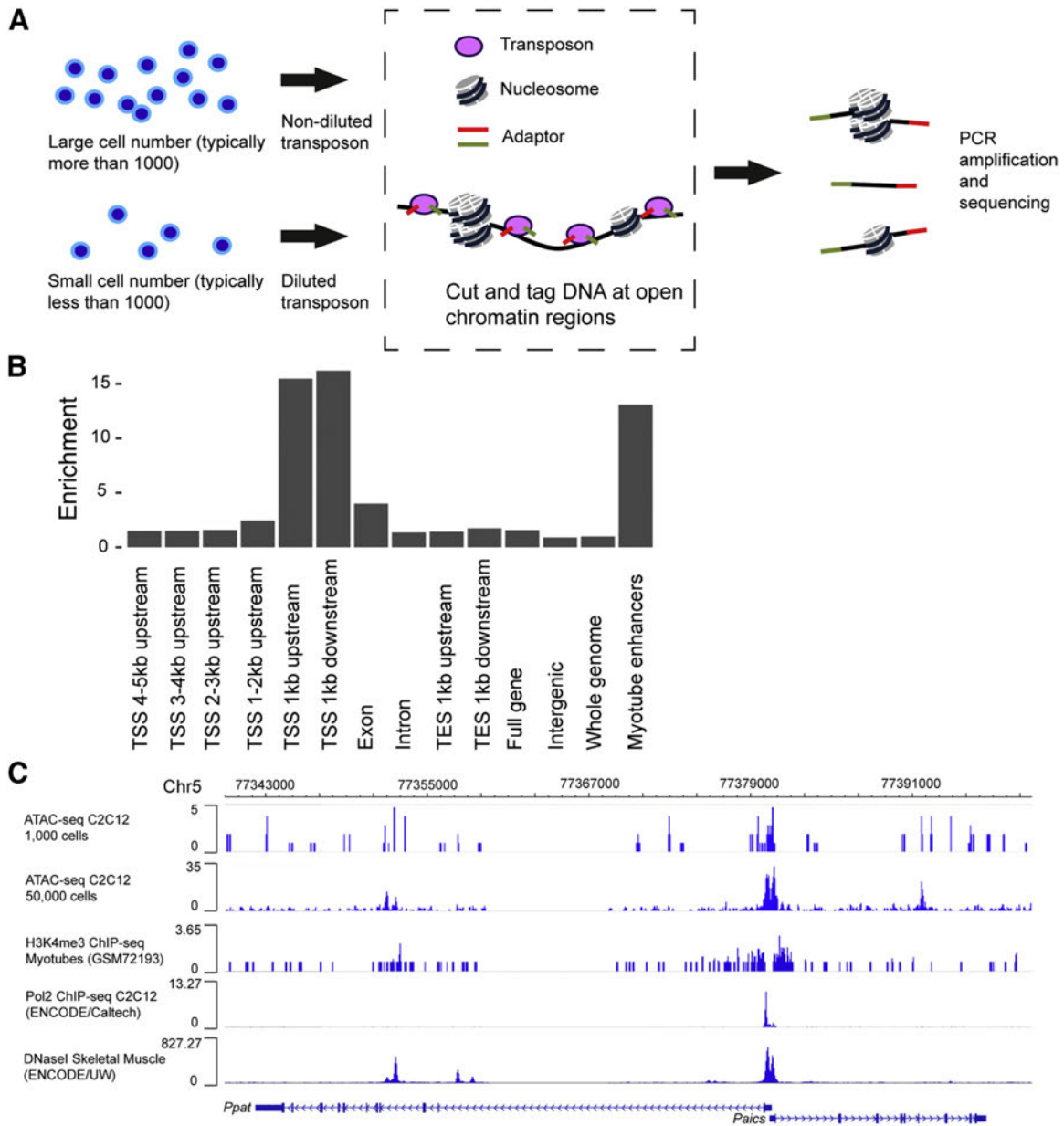


Figure 3-1: ATAC-seq enables the identification of open chromatin regions. A) The modified ATAC-seq protocol for our experiments. B) The genomic distribution of ATAC-seq peaks in C2C12 mouse myoblasts, representing open chromatin. The y axis represents the enrichment of peaks in each type of genomic region relative to the whole genome. Two independently prepared ATAC-seq libraries were used for the analysis. C) A track image of ATAC-seq from 1,000 and 50,000 cells, chromatin immunoprecipitation sequencing (ChIP-seq) of H3K4me3 (GSM72193) and RNA polymerase II (GSM915176), and DNase-seq (GSM1014189) at the *Ppat* and *Paics* genes in mouse muscle cells. Open chromatin at the TSS is adjacent to H3K4me3 marks.

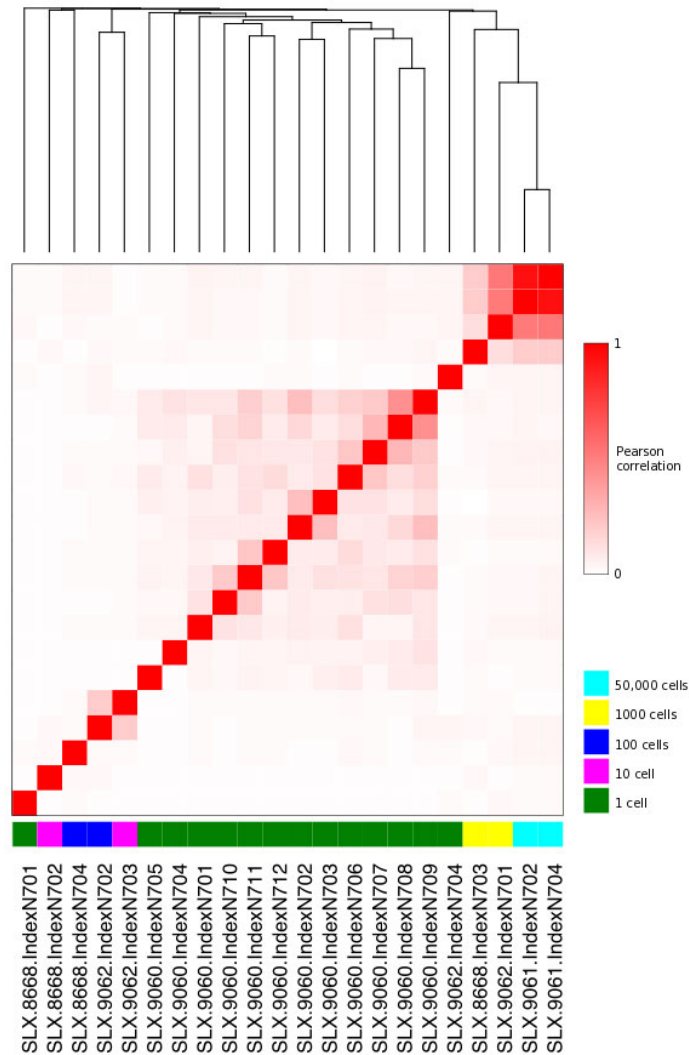


Figure 3-2: Correlation between mapped ATAC-seq reads in C2C12 mouse myoblasts when 1, 10, 100, 1,000 and 50,000 cells were analyzed (22 samples). Intensity of red color denotes strong correlation (Pearson). A minimum of two replicates were performed for each group of cell numbers.

3.3.2 ATAC-seq reveals changes in chromatin accessibility after NT of somatic cells to *Xenopus* oocytes

We then examined chromatin accessibility dynamics during transcriptional reprogramming of somatic cells in oocytes by using the modified ATAC-seq protocol. We took advantage of the direct transcriptional reprogramming system in *Xenopus* oocytes, in which hundreds of mouse somatic nuclei transplanted into the germinal vesicle of a *Xenopus* oocyte undergo extensive transcriptional reprogramming toward an oocyte-like state within 2 days, without cell divisions and DNA replication (Figure 3-3A) (Jullien et al., 2014). Therefore, any changes in chromatin accessibility observed in our reprogramming system are accomplished in a replication-independent manner. We performed ATAC-seq on donor mouse embryonic fibroblasts (MEFs) and on MEFs 48 hr after NT. Two biologically independent NT experiments were performed for the subsequent analyses (10 NT oocytes, equivalent to 3,000 cells, were pooled in each experiment) (Figure 3-3A).

Our ATAC-seq reads of donor MEFs resembled reads from published DNase-seq of mouse 3T3 embryonic fibroblast cells (GSM1003831, Spearman correlation with our reads = 0.84) and mouse headless embryos at day 11.5 (GSM1014172, Spearman correlation with our reads = 0.85; Figure 3-4). Examination of the genomic distribution of peaks showed enrichment around TSSs both before and after NT (Figure 3-3B). Open chromatin peaks were also enriched at enhancers for MEFs before NT, but this enrichment was not observed after NT (Figure 3-3C), suggesting the closing of open chromatin at MEF-specific enhancers after NT. The relative abundance of peaks for embryonic stem cell (ESC) enhancers did not increase much after NT, but the closing of enhancers was not observed, unlike MEF enhancers (Figure 3-3B, right two bars). These results suggest that MEFs, which have undergone NT, may use a different set of enhancers from those before NT to regulate gene transcription.

We next investigated genomic regions newly opened after NT. We found that 3,146 peaks appeared only after NT; those peaks present after NT were at least 5,000 bp away from a peak in MEF before NT. Such newly appeared open chromatin regions

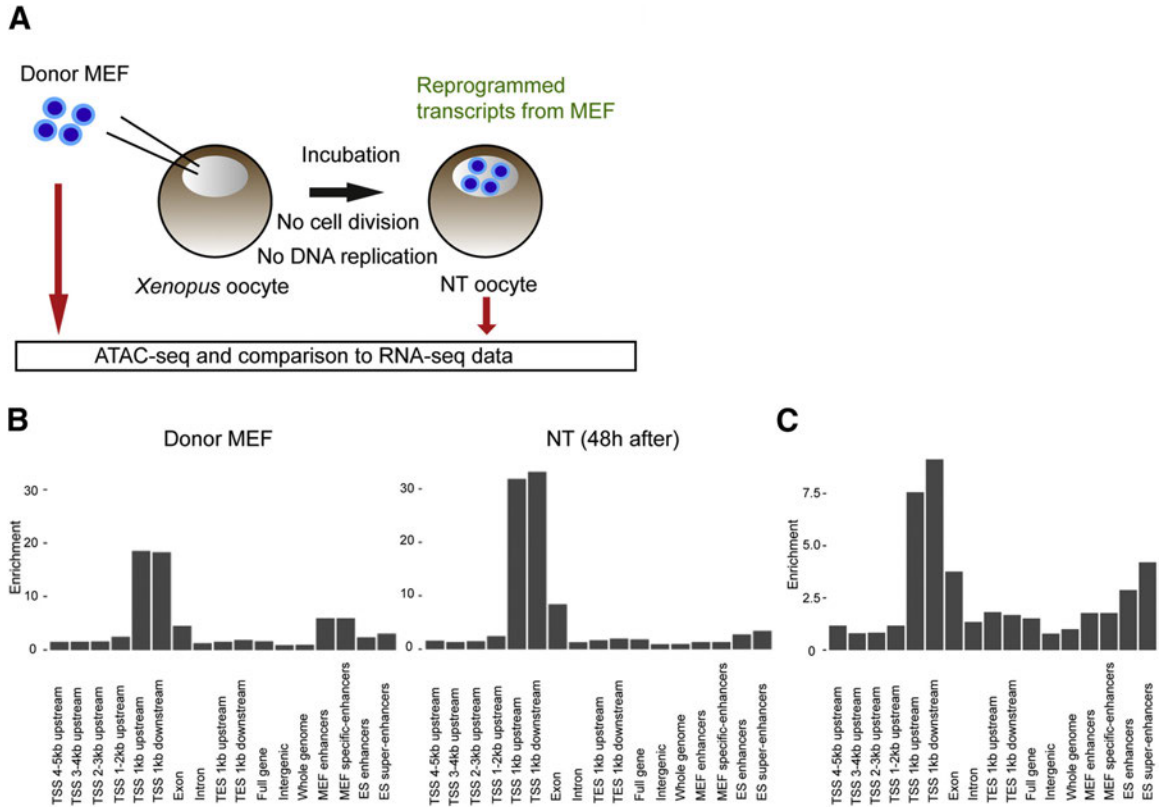


Figure 3-3: Chromatin accessibility changes in oocyte reprogramming. (A) Experimental workflow. (B) Genomic distribution of ATAC-seq peaks representing open chromatin before and after NT. The y axis represents the enrichment of peaks in each type of genomic region relative to the whole genome. (C) Genomic distribution of newly appeared ATAC-seq peaks after NT.

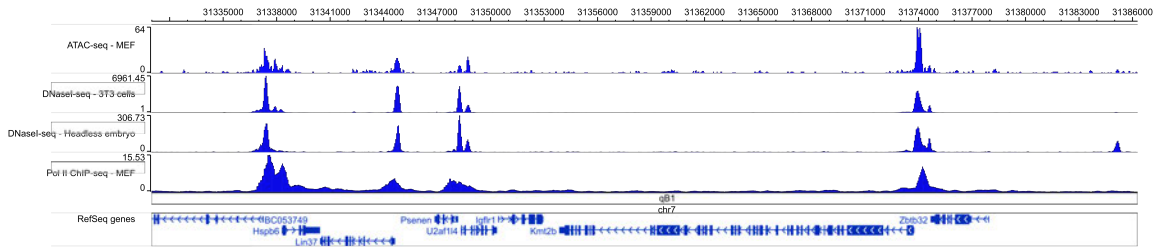


Figure 3-4: Track images of ATAC-seq using 3,000 MEFs, DNase-seq using 3T3 mouse embryonic fibroblasts (GSM1003831) and Day 11.5 headless embryo (GSM1014172), and ChIP-seq analyses for RNA polymerase II using MEFs (GSM918761). The same genomic region is shown for all assays.

were enriched at gene regulatory regions including TSSs and embryonic stem (ES) super-enhancers (Figure 3-3C). The mean distance of these newly opened chromatin regions from a TSS is 61,049 bp, compared with 23,813 bp for all peaks after NT, suggesting that newly open regions are more likely to be distant from promoters. These results suggest that dynamic changes in chromatin accessibility are induced in a replication-independent manner during oocyte-mediated reprogramming.

3.3.3 Open promoters of silent genes are permissive for transcriptional activation in oocytes

Distribution of open chromatin peaks clearly indicates that accessible chromatin sites are mainly located at gene regulatory regions (Figure 3-3BC). We then investigated the relationship between open chromatin at gene regulatory regions and transcriptional activation. We took advantage of our previous RNA-seq dataset in which the transcriptome of donor MEFs before and after NT was revealed [42]. Mapped ATAC-seq reads were plotted around TSSs and, as expected, genes expressed after NT showed more open chromatin than non-expressed genes in these samples (Figure 3-5B). Intriguingly, genes expressed after NT exhibited more open chromatin than genes not expressed after NT at TSSs in the donor cells as well (Figure 3-5A), suggesting that pre-existing open chromatin might affect subsequent transcriptional reprogramming. To further pursue this, we examined genes newly activated after NT. These genes showed more open chromatin at TSSs than non-expressed genes in

MEFs both before and after NT ("only after NT" gene group versus "neither before nor after NT" gene group; Figure 3-5CD).

We visualized the chromatin accessibility of promoters in MEFs for each gene newly activated after NT (Figure 3-6, upper panels). This shows that a large proportion of genes activated after NT contained open chromatin already in donor cells and were also associated with H3K4me3 and unphosphorylated RNA polymerase II (Figure 3-6, upper panels). This is in contrast to genes that our previous study identified as resistant to transcriptional activation in MEFs during oocyte-mediated reprogramming (Figure 3-6, lower panels; reprogramming-resistant genes are discussed below) [44]. In conclusion, genes marked at their TSSs by open chromatin and other priming-associated factors such as the loading of poised RNA polymerase II [3] and H3K4me3 [92] are more likely to undergo transcriptional activation during reprogramming in oocytes compared to genes whose TSSs are not marked by these factors. These results also imply that pre-existing open chromatin states in donor cells affect transcriptional reprogramming, further extending our previous finding that somatic cell genes abnormally maintain their expression in *Xenopus* NT embryos as somatic memory genes [37].

3.3.4 Open chromatin is often maintained near genes downregulated after NT

Genes downregulated after NT (genes that were only expressed in donor MEFs) had a higher level of open chromatin than non-expressed genes in donor cells, as expected, but strong open chromatin states were unexpectedly maintained in downregulated genes even after NT (Figure 3-5CD, only before NT).

On further inspection, we observed that, while downregulated genes retained open chromatin after NT at the global level, the open regions shifted their locations slightly (Figure 3-7). In fact, some motifs were more enriched in ATAC-seq peaks detected in MEFs near these downregulated genes than in peaks detected in MEFs after NT around the same genes, and vice versa (Figure 3-8). As expected, the TFs with motifs

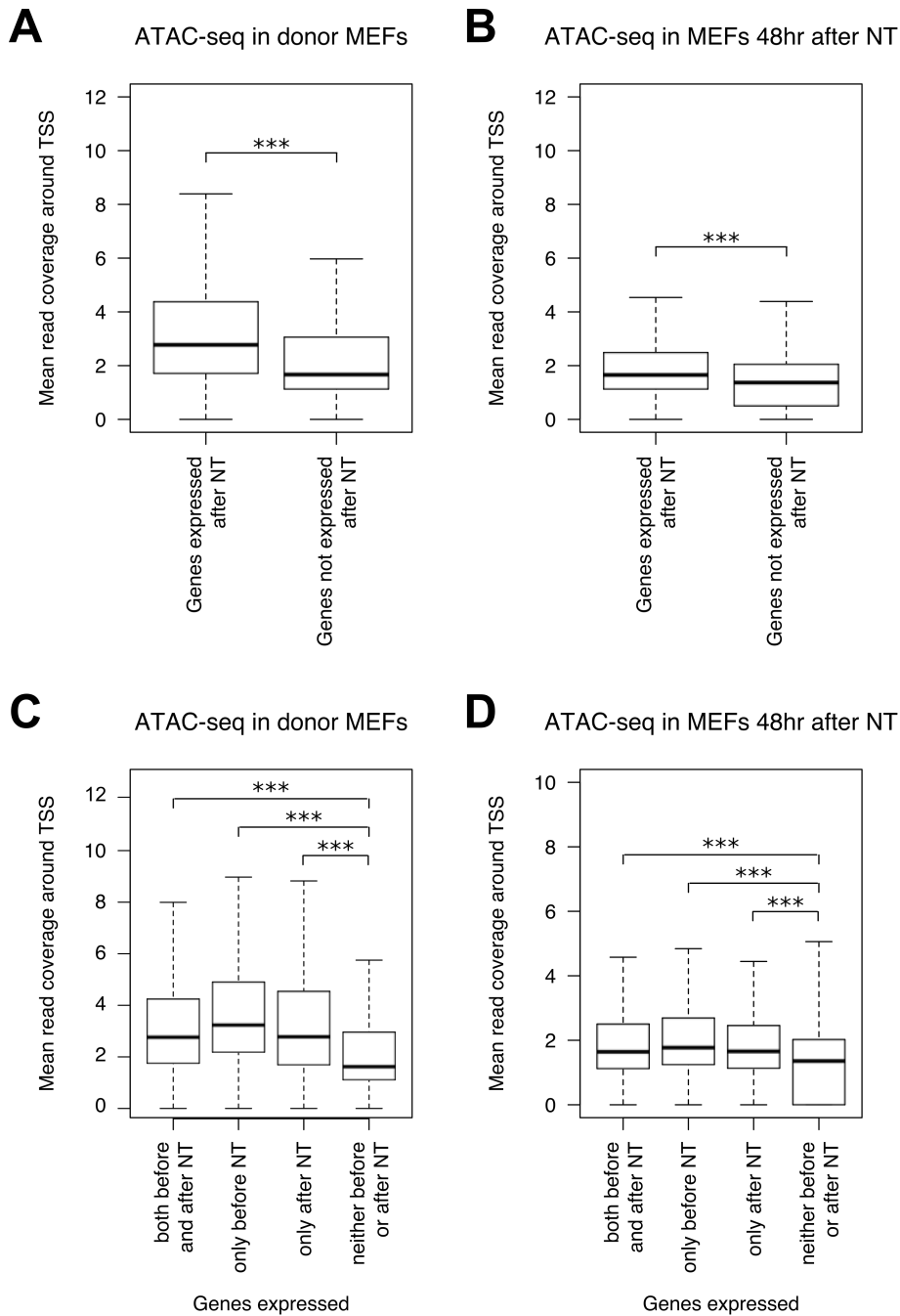


Figure 3-5: ATAC-seq reads in MEFs before and after NT were compared around TSSs in groups of genes. The y axis represents the mean read coverage in a 1-kb window centered on the TSS. *** indicates p-value < 1E-6 by the Mann-Whitney U-test. A) Donor MEFs. Genes were divided into two categories: expressed in NT oocytes and not expressed in NT oocytes. B) MEFs 48 hours after NT. Genes were divided into two categories as in A. C) Donor MEFs. Genes were divided into four categories: genes expressed before and after NT, those expressed only before NT, those expressed only after NT, and those expressed at neither time point. D) MEFs 48 hours after NT. Genes were divided into four categories as in C.

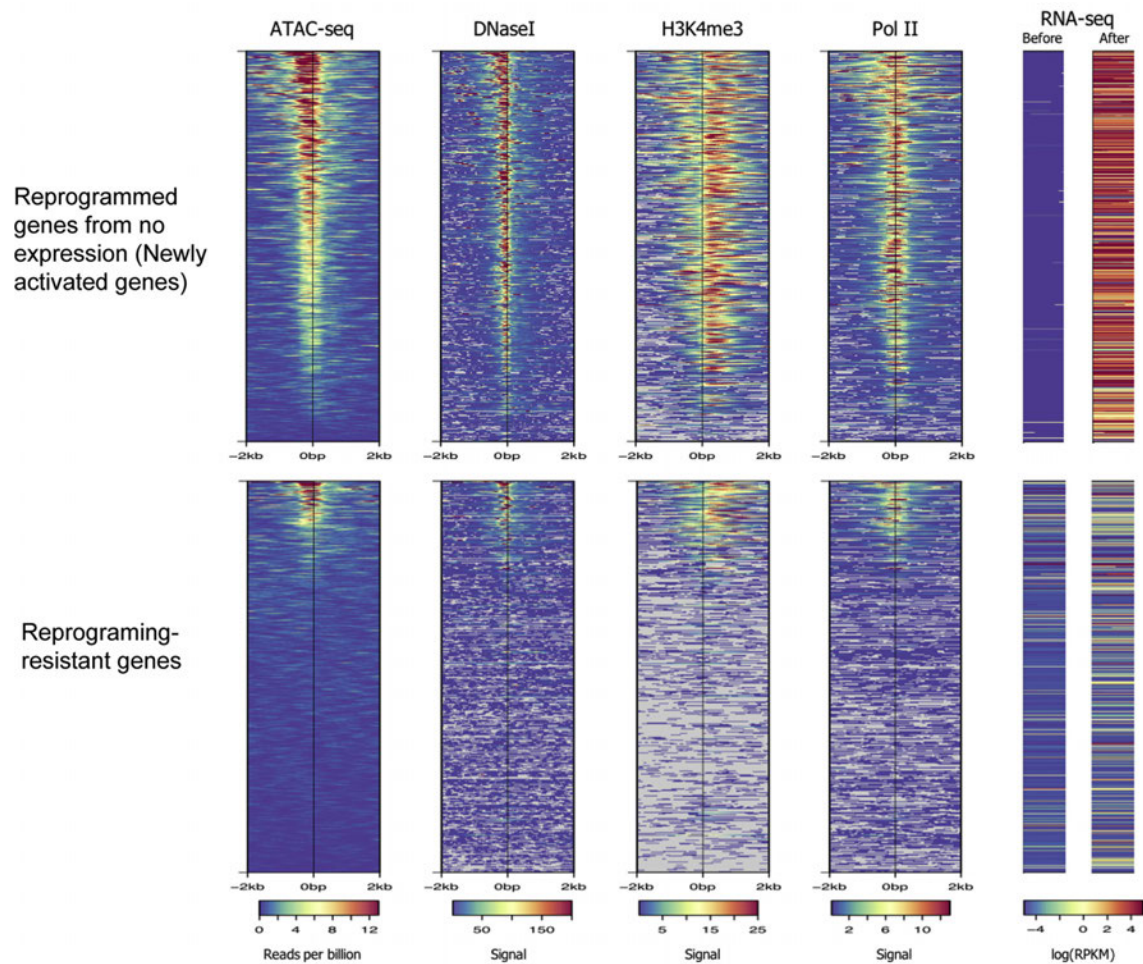


Figure 3-6: Representation of signal associated with open chromatin (ATAC-seq, DNase-seq, H3K4me3 ChIP-seq, and Pol II ChIP-seq) at TSSs of MEF genes reprogrammed in *Xenopus* oocytes. RNAseq results (Jullien et al., 2014, 2017) are shown at the right panel. Accession numbers for the DNase-seq, H3K4me3, and Pol II data are GSM1014172, GSM769029, and GSM918761, respectively. *** $p < 1E-6$ by the Mann-Whitney U test.

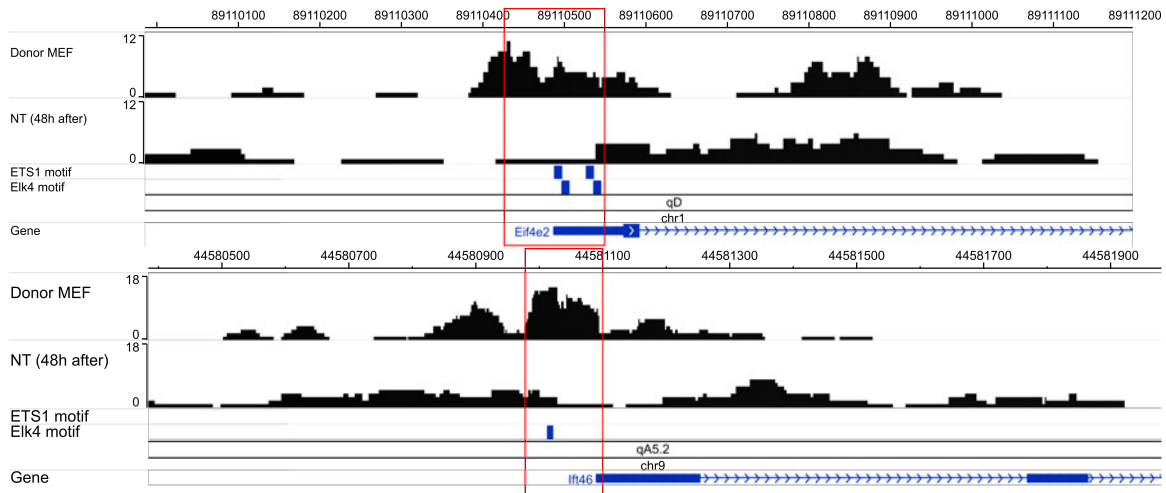


Figure 3-7: ATAC-seq reads for MEFs and MEFs 48 hours after NT around the genes *Eif4e2* and *Ift46*, which were identified as genes whose expression is downregulated after NT. Locations of ETS1 and ELK4 motifs are also shown, which are motifs associated with MEF peaks near downregulated genes. Red parts indicate open chromatin regions that disappeared after NT.

more enriched in MEF peaks tended to be more expressed in MEFs and not expressed in *Xenopus* oocytes, while the TFs with motifs more enriched in post-NT peaks tended to be more expressed in MEFs after NT and/or in *Xenopus* oocytes (Figures 3-8). These results suggest that oocytes utilize different TFs from somatic cells to maintain the open chromatin near genes expressed only before NT, in accordance with our previous finding that the removal of somatic transcriptional machinery and the loading of the oocyte counterpart is observed after NT to oocytes[42].

Finally, this suggests that open chromatin is not sufficient to predict transcription after nuclear transfer, since these downregulated genes are open both before and after NT but are not expressed after NT.

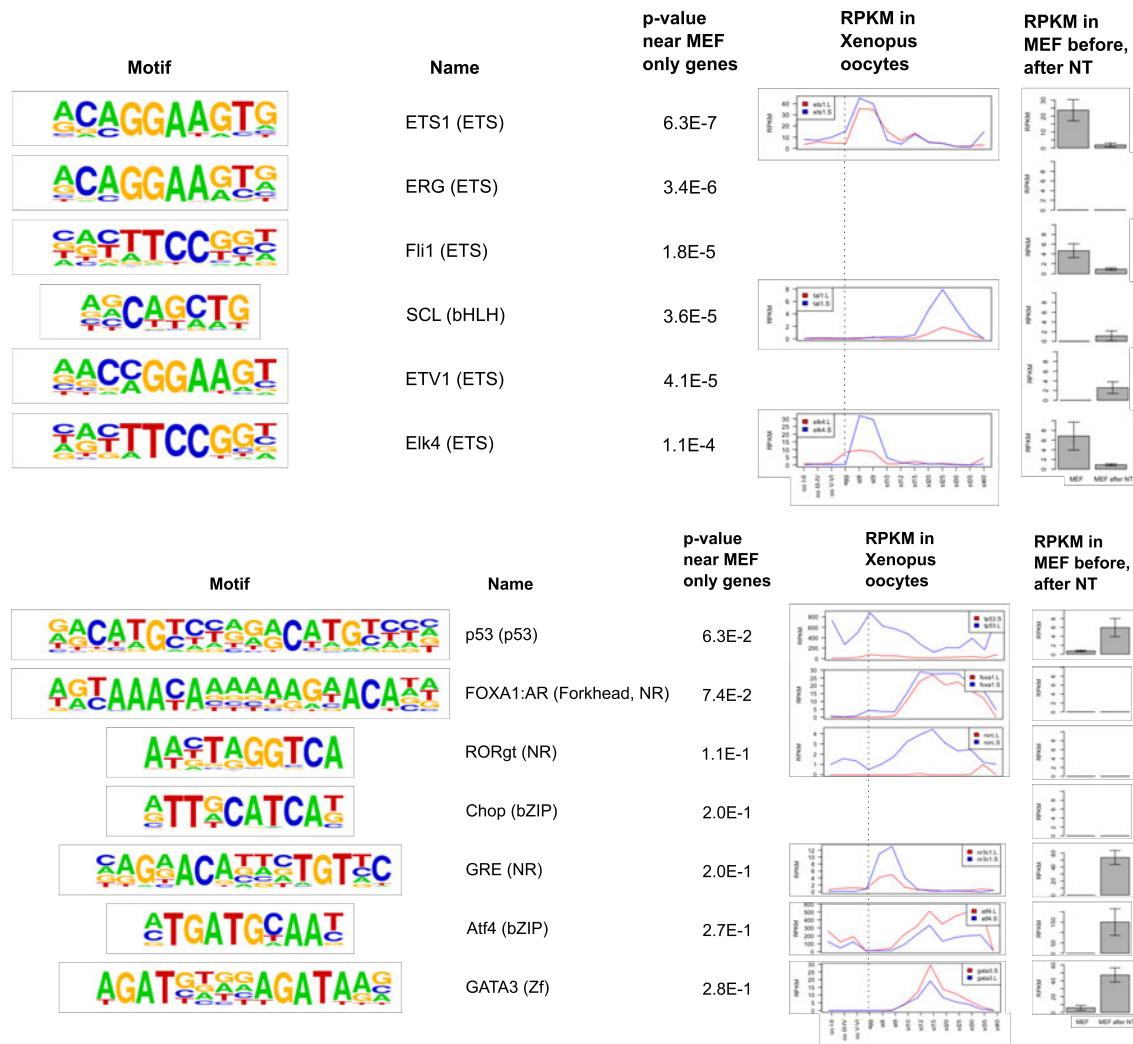


Figure 3-8: Transcription factor motif analysis of genes downregulated after NT. Top) Sequence motifs enriched in open chromatin regions in donor MEFs near genes downregulated after NT (overlapping or within 1 kb, relative to a background of NT peaks near downregulated genes). P-values are calculated by comparing the number of foreground and background sequences with each motif using the binomial test. Expression levels of the identified transcription factors during *Xenopus* oogenesis and embryonic development are shown as line graphs (Session et al., 2016). The dotted black line separates oogenesis and embryogenesis. In case the corresponding MEF genes are not found in *Xenopus*, no graphs are shown (n.a.). RPKM values of the identified TFs in MEFs before and after NT to *Xenopus* oocytes are also shown as bar graphs, with error bars representing \pm SEM (Jullien et al., 2014). Bottom) Sequence motifs enriched in open chromatin regions in MEFs after NT near genes downregulated after NT (overlapping or within 1 kb, relative to a background of MEF peaks near downregulated genes).

3.3.5 Oocyte-Mediated Reprogramming Involves Opening of Closed Chromatin

Although many genes are transcribed from pre-existing open TSSs (Figure 3-6), some genes acquire open chromatin from the closed state during oocyte reprogramming. 3,146 ATACseq peaks were detected only in NT samples, but not in donor MEFs, and 497 genes without an ATAC-seq peak in MEFs gained a peak after NT, such as *Utf1* and *Hoxc8*, both of which have been shown to be reprogrammed after NT [69] (Figure 3-9). We also found 1,245 ATAC-seq peaks after NT that do not overlap genes and are located at least 5 kb from a TSS or a peak before NT. These newly formed regions of open chromatin may serve as enhancers that drive transcriptional reprogramming in this system, analogous to the enhancers that open during cellular differentiation and iPS reprogramming [94, 97]. Newly opened regions were associated with genes related to fatty acid biosynthesis, regulation of stem cell maintenance, and protein targeting to membrane, as revealed by the Genomic Regions Enrichment of Annotations Tool (GREAT) [60] (Figure 3-10). Furthermore, several motifs for binding of TFs were enriched in ATAC-seq peaks after NT relative to before NT, including GATA binding protein 3 (GATA3), retinoic acid receptor gamma (RAR), and RE1-silencing transcription factor (REST) (Figure 3-11). Interestingly, REST has been shown to play a key role in NT-mediated nuclear reprogramming in porcine oocytes [49], supporting the validity of our approach.

3.3.6 Transcription Factors in Oocytes Are Involved in Transcriptional Reprogramming

Because chromatin sites newly opened after NT are enriched with specific TF motifs (Figure 3-11), we sought to test roles of TFs in oocytes for transcriptional reprogramming. Among the TFs identified, we focused on RAR because it showed one of the most significant hits and because RAR functions in *Xenopus* oocytes [64]. We first selected genes to test for further analyses. *Pou5f1* (known as *Oct4*) and *Utf1* have been shown to be regulated by RARs [24, 72] and to be newly activated after NT to

A

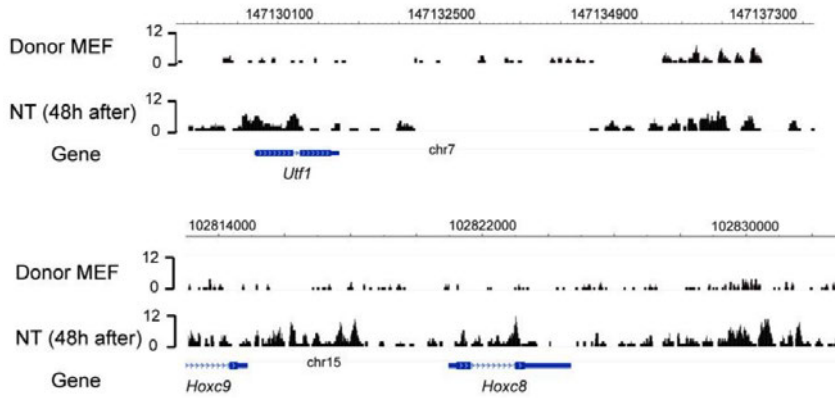


Figure 3-9: Newly produced open chromatin during reprogramming in *Xenopus* oocytes. ATAC-seq reads for donor MEFs and MEFs 48 hours after NT.

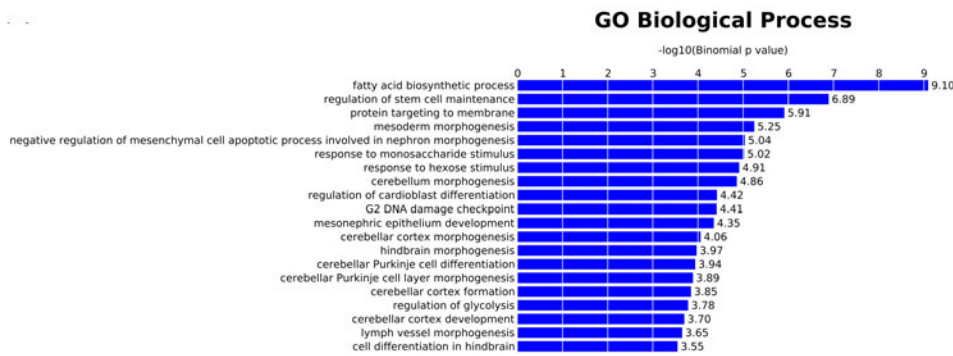


Figure 3-10: Gene categories associated with newly opened chromatin peaks after NT are identified by the Genomic Regions Enrichment of Annotations Tool (GREAT). Shown are the top 20 Gene Ontology Biological Process terms that are significant by both the region-based binomial test and the gene-based hypergeometric test and have at least a 2-fold region-based enrichment, ranked by region-based binomial p-value.

Rank	Motif	Name	p-value	q-value
1	AGATGTTTAAACAGATAAAG	GATA3	1.00E-05	0.0013
2	AGGTC AAGGICA	RARg	1.00E-05	0.0013
3	GGAGCTGTCCAIGGTGCTGA	REST-NRSF	1.00E-02	0.587

Figure 3-11: Sequence motifs enriched in open chromatin regions after NT, relative to open chromatin regions in donor MEFs. Transcription factors that recognize the indicated motifs are also shown. P-values are calculated by comparing the number of foreground and background sequences with each motif using the binomial test.

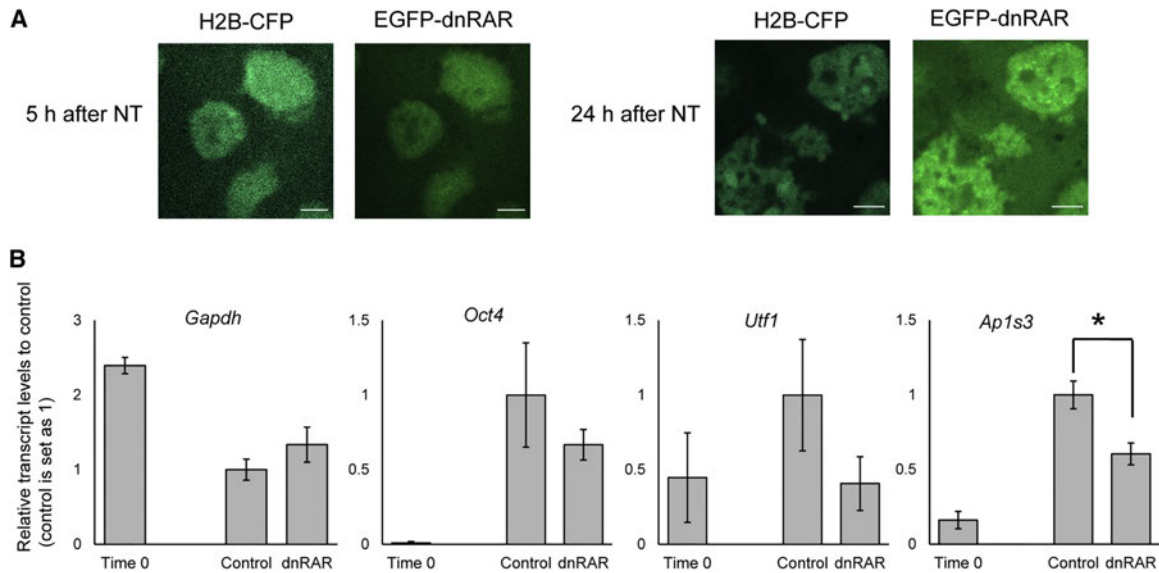


Figure 3-12: RAR Influences Transcriptional Reprogramming in Oocytes (A) NT oocytes overexpressed with EGFP-dnRAR and histone H2B-CFP were subjected to confocal microscopy. EGFP-dnRAR was accumulated in the injected nuclei. Scale bars indicate 5 mm. (B) Expression of RA-regulated genes was downregulated after overexpression of EGFP-dnRAR in NT oocytes. n = 3 (time 0) or 4 (control and dnRAR). Time 0 represents NT oocytes just after NT. Error bars represent \pm SEM. *p < 0.05 by the Student's t test.

oocytes [69]. We also examined *Ap1s3*, which contains Retinoic Acid Response Element (RARE) at the TSS and is activated after NT. A RAR-*alpha* dominant-negative form (dnRAR) [95] was overexpressed in NT oocytes. dnRAR was localized in transplanted MEF nuclei several hours after NT (Figure 3-12A). The overexpression of dnRAR in NT oocytes inhibited transcription from *Oct4*, *Utf1*, and *Ap1s3* (1.5- to 2.5-fold decrease), whereas expression of *Gapdh*, a housekeeping gene, was not affected (Figure 3-12B). These results suggest that newly activated genes are indeed regulated by RAR. Therefore, TFs in oocytes impact transcriptional reprogramming, possibly through opening of inaccessible sites in somatic chromatin.

3.3.7 Closed Chromatin Is Associated with Reprogramming-Resistant Genes

We next asked whether an opening of somatic chromatin is efficiently carried out or not after NT to oocytes. Our previous study identified genes resistant to oocyte-mediated reprogramming in MEFs [44]. These genes are not activated after nuclear transplantation of MEFs to oocytes but are transcribed after transplantation of mouse ESCs [44]. This suggests that although transcriptional activators for the MEF reprogramming-resistant genes are present in oocytes, some feature of MEFs is preventing successful activation. Interestingly, both resistant and activated genes were transcriptionally silent in donor MEFs, but the TSSs at resistant genes in MEFs were clearly more closed compared with successfully activated genes (Figure 3-13). The opening of chromatin is a central factor for successful transcriptional reprogramming, but is inefficient as is evident from the presence of the MEF reprogramming-resistant genes, suggesting that the opening of closed chromatin is a barrier for reprogramming.

3.3.8 Chromatin Opening Facilitates Transcriptional Reprogramming

We finally tested whether accessible chromatin states enable efficient transcriptional reprogramming. We have previously reported that enhancing nuclear actin polymerization results in efficient transcriptional reprogramming [67], and that nuclear actin polymerization plays a role in chromatin decondensation [7]. Toca1/Fnbp11, a factor that enhances actin polymerization and transcriptional reprogramming [67], was overexpressed in oocytes. MEFs transplanted to the nucleus of these Toca1/Fnbp11-overexpressed NT oocytes were subjected to ATAC-seq. More open chromatin peaks were found in MEFs transplanted to Toca1/Fnbp11-overexpressed NT oocytes, compared with control NT oocytes (2.6-fold increase), and 2.1-fold more genes were associated with ATAC-seq peaks when Toca1/Fnbp11 was overexpressed. Toca1/Fnbp11 overexpression has been shown to enhance *Oct4* activation [67]. Indeed, more ATAC-seq reads were found at the *Oct4* gene locus in MEFs transplanted to Toca1/Fnbp11-

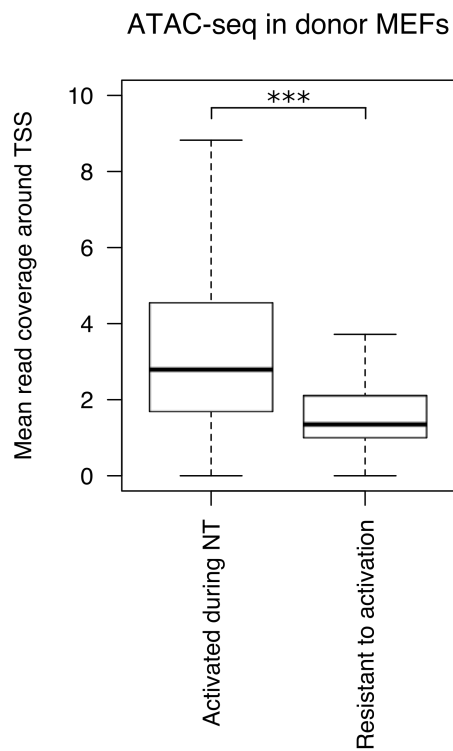


Figure 3-13: ATAC-seq reads around TSSs in donor MEFs were compared among different gene categories; genes expressed after NT and the MEF reprogramming-resistant genes[44]. The y-axis represents the mean read coverage in a 1 kb window centered on the TSS. *** indicates p-value < 1E-6 by the Mann-Whitney U-test.

overexpressed NT oocytes than in MEFs transplanted to control NT oocytes (Figure 3-14A), suggesting that enhanced transcriptional reprogramming is associated with the accessible chromatin state. On the other hand, genes whose transcription is not affected by Toca1/Fnbp11 overexpression [67] also showed more ATAC-seq reads (Figure 3-15), supporting our contention that open chromatin at TSSs is not sufficient to determine transcriptional activities. The genomic distribution of peaks found in MEFs transplanted to Toca1/Fnbp11-overexpressed NT oocytes showed enrichment around TSSs (Figure 3-16), like MEFs transplanted to normal NT oocytes (Figure 3-3B). Together, accelerated nuclear actin polymerization by Toca1/Fnbp11 overexpression seems to induce accessible chromatin globally, rather than locally, at specific loci. We next treated NT oocytes with 50 nM Trichostatin A (TSA), a histone deacetylase inhibitor, which can induce the accessible chromatin state in somatic nuclei [17, 58, 68]. Transplanting MEFs to oocytes incubated with TSA increased *Oct4* and *Utf1* expression by 19.6- and 30.2-fold, respectively (Figure 3-14B). These results suggest that somatic chromatin is not fully opened to support transcriptional reprogramming, even after NT to oocytes, and reprogramming can be enhanced by modulating chromatin accessibility.

3.4 Conclusions

In summary, our study reveals genome-wide changes in chromatin accessibility in a DNA-replication independent manner after NT of somatic cells to *Xenopus* oocytes. Transcriptional activation from the silenced genes in donor somatic cells can be accomplished from both pre-existing open chromatin and closed chromatin. However, we observed that silent genes that are primed with pre-existing open chromatin are more likely to be transcriptionally activated during oocyte-mediated reprogramming than silent genes without open chromatin. Presumably, having open chromatin, H3K4me3, and poised RNA polymerase II (Pol II) at these genes helps ensure that these genes can be promptly and efficiently activated, compared with genes with closed chromatin. In contrast, the opening of closed chromatin seems to be an inefficient process. In-

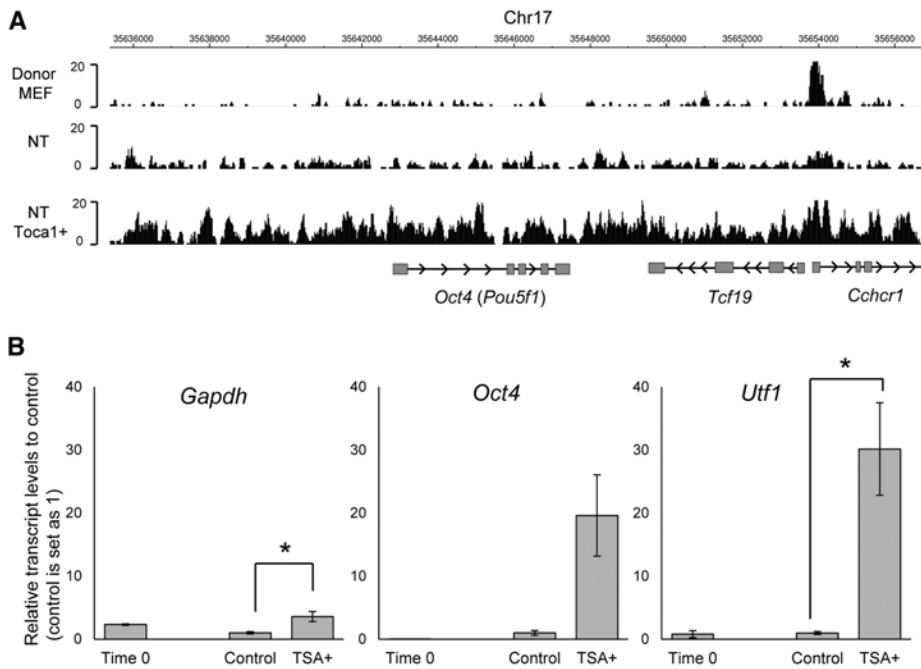


Figure 3-14: Efficient Transcriptional Reprogramming Is Associated with Increased Chromatin Accessibility (A) ATAC-seq reads around the *Oct4* gene locus for donor MEFs, MEFs after NT to control oocytes, and MEFs after NT to oocytes overexpressed with Toca1/Fnbp1l. (B) TSA treatment enhances transcriptional reprogramming. MEFs after NT to oocytes incubated with or without 50 nM TSA for 24 hr were used for qRT-PCR. Time 0 represents cells just after NT. $n = 3$ (time 0) or 4 (control and dnRAR). Error bars represent \pm SEM. * $p < 0.05$ by the Student's t test.

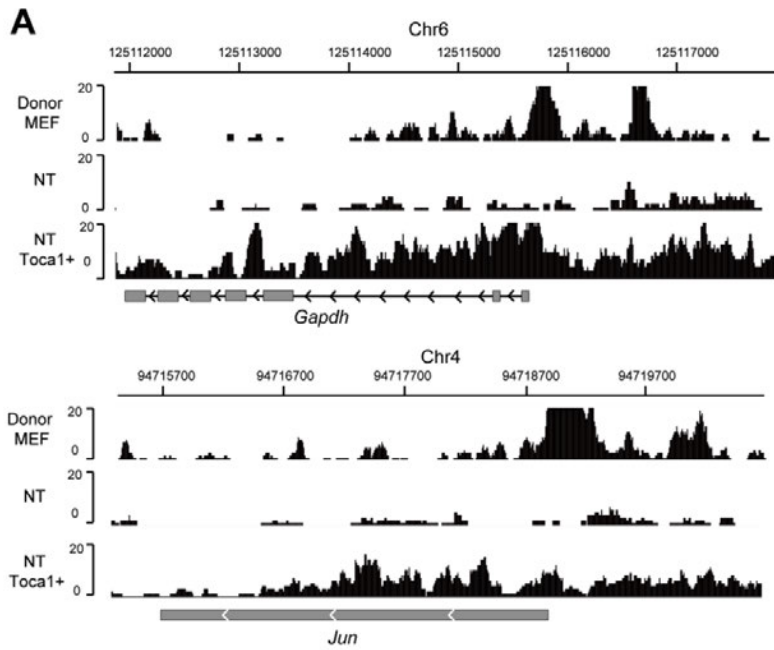


Figure 3-15: ATAC-seq reads around the *Gapdh* and *Jun* gene loci for donor MEFs, MEFs after NT to control oocytes, and MEFs after NT to oocytes overexpressed with Toca1/Fnbp1l.

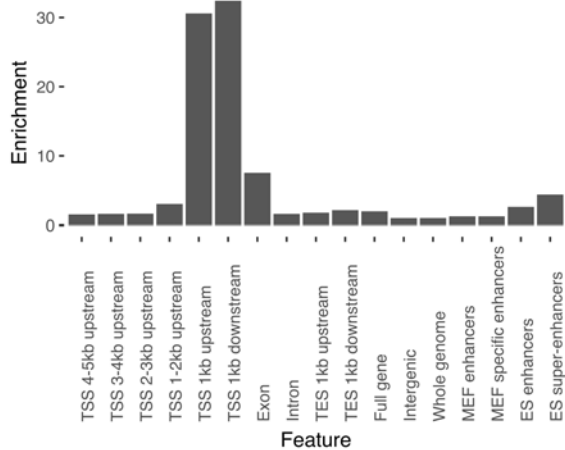


Figure 3-16: The genomic distribution of ATAC-seq peaks detected in MEFs after NT to Toca1/Fnbp1l-overexpressed oocytes. The y-axis represents the enrichment of peaks in each type of genomic region relative to the whole genome.

deed, we found that a closed chromatin configuration at TSSs is a barrier to some genes successfully being activated during reprogramming. Chromatin opening related to transcriptional reprogramming is likely mediated by transcription factors (TFs). Although oocytes possess many TFs and other machinery to relax chromatin [43, 42], experimental interventions can still influence chromatin opening and transcriptional reprogramming of somatic cells. We showed that inhibiting the transcription factor RAR in oocytes inhibits NT-mediated transcriptional reprogramming, and that over-expressing Toca1/Fnbp11 in oocytes can increase chromatin accessibility in donor cells after NT. These observations suggest that somatic nuclei acquire stable chromatin signatures in the course of differentiation that shape the transcriptional reprogramming induced by nuclear transfer. Increased chromatin accessibility contributes to the destabilization of somatic chromatin features and hence successful reprogramming.

Chapter 4

Multi-omic mediation analysis across multiple disease traits and cell types

4.1 Introduction

Genome-wide association studies (GWAS) have identified many genetic loci that are associated with disease, but it remains difficult to translate those associations into therapeutic targets[91, 96]. The majority of risk loci that have been identified by GWAS are non-coding[91, 96], and linkage disequilibrium (LD), the correlation structure between variants in the genome, makes it difficult to identify truly causal variants[83]. Furthermore, the disease mechanism is not always obvious, even given the causal variant(s)[83, 36].

Mediation analysis provides a powerful approach for linking risk variants to disease mechanisms, by identifying intermediate molecular phenotypes, such as gene expression, that also vary significantly with disease-associated variants. Further, this approach enables the interpretation of non-coding, regulatory variants[26]. Several computational methods have been proposed to integrate GWAS data with data from studies measuring how the expression of genes vary with genotype, also known as expression quantitative trait locus or eQTL studies. These include Mendelian randomization[23], PrediXcan[26], TWAS[29], coloc[27], and CaMMEL[74], among others. Uniquely, CaMMEL is able to account for linkage disequilibrium between

variants, allows for the possibility that the GWAS effect of a variant acts separately from the QTL features under consideration, and considers multiple genes in a locus jointly in order to select a sparse set of causal intermediate phenotypes that can explain the GWAS signal in a region[74].

Most previous studies integrating GWAS with QTLs have focused on expression QTLs, though a few studies have also examined splicing QTLs[57, 89, 55] and DNA methylation QTLs[30]. The gene-trait associations implicated using these non-expression QTLs are often corroborated by eQTLs, but in some cases they suggest mechanisms that are not detected in eQTL analysis[30]. Therefore, investigating additional types of intermediate phenotypes may reveal new insights about disease etiology.

Recent technological developments and dramatic reductions in the cost of sequencing have seen the number of molecular phenotypes that can be systematically profiled dramatically increase. However, current analyses typically focus on one molecular phenotype at a time, making it difficult to understand where genetic variants primarily act, and whether the same genetic variants affect multiple intermediate molecular phenotypes, such as gene expression and enhancer activity.

In this work, we discover a wide range of potential disease mechanisms by integrating GWAS for multiple traits (Crohn’s disease, ulcerative colitis, coronary artery disease, Alzheimer’s disease, and systemic lupus erythematosus) with QTL studies in immune cell types for several transcriptional (gene expression, splicing) and epigenetic mechanisms (DNA methylation, histone marks H3K27ac and H4K4me1) using CaMMEL. We introduce a method to summarize the impact of genetic variants across multiple molecular phenotypes for each gene, and explore the relationships between these molecular mechanisms in the context of disease.

4.2 Materials and Methods

4.2.1 Datasets and processing

Linkage disequilibrium blocks

To define genomic regions that may be linked to a given causal variant, we segment the genome into blocks that are approximately independent genetically (i.e. minimal linkage disequilibrium between blocks). We used the LD block boundaries calculated by Berisa et al. for European populations for our downstream QTL and mediation analyses[10].

QTL data

We analyzed 15 datasets from Chen et al., including genome-wide profiles of five molecular traits (gene expression, splicing, H3K27ac, H3K4me1, and DNA methylation) in three immune cell populations (monocytes, neutrophils, and naive CD4+ T-cells)[20]. We used the percent spliced in (PSI) metric of alternative splicing and the M-value metric of DNA methylation. These data are publically available at: ftp://ftp.ebi.ac.uk/pub/databases/blueprint/blueprint_Epivar/Pheno_Matrix/. Genotypes for the donors corresponding to these molecular measurements are available through the European Genome-Phenome Archive (accession EGAD00001002663). To map cis-QTLs, we first applied PEER with 10 factors to the molecular trait data to estimate potential confounders[85]. We then used the matrixQTL package to estimate the effect size and p-value of variant-QTL associations, fitting the additive linear model with the PEER factors as covariates[79]. We considered features within 1 Mb of the LD block corresponding to each genetic variant. For the downstream mediation analysis, we filtered out molecular traits (e.g. genes, peaks) with no meaningful QTL associations ($p < 1e-3$).

GWAS data

The GWAS that we analyzed are listed in Table B.1. We converted any GWAS statistics that were originally mapped to the hg38 genome assembly to hg19 with the liftOver tool[47].

Estimation of linkage disequilibrium matrix

To create a reference genotype matrix, we combined genotypes from phase 3 of the 1000 Genomes project [61] and from the UK10K project [93]. We excluded rare variants (minor allele frequency $MAF < 0.05$). Within each LD block, we estimate the truncated-SVD regularized linkage disequilibrium (LD) matrix as described by Shi et al[81]. Briefly, we perform singular value decomposition on the standardized genotype dosage matrix and drop the SVD components corresponding to small singular values. We use the remaining SVD components to reconstruct the regularized genotype and LD matrices.

4.2.2 Mediation analysis

General algorithm

Our goal is to identify cis-regulatory features (genes, histone peaks, DNase sites, etc.) that may causally mediate a specific GWAS trait. For each GWAS trait, we select the LD blocks containing at least one variant with an association p-value of $1e-4$ or lower for further analysis. We also filter QTL features, only including genes or epigenetic features with at least one QTL association p-value of $1e-3$ or lower. Then, we perform causal inference on each LD block (as described below) for each GWAS and QTL study combination. This estimates the mediation effect size and posterior probability of mediation for each feature-trait pair.

Causal inference

We performed causal mediation analysis using the CaMMEL framework as previously described[74]. The method is summarized briefly below.

Before we can fit the mediation effects, we need to estimate unmediated effects. These are the components of GWAS effect that act independently of the QTL features under consideration (e.g. through unmeasured QTLs or through pleiotropic effects). To do this, we first learn genetic components C that explain how the genotype matrix X relates to the effects of the variants, by solving the sparse factorization problem $E[z] = n^{-\frac{1}{2}}X^T C W$, where z is the matrix combining z-scores from both the QTL study and the GWAS, X is the reference genotype matrix, and n is the number of individuals in represented in X . Then we exclude components (columns of C) that are associated with both GWAS and QTL effect, i.e. c_l with non-zero weights w_{kl} (posterior inclusion probability of $w_{kl} > 0.5$) for both a GWAS trait k_1 and a QTL trait k_2 . The remaining components are those that may contribute to unmediated effects.

We can now fit a summary-based multivariate regression analysis of the GWAS z-scores on the QTL z-scores and the unmediated components:

$$z_{GWAS} \sim \mathcal{N}\left(\sum_k \beta_k z_k^{QTL} + \sum_l \gamma_l z_l^{unmed}, R\right)$$

z_l^{unmed} can be derived from c_l as follows: $z_l^{unmed} = n^{-\frac{1}{2}}X^T c_l$. Mediation effect sizes β and unmediated effect sizes γ are estimated through Bayesian variable selection, with a spike-slab prior[65, 18]. Parameters α , β , and γ are estimated by posterior mean of variational Bayes inference. The significance of the β parameters is quantified by the Bayesian posterior inclusion probability, $\Pr(\beta \neq 0|data)$.

This method is implemented using the zQTL software package, version 1.3.5 (<https://github.com/YPARK/zqtl/>).

Estimation of heritability explained

For a given GWAS trait, we estimate the local total SNP heritability and the local SNP heritability explained through mediation following the method of Shi et al[81]. Specifically, given the reference LD matrix R , the GWAS effect sizes θ , the QTL effect sizes α , and the mediation effect sizes β , the local total SNP heritability is

$\theta^T R \theta$ and the heritability explained through mediation is $(\alpha\beta)^T R(\alpha\beta)$. We focused on LD blocks with at least one genome-wide significant GWAS SNP for the heritability explained analyses.

Scaled effect size

For a given QTL feature k , we can estimate $z_k^{med} = \beta_k z_k^{QTL}$ for each of the variants associated with k . We define the scaled effect size $B_k = \max(\text{abs}(z_k^{med}))$.

Associating genes with epigenetic features

Epigenetic features (DNAm, H3K4me1, and H3K27ac) are associated with the protein coding gene whose transcription start site (TSS) is closest to the feature. If a histone peak overlaps the TSS of multiple genes, it is associated with all of these genes. Gene annotations (protein-coding status, TSS locations) are derived from the GENCODE v15 human gene annotation set.

Gene level scores

To construct the normalized mean difference matrix D , we define the value of an entry

$$d_{gq} = \frac{\overline{x_{gq}} - \mu_q}{\frac{\sigma_q}{\sqrt{n_{gq}}}}$$

where $\overline{x_{gq}}$ is the mean log odds of any feature of type q associated with gene g , μ_q is the mean log odds of any feature of type q , σ_q is the standard deviation of log odds for features of type q , and n_{gq} is the number of features of type q associated with gene g . We pooled features across all diseases and cell types in order to estimate μ_q and σ_q . Empty entries (i.e. d_{gq} for g with no associated features of type q) take the value 0. d_{gq} summarizes how different the log odds of features assigned to the gene g are from the expectation for features of type q . We also calculated a max score matrix M , where m_{gq} is equal to the maximum log odds of any feature of type q associated with gene g .

Clustering genes by mediation scores

We filtered the gene score matrix D for each GWAS to genes that were significantly linked to disease through at least one molecular phenotype (2 standard deviations above expectation). We concatenated these filtered matrices (such that each row is a GWAS-gene pair and each column is a QTL type), and performed k-means clustering on this matrix, with the Lloyd algorithm. To facilitate the visualization of the data, we performed UMAP on this matrix with the R package uwot.

4.2.3 Functional analysis

Gene set analysis

We used Gene Set Enrichment Analysis (v3.0) to interpret mediation scores for genes with respect to previously defined functionally related sets of genes[87, 70]. We use the GSEAPreranked tool to perform GSEA on with the classic scoring scheme. We tested gene derived from three sources: MSigDB Hallmark gene sets, Gene Ontology Biological Process terms, and KEGG[53, 5, 1, 45]. The pathways were downloaded from MSigDB v6.2 (<http://software.broadinstitute.org/gsea/downloads.jsp>). The MSigDB KEGG pathway set was supplemented with KEGG pathways included in the Homer software package (<http://homer.ucsd.edu/homer/>)[33]. In a given analysis, we exclude gene sets with fewer than 15 genes or more than 500 genes present in the current ranked gene list.

4.3 Results and Discussion

4.3.1 Discovering molecular mediators of disease

We jointly analyze GWAS and QTL studies to find molecular mechanisms in immune cells that may be causally linked to the development of disease. We start with publicly available data from GWAS for multiple traits (Crohn’s disease, ulcerative colitis, coronary artery disease, Alzheimer’s disease, and systemic lupus erythematosus), and

15 QTL studies, measuring the genetic association of 5 different molecular phenotypes (gene expression, splicing, H3K27ac, H3K4me1, and DNA methylation) in 3 immune cell types (monocytes, neutrophils, and naive CD4+ T-cells), and perform CaMMEL for each GWAS and QTL study combination. CaMMEL estimates the likelihood that each tested expression or epigenetic feature casually mediates the disease of interest, as well as the effect size of this mediation (Figures 4-1 – 4-5, Tables B.2 – B.6).

For each of the diseases, we see many predicted mediators, across all molecular phenotypes (Figures 4-1 – 4-5). For Crohn’s disease, the top predicted mediators by both posterior probability and effect size are near the PTGER4 gene on chromosome 5 (Figure 4-1). In our analysis, the role of PTGER4 is supported by gene expression, H3K27ac, H3K4me1, and DNAm mediators, although interestingly, the epigenetic mediators are predicted to be stronger than the expression mediators. For ulcerative colitis, the top predicted mediators correspond to the ERAP1 gene, on chromosome 5, and the HLA-DQB1 gene, on chromosome 6 (Figure 4-2). ERAP1 is associated with DNAm and expression mediators, while HLA-DQB1 is associated with H3K4me1, DNAm, and H3K27ac mediators. For coronary artery disease, the top mediators are less clear (Figure 4-3). The CRT3 and AKAP13 genes, on chromosome 15, are both associated with DNAm mediators with high posterior probability and effect size, and AKAP13 is also associated with a H3K27ac mediator. The gene TMEM106B, on chromosome 7, also has relatively high signal across multiple QTL types (H3K27ac, H3K4me1, and gene expression). For Alzheimer’s disease, several of the top predicted mediators are associated with the BIN1 gene, on chromosome 2 (Figure 4-4). BIN1 is associated with DNAm, H3K4me1, H3K27ac, and gene expression mediators. Finally, for systemic lupus erythematosus (SLE), there is a strong cluster of mediation signal around the HLA locus on chromosome 6, across all QTL types (Figure 4-5).

These top mediation loci generally correspond to genome-wide significant GWAS loci (Figures 4-1 – 4-5). This includes PTGER4 for Crohn’s disease, the HLA locus for ulcerative colitis, CRT3 for coronary artery disease, BIN1 for Alzheimer’s disease, and the HLA locus for SLE. However, some top GWAS loci are not highly ranked in this mediation analysis or do not correspond to predicted mediators in our analysis.

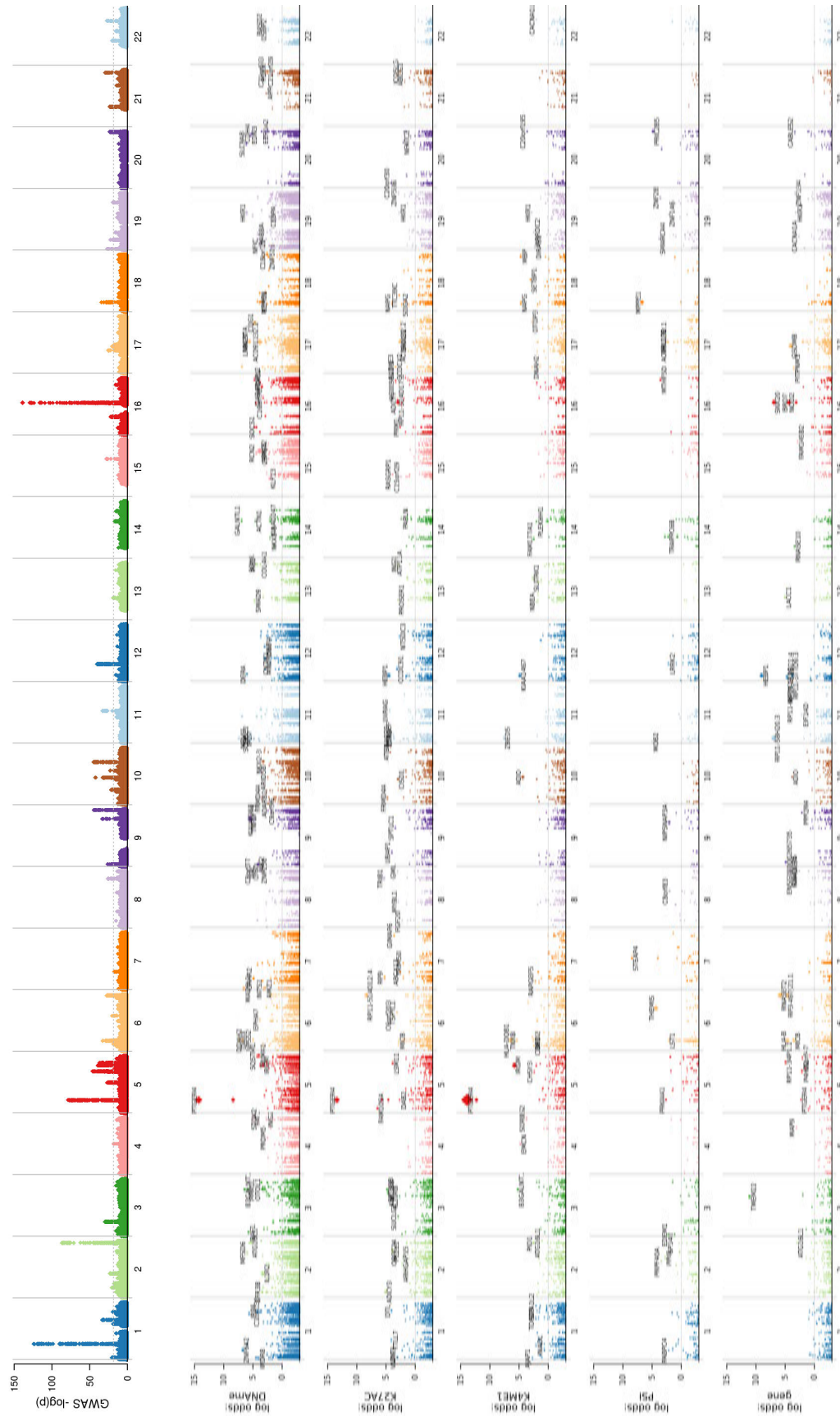


Figure 4-1: Features associated with Crohn’s disease. Each row shows signal for a different data type: GWAS, DNA methylation mediation, H3K27ac mediation, H3K4me1 mediation, splicing mediation, and gene expression mediation. Each point represents one feature (SNP, CpG site, region, splice site, gene). For the mediation signal, the size of each point reflects the effect size, and we display the maximum signal over all of the cell types (monocytes, neutrophils, CD4+ T-cells) for the given feature. Each feature is labeled with the name of the corresponding gene (for expression and splicing features) or the closest gene (for epigenetic features).

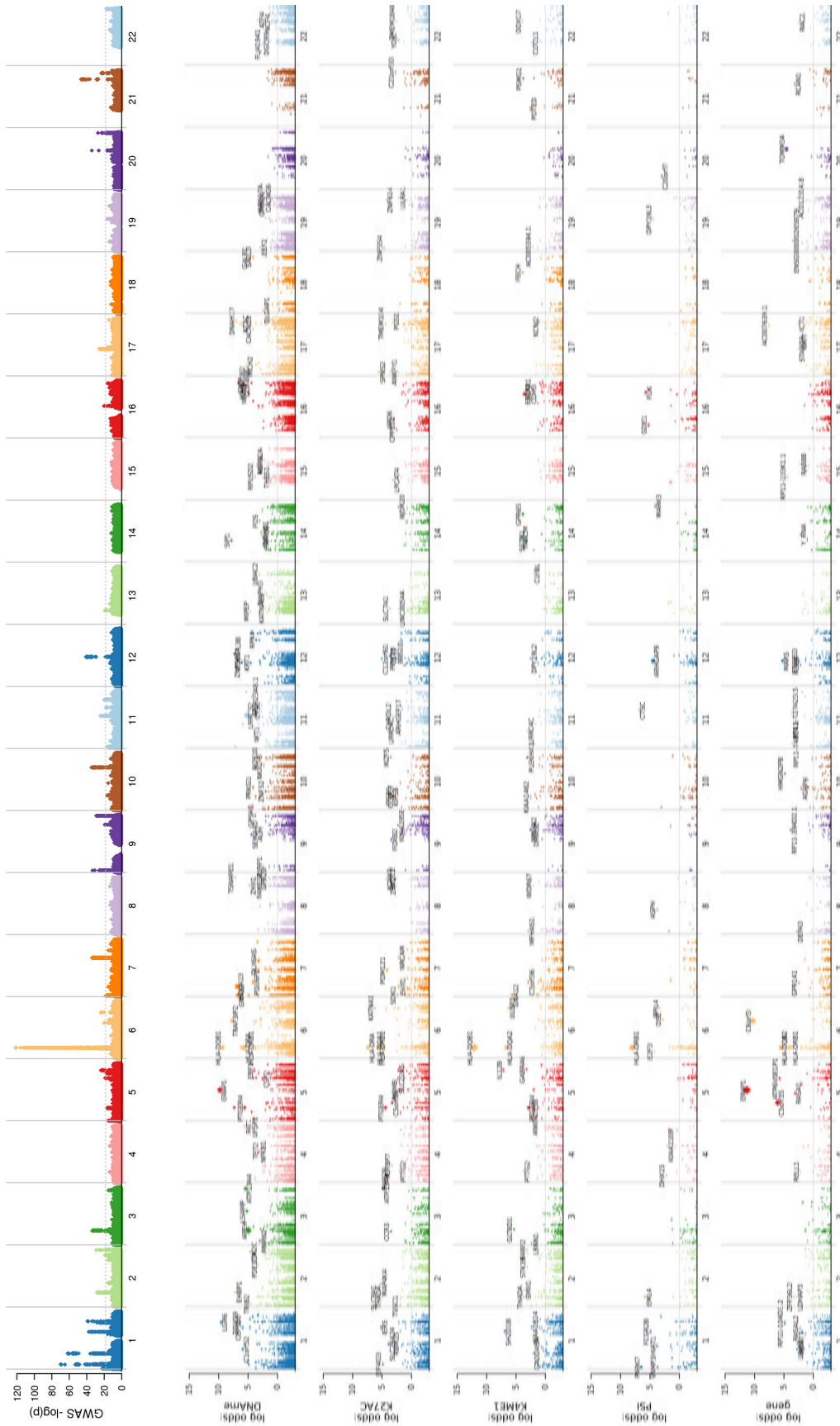


Figure 4-2: Features associated with ulcerative colitis.

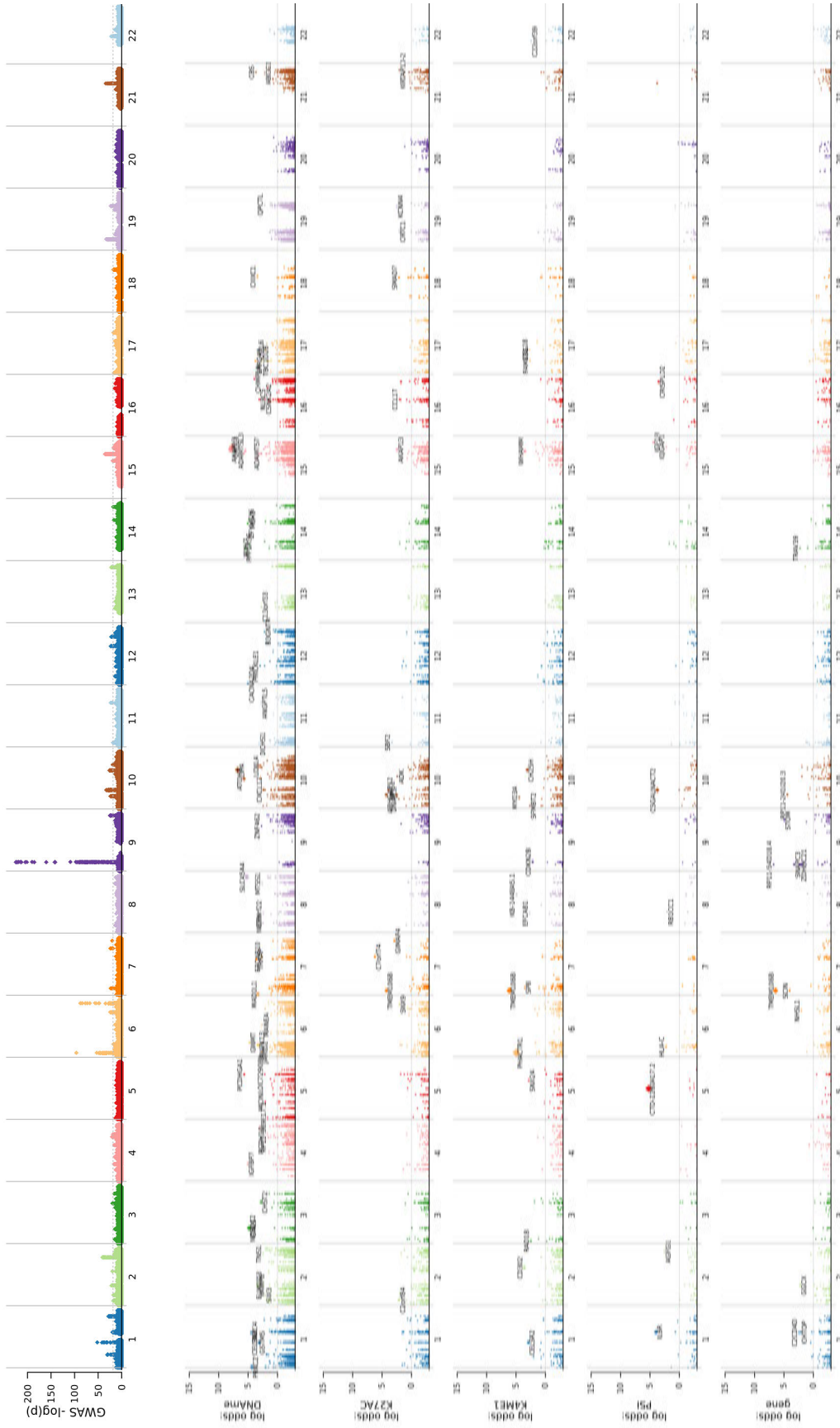


Figure 4-3: Features associated with coronary artery disease.



Figure 4-4: Features associated with Alzheimer's disease.

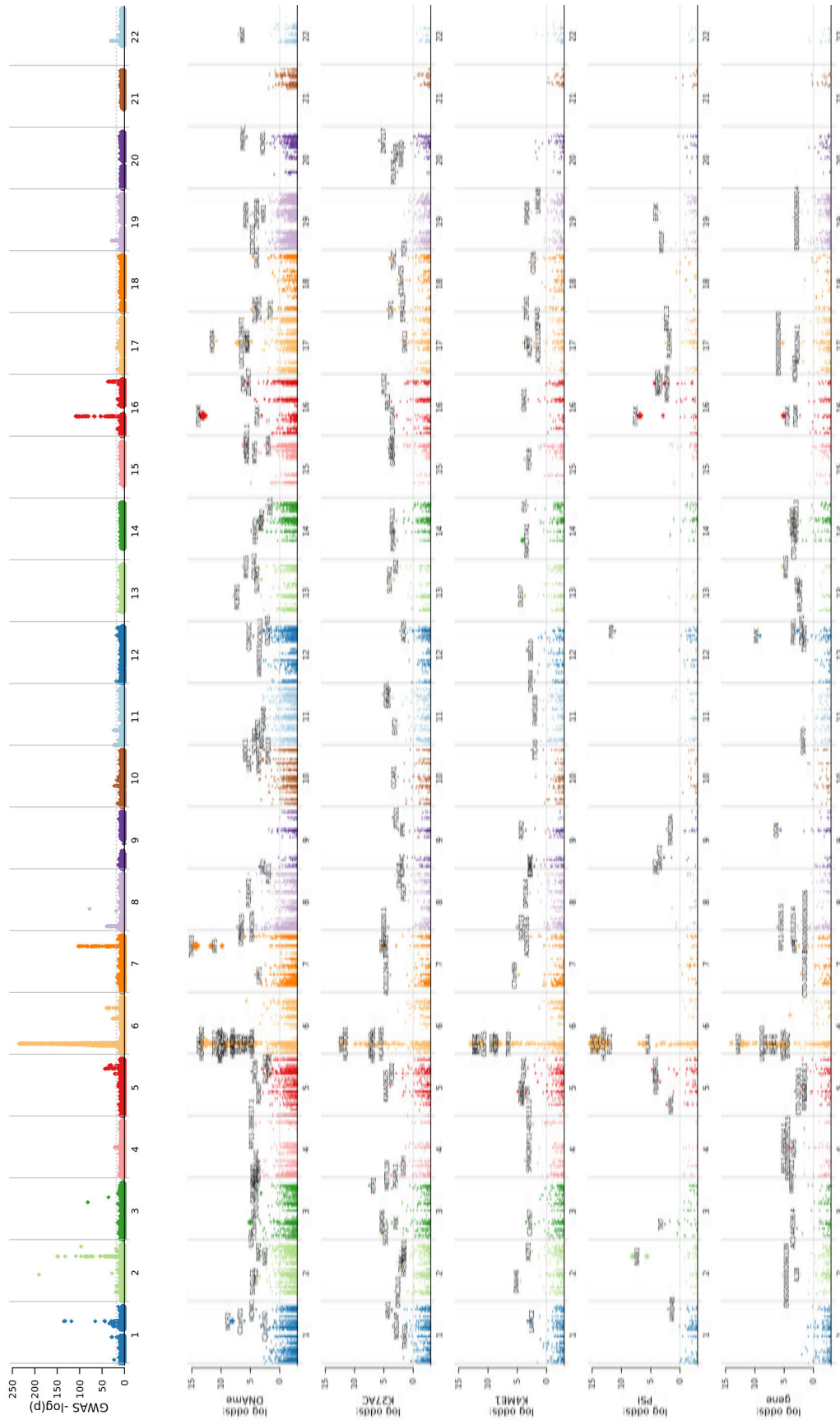


Figure 4-5: Features associated with systemic lupus erythematosus.

This suggests that the GWAS signal in these loci may be mostly due to effects that were not tested in our analysis, such as coding sequence mutations, QTLs active in other cell types, or other molecular phenotypes.

At the top loci, we often observed that the role of a given gene in disease was supported by mediators of multiple molecular phenotypes. However, at other loci, we also often observe strong signals in one molecular phenotype or cell type that are not echoed in others (Figures 4-1 – 4-5). For example, the expression of the *LACC1* gene on chromosome 13 is predicted to be a strong mediator of Crohn’s disease, but there is no corresponding signal in the other molecular phenotypes tested. Another example is the *PHACTR1* locus on chromosome 5, which contains a predicted H3K4me1 mediator for coronary artery disease.

4.3.2 Every QTL type contributes to heritability explained

From the QTL effect sizes and the mediation effect sizes estimated by CaMMEL, we can quantify how much of the local SNP heritability of each trait is captured by each of these molecular phenotypes (Figure 4-6).

For each of the GWAS traits tested, DNA methylation is able to explain more heritability than any of the other molecular phenotypes (Figure 4-6). The heritability explained by H3K27ac, H3K4me1, and expression are relatively similar to each other in each of the diseases. Splicing often explains the least heritability, relative to the other QTL types, but this is not always true. For example, in SLE and CAD, splicing explains a similar amount of heritability to DNAm, H3K27ac, and H3K4me1.

The relative amount of heritability explained by cell type also varies within each trait. In Crohn’s disease, monocyte QTLs explain more DNAm, H3K27ac, H3K4me1, and expression heritability than neutrophil or CD4+ T-cells (Figure 4-6).

Across traits, the heritability explained by these immune cell QTLs also varies (Figure 4-6). The heritability of coronary artery disease is explained the least, followed by Alzheimer’s disease, Crohn’s disease and ulcerative colitis, and finally, systemic lupus erythematosus, which is the best explained by these QTLs. This makes sense given what we know about the immune system involvement of these diseases.

SLE is primarily an autoimmune immune disorder, while Crohn’s disease and ulcerative colitis are autoimmune disorders that also involve other body systems. Further, Alzheimer’s disease is recognized to have an important immune component but is primarily neurodegenerative disease, and the role of inflammation and immunity in cardiovascular disease is still being investigated.

We observed that DNA methylation tends to explain the most trait heritability, relative to the other QTL types (Figure 4-6). However, one factor that may be associated with the amount of heritability a QTL type is able to explain is the number of potential features of that type. Therefore, we calculated the amount of variance explained relative to the number of potential features (Figure 4-7, 4-8).

We observe that the contribution of DNA methylation per potential mediator is the lowest, across all traits, despite the fact that it explains the most variance overall. This implies that while there are many DNA methylation mediators, each individual DNA methylation mediator tends to have a smaller effect than mediators of other QTL types. However, for a given trait, the degree to which DNA methylation mediators have smaller effects also seems to reflect the immune involvement of each trait, where DNAm in CAD explains relatively more heritability per feature than DNAm in SLE.

For the other QTL types (gene expression, splicing, H3K27ac, and H3K4me1), the relative contribution per potential mediator varies across traits. For example, splicing explains the most heritability per mediator in lupus, but the 4th most heritability per mediator in Alzheimer’s disease. In general, however, gene expression and H3K4me1 are enriched in heritability to a similar degree across all five traits, and epigenetic mediators (DNAm, H3K27ac, H3K4me1) are at least as impactful as gene expression mediators, on either a total or per-mediator basis (Figure 4-6, 4-7, 4-8).

For one trait, we explored the overlap in heritability explained by combinations of the QTL studies. The QTLs we have tested here, when all combined, explain 47.6% of the genetic variance in Crohn’s disease susceptibility captured in the GWAS (Figure 4-9). Next, we calculate variance explained through mediation with each QTL type combined across cell types. No single QTL type can recapitulate the total

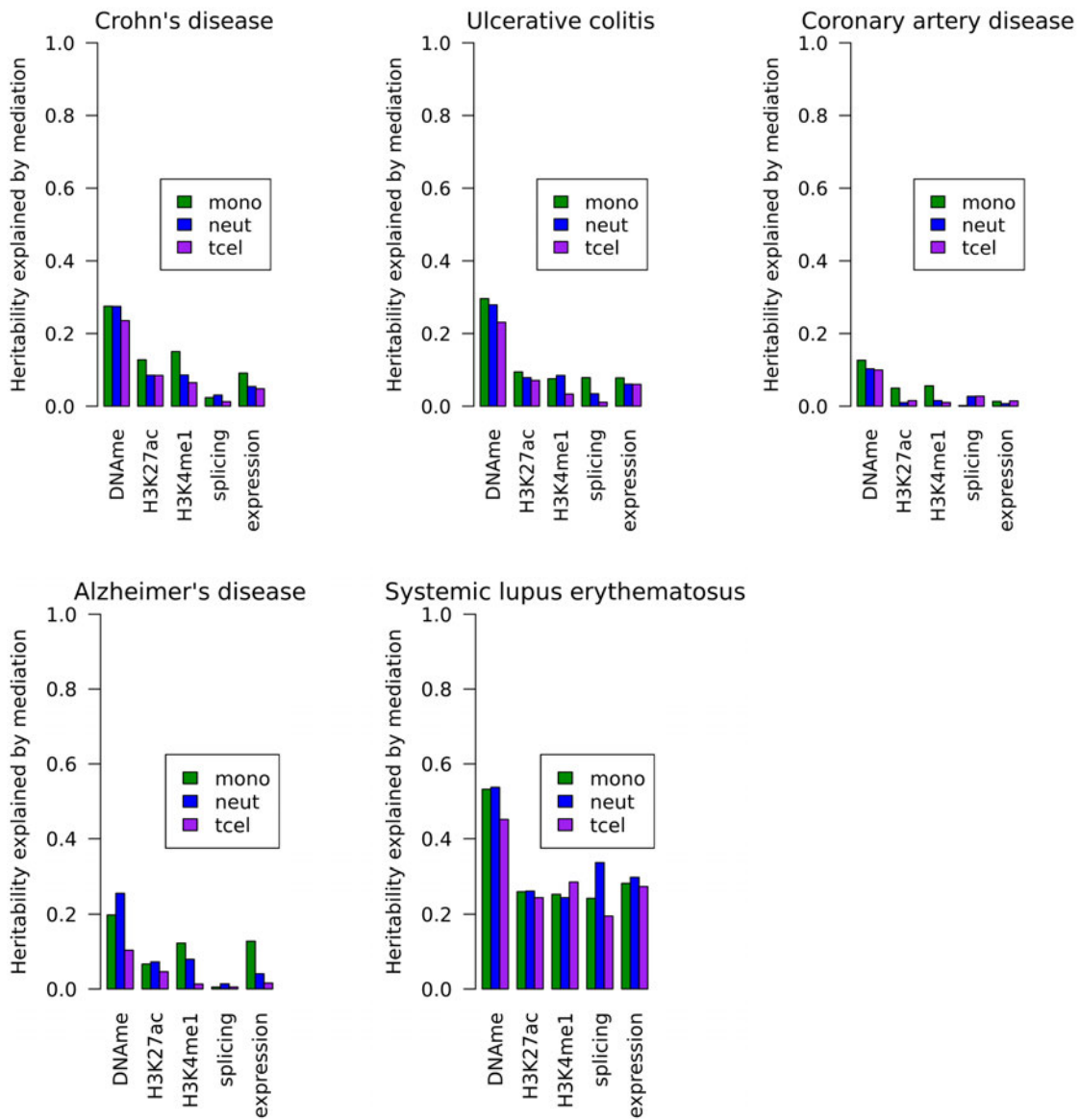


Figure 4-6: Heritability explained by each QTL study, for each GWAS.

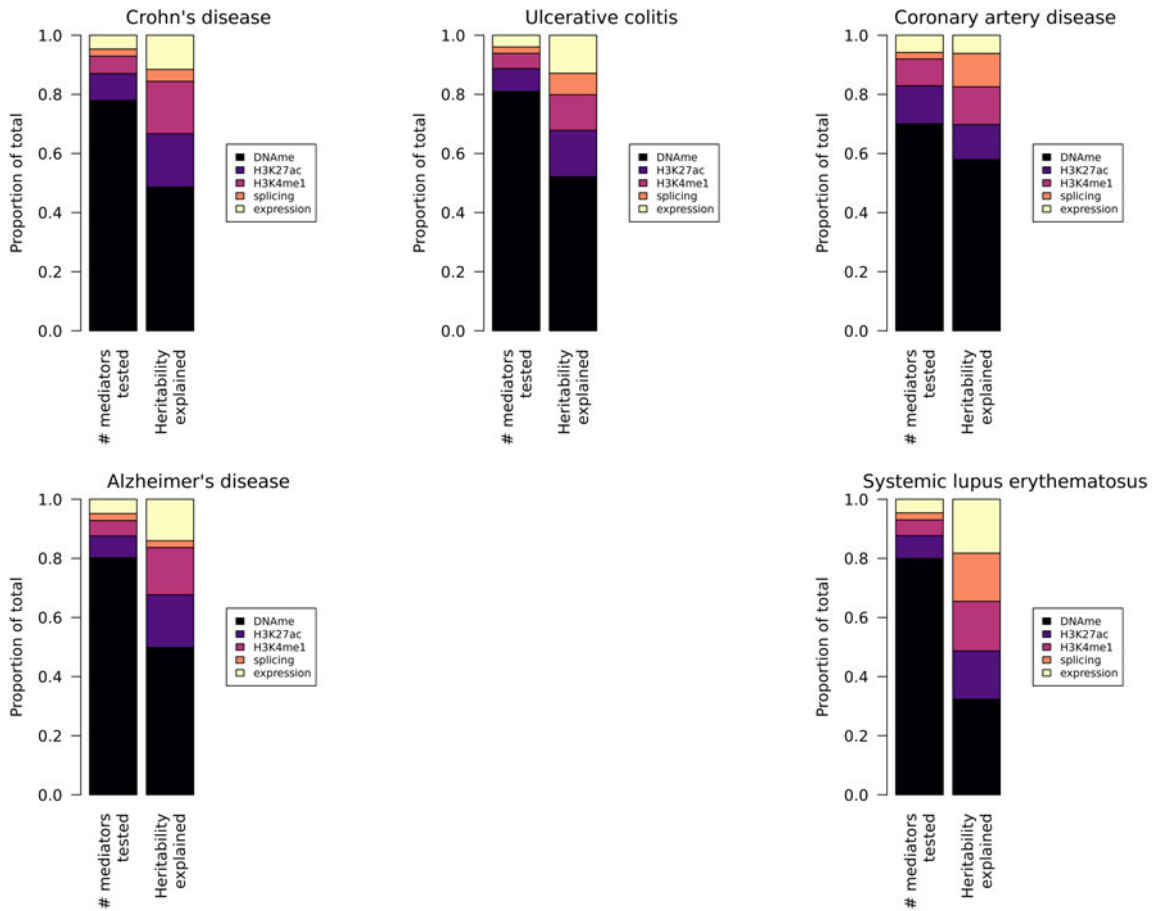


Figure 4-7: Heritability explained by each QTL type, compared to number of features tested, for each GWAS.

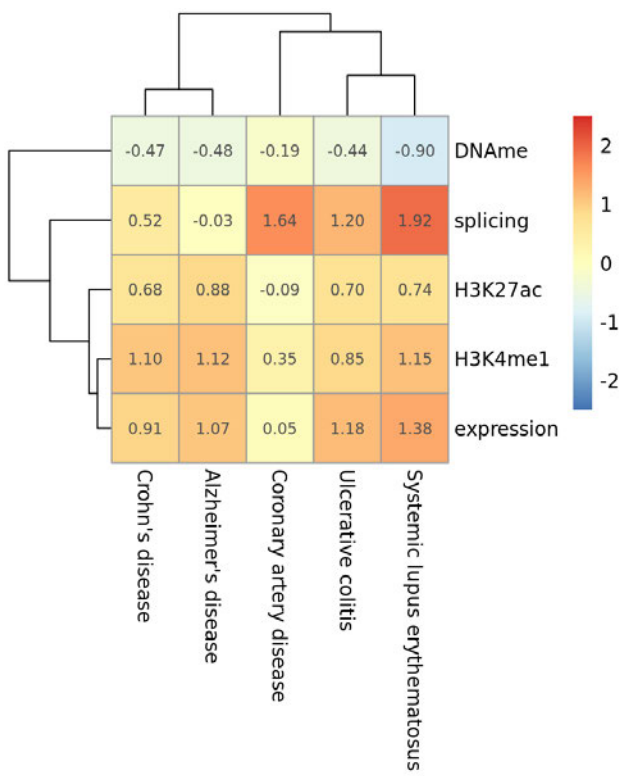


Figure 4-8: Enrichment of heritability explained through mediation relative to number of features tested, for each QTL type and GWAS. The value shown is the natural logarithm of the enrichment.

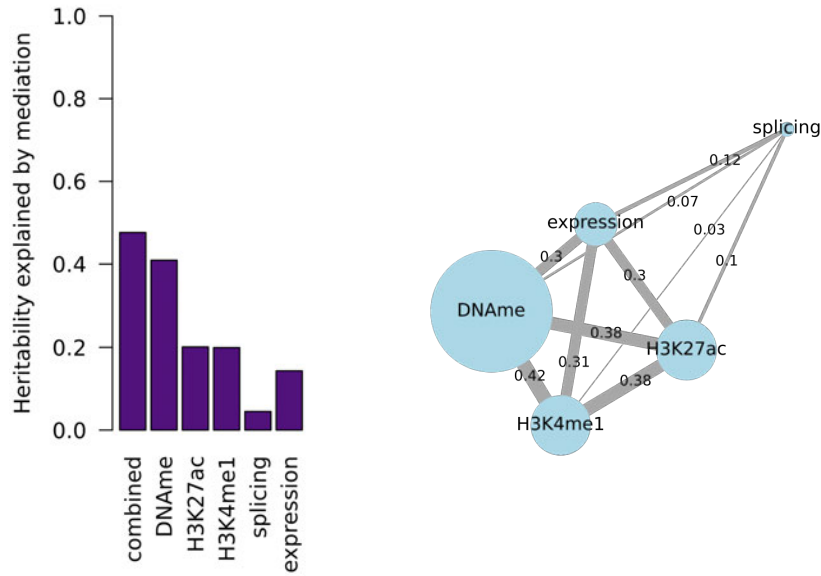


Figure 4-9: Heritability explained by combinations of QTL studies, for Crohn’s disease. Left: Heritability explained by QTL types, combining across cell types. Right: Overlap in heritability explained by QTL types.

variance explained, meaning that using multiple QTL types helps us explain more of the disease genetics. DNA methylation is able to explain the most genetic variance (40.9%), followed by H3K27ac (20.1%), H3K4me1 (19.9%), gene expression (14.3%), and splicing (4.5%).

We observed that in Crohn’s disease, the sum of the heritability explained by each individual QTL type is 99.8%. This is higher than the variance explained by these QTLs combined, indicating that the variance explained by these modalities is overlapping. So, in order to quantify the relationships between the QTL types, we tested how much variance is explained by pairs of QTL types. Using the combined variance and the variance explained by each QTL independently, we can calculate how much the variance explained by each QTL type overlaps, and calculate a similarity metric analogous to Jaccard similarity. We can then plot the QTL types graphically, using the Jaccard similarities as edge weights (Figure 4-9). This shows that the epigenetic QTLs (DNAm, H3K4me1, and H3K27ac) are the most related, while splicing is the least related to the rest of the QTL types. Further, we asked how much additional

heritability can be explained relative to gene expression. Our results indicate that adding DNAm explains an additional 28.2% of heritability, H3K27ac adds 12.2%, and H3K4me1 adds 11.8%. Compared to the baseline of 14.3% of heritability explained by gene expression QTLs, including epigenetic mediators is quite useful.

4.3.3 Genes are linked to disease through varied mechanisms

Earlier, we examined a few of the top predicted mediators. These mediators are individual features, such as histone peaks, CpG sites, splice sites, and genes, which are each associated with a gene. However, to accurately quantify the degree to which a specific gene g is associated with disease through a certain type of molecular phenotype q , we must account for the number of features of type q associated with gene g . The number of associated features varies widely from gene to gene, and genes with more features are more likely to be associated with a high scoring mediator, just by random chance. Therefore, we calculate a normalized score matrix D , where each entry D_{gq} describes how much the mean likelihood of all of the features of type q near gene g deviates from expectation, given the number of tested features. We can calculate D for each cell type separately, as well as for the pooled results.

We can summarize the total mediation signal linking gene g to a trait by adding the normalized mediation score associated with each QTL type for that gene. In general, we observe that genes with higher total mediation signal across QTL types also tend to have higher associated GWAS signal, though the pattern is clearer in some traits (e.g. lupus) than in others (e.g. coronary artery disease) (Figure 4-10).

To investigate the patterns in the way that genes are linked to disease, we first filtered the pooled D matrix for each GWAS to genes that were significantly linked to disease through at least one molecular phenotype (2 standard deviations above expectation). Next, we clustered these data points using k-means clustering, and performed dimensionality reduction with UMAP in order to more easily visualize the data. The largest scale divisions in the data correspond to the scores for each of the QTL types, although there is some overlap between them (Figure 4-11).

The k-means clustering reveals additional structure in the data. The four largest

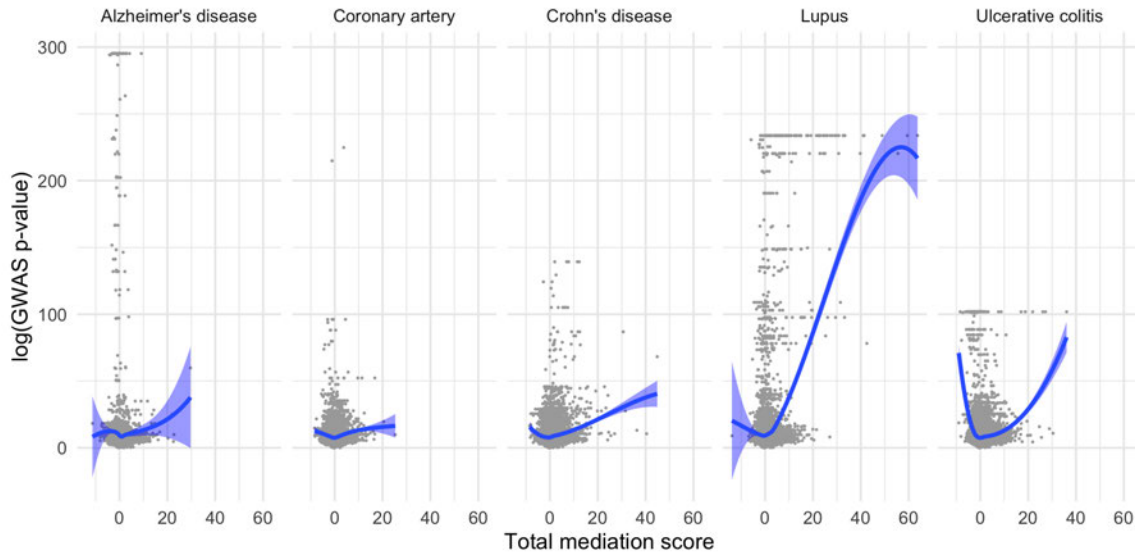


Figure 4-10: GWAS significance as a function of total mediation signal, for each GWAS. Trend line is generated using loess smoothing.

k-means clusters (cluster 11, 5, 6, and 9) correspond to genes linked by only DNA methylation, only H3K27ac, only H3K4me1, and only gene expression, respectively (Figure 4-12,4-13). Cluster 4 also corresponds to genes linked by only DNA methylation, but typically has higher scores than genes in cluster 11. Other clusters correspond to genes linked by both DNA methylation and H3K27ac (cluster 1), genes linked by both DNA methylation and H34me1 (cluster 2), and genes linked by alternative splicing (cluster 7). The smallest k-means clusters (cluster 3, 10, 8, and 12) correspond to genes linked to disease through multiple molecular mechanisms. Cluster 3 is characterized by relatively higher gene expression scores, cluster 10 is characterized by relatively higher epigenetic scores, particularly H3K4me1, and cluster 8 is characterized by relatively higher DNA methylation scores. Genes in cluster 12 have the highest scores overall, with scores more evenly balanced across the QTL types. Taken together, this suggests that many more genes are linked to disease through a single molecular mechanism rather than through combinations of mechanisms, though a small subset of genes are strongly linked through multiple mechanisms.

There are several possible reasons that these molecular mediators overlap less than we might expect, given that DNA methylation, H3K27ac, and H3K4me1 are

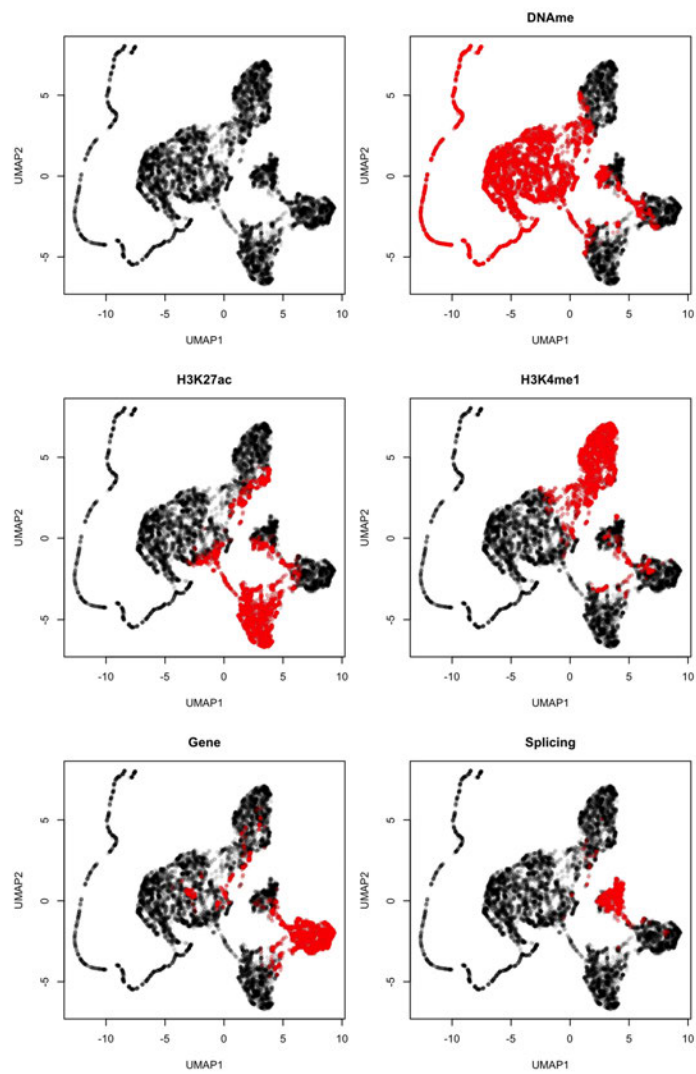


Figure 4-11: UMAP, by QTL.

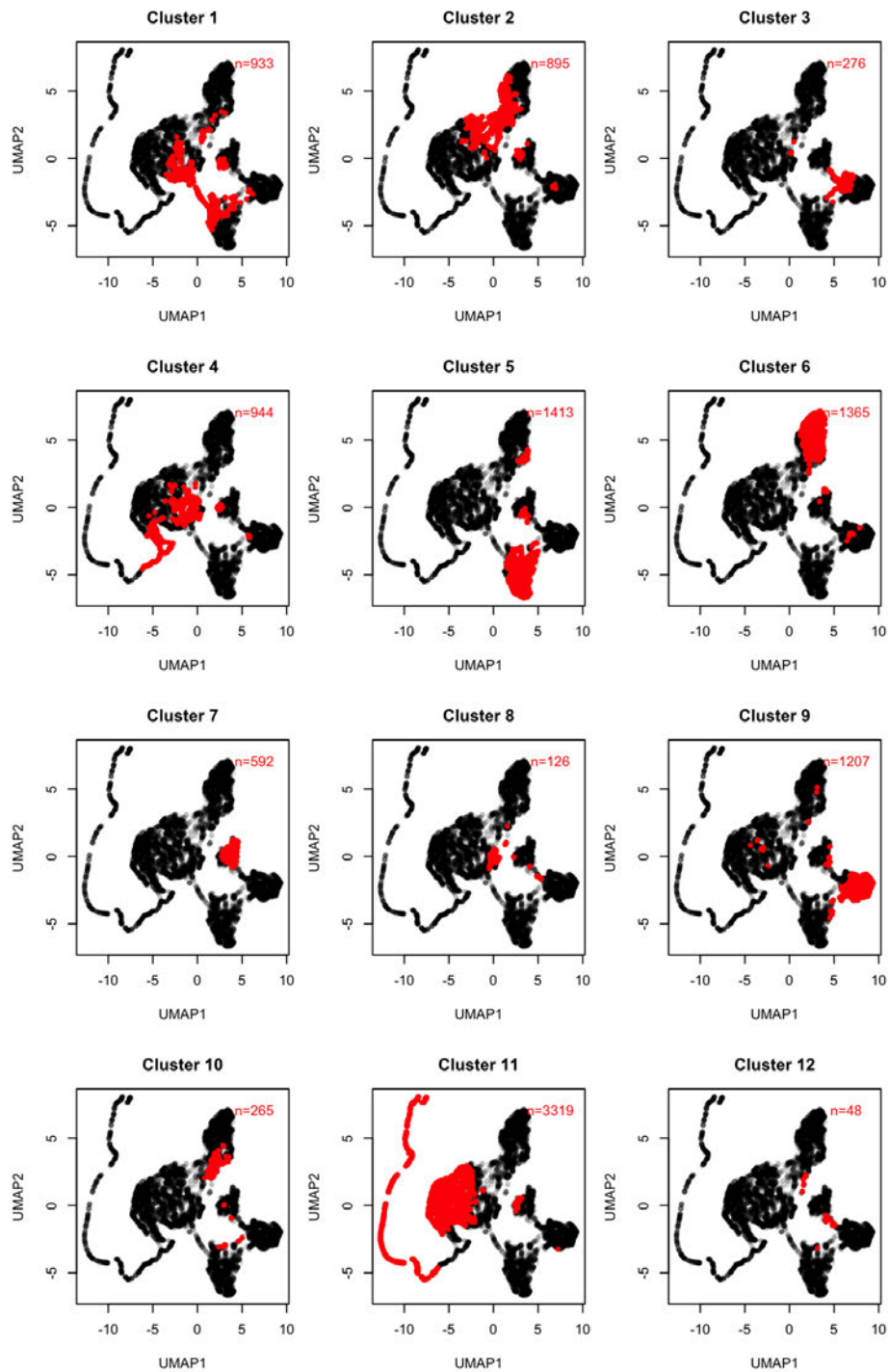


Figure 4-12: UMAP, by cluster.

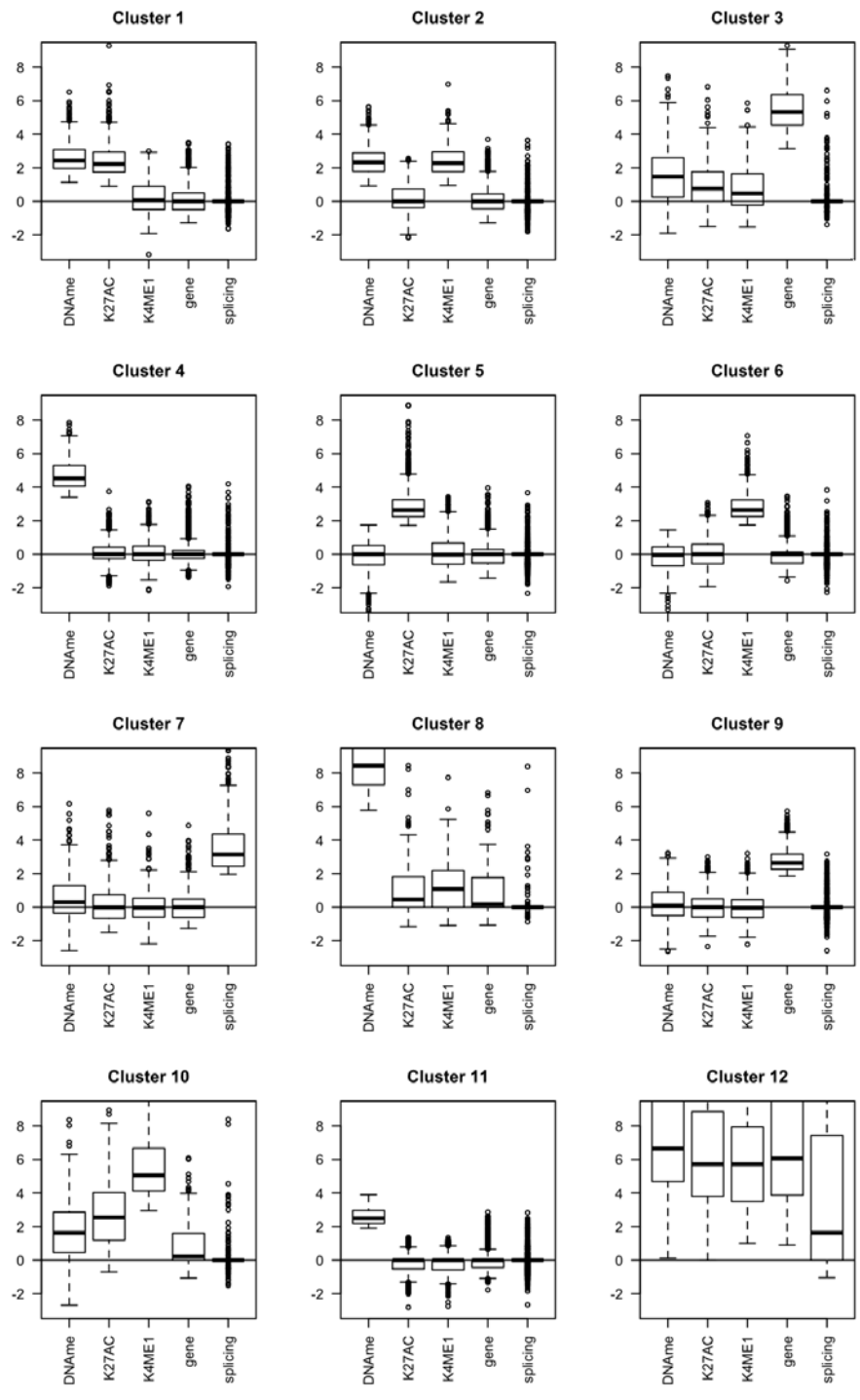


Figure 4-13: Score of points belonging to each cluster

known to be involved in gene expression regulation. First, we assigned epigenetic mediators to the closest gene, but this heuristic is not completely accurate. We may observe greater gene-level overlap in molecular mediators if we are able to use a more accurate method for linking regulatory regions to genes. Even with perfect linking, however, QTLs are often celltype-specific and can even be condition-specific, so we may not observe significant overlap unless we have measured the QTLs in the cell state most relevant to the disease mediation. For instance, genes associated primarily with one type of molecular mediator may not act through the immune cell types we examined, but rather act in a cell type from a different cell lineage or developmental stage. An example of this is AKAP13, which is the gene associated with the mostly likely DNAm mediator for coronary artery disease. AKAP13 has been reported to be essential for cardiac development and function, and to play a role in blood pressure regulation [59, 41, 35]. Another possibility is that genes may be linked to disease only through epigenetic mechanisms in our analysis because the gene expression QTLs are only observable under specific environmental conditions.

4.3.4 Functional dissection of GWAS loci

Molecular mediators provide insight into how GWAS variants are functionally linked to disease. In a GWAS locus with many genes, we may be able to determine which gene(s) are likely to play a role in disease etiology by considering the total molecular mediation evidence for each gene. Furthermore, predicted epigenetic disease mechanisms point to specific regulatory regions that are associated with disease. We can analyze these disease-associated regulatory regions to predict which transcription factors tend to drive disease-associated processes, and what cell types these processes are active in.

For each of our disease traits, we examined genes in GWAS-significant loci that were linked to molecular mediators of several types. In these cases, we observe that molecular mediators of different types do overlap and implicate specific genes and regulatory regions. For example, in the Crohn's disease risk locus near the PTGER4 gene, we see that the expression of the PTGER4 gene, as well as the activity of

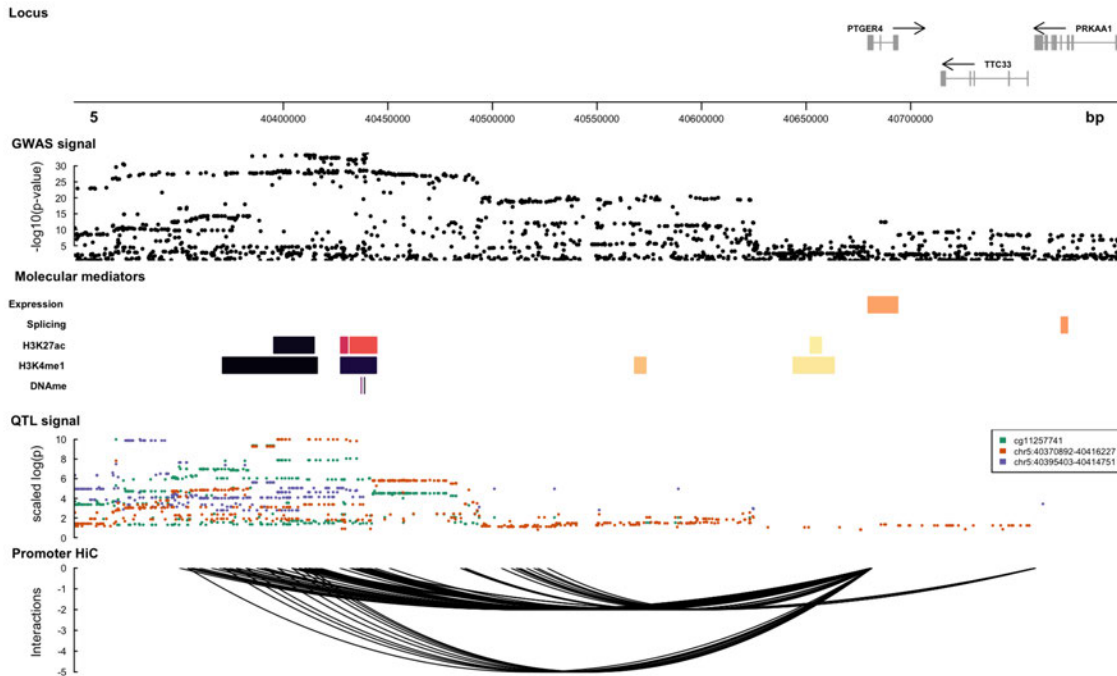


Figure 4-14: Crohn's disease risk locus PTGER4.

several H3K27ac, H3K4me1, and DNAm sites is linked to disease risk (Figure 4-14). The epigenetic mediators overlap each other near chr5:40370000-40440000, and are linked to the PTGER4 gene in the Hi-C data, suggesting that this genomic region regulates Crohn's disease risk through the PTGER4 gene. In ulcerative colitis, we see that both the expression of the gene RBM6 and the H3K4me1 activity of a region that overlaps the promoter of RBM6 are predicted to mediate disease (Figure 4-15). In coronary artery disease, the expression of the gene LIPA, as well as a H3K4me1 peak and two DNA methylation sites overlapping the gene, are predicted to mediate disease (Figure 4-16). Similarly, the role of BIN1 in Alzheimer's is supported by a gene expression mediator and multiple epigenetic mediators, several of which overlap the BIN1 promoter (Figure 4-17), and the role of the integrins ITGAM and ITGAX in systemic lupus erythematosus are supported by several mechanisms (Figure 4-18).

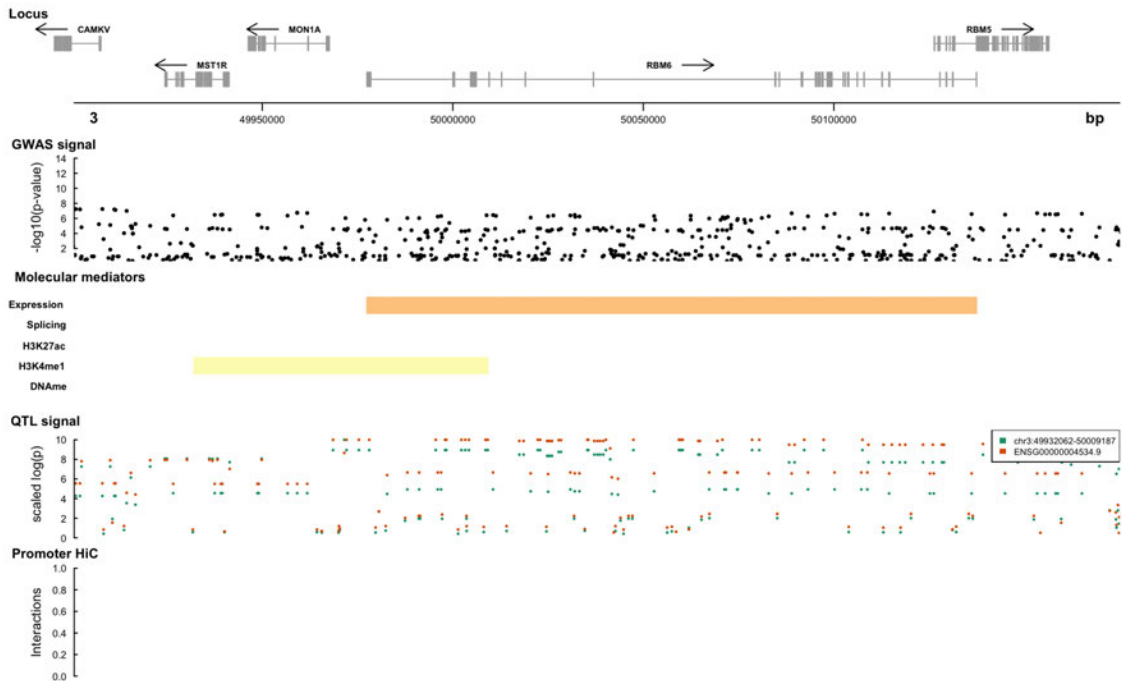


Figure 4-15: Ulcerative colitis risk locus RBM6.

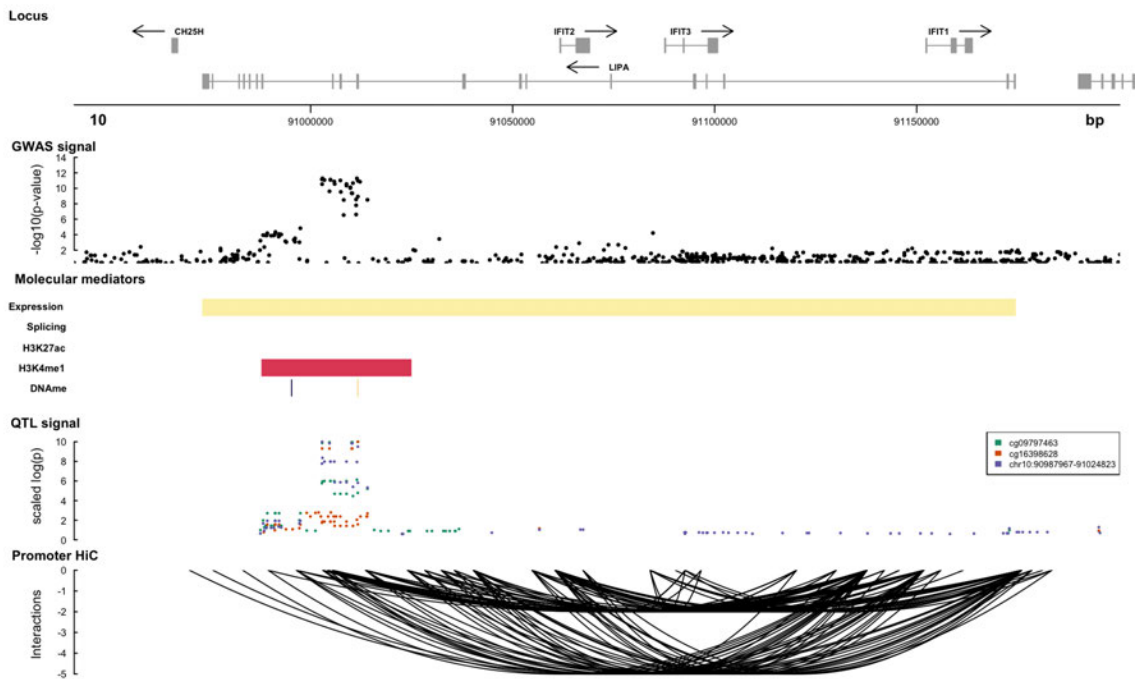


Figure 4-16: Coronary artery disease risk locus LIPA.

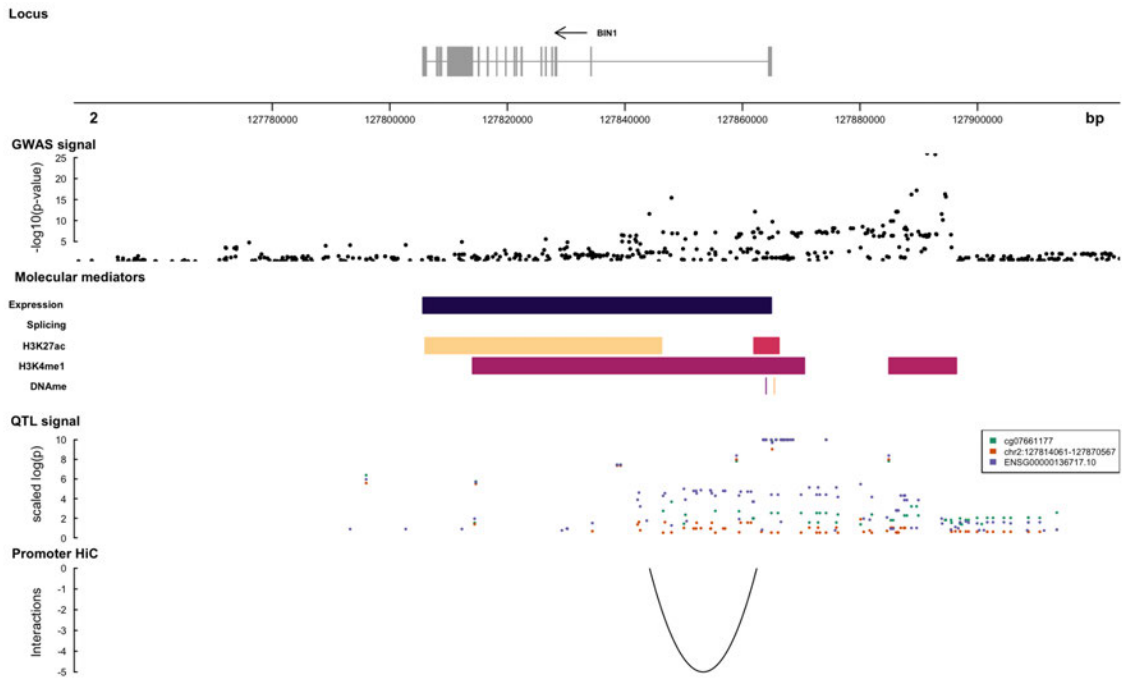


Figure 4-17: Alzheimer's disease risk locus BIN1.

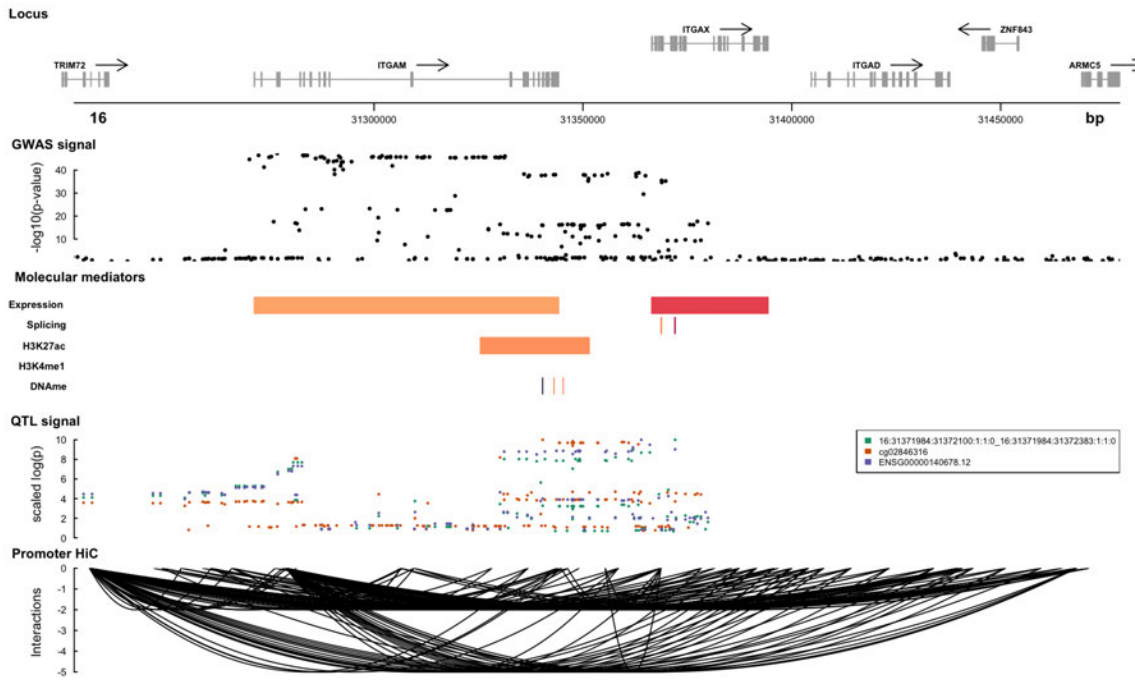


Figure 4-18: Lupus risk locus ITGAX, ITGAM.

4.4 Conclusion

We discovered many potential disease mechanisms by integrating genome-wide association studies for multiple traits with QTL studies in immune cell types for both transcriptional (gene expression, splicing) and epigenetic (DNA methylation, histone marks H3K27ac and H4K4me1) molecular phenotypes using CaMMEL. To our knowledge, this is the first time mediation analysis has been performed with QTLs for histone marks. Our work provides evidence that risk variants can exert their influence on disease through epigenetic mechanisms, as well as through transcriptional mechanisms.

We introduce a method to quantify the impact of disease variants on a given gene through each type of QTL, accounting for the varying number of features assigned to each gene. This allows us to identify genes that are significantly associated with these molecular mechanisms, and explore the relationships between these mechanisms in the context of disease. We find that in most cases, we find evidence to link a gene to disease through only one molecular phenotype, although there are examples of genes that are linked to disease through both expression and histone mark activity, through both DNA methylation and histone marks, and other combinations of mechanisms.

We show that analyzing QTLs associated with different facets of biological regulation improves our ability to identify mechanisms underlying disease compared to expression QTLs alone. In our analysis, some genes known to be associated with disease are only linked through epigenetic mechanisms and not transcriptional mechanisms. Further, we see through analysis of heritability explained that each QTL type provides some information about disease that the others do not. Finally, in loci with epigenetic mediators, we can identify the specific regulatory regions that drive disease and therefore learn about driver regulatory factors and cell types.

This is a promising area for further study. Improved linking of regulatory regions to genes will help us refine our understanding of the relationship between these molecular mediators and prioritize disease-related genes. We may be able to identify genetic variants that influence disease risk by changing how cells respond to environmental

stimuli.

Chapter 5

Conclusion

5.1 Summary of contributions

In this thesis, we have applied computational methods to analyze genetic and epigenetic data, contributing to knowledge about the role of epigenetics in governing phenotypic variance.

In Chapter 2, we discovered a novel chromatin state signature of adult stem cell differentiation, by analyzing chromatin state segmentations generated from genome-wide maps of five histone modification marks in 45 epigenomes, including samples from neural progenitors, fetal and adult brain, mesenchymal and hematopoietic adult stem cells, and mesenchymal and hematopoietic tissues. We found that the "ZNF/Rpts" chromatin state, which is characterized by the H3K9me3 and H3K36me3 histone modifications, is more abundant in multipotent cell types than in fully differentiated cells, particularly on actively transcribed genes associated with biological processes including protein degradation and organelle assembly. This work provides new insight into the function of these epigenetic marks and suggests a new mechanism by which stem cells may be able to maintain their unique cellular state.

In Chapter 3, we profiled the genome-wide changes in chromatin accessibility that occur when somatic cells are transferred into the nucleus of *Xenopus* oocytes, a process that causes the somatic nuclei to undergo transcriptional reprogramming in a DNA-replication-independent manner, and explored the relationship between

chromatin accessibility and transcriptional activation in this system. We found that genes with open chromatin in the somatic cells are more likely to be transcriptionally activated and that closed chromatin can be a barrier to some genes successfully being activated during reprogramming. We also identified two regulators of transcriptional reprogramming from these data, the transcription factor RARA and the protein TOCA1/FNBP1L. This work suggests that the initial epigenetic state of a cell influences its response to nuclear transfer, and adds to our understanding of how to engineer cellular responses.

In Chapter 4, we discovered many potential disease mechanisms by applying mediation analysis to integrate studies measuring the effects of genetic variants on disease (GWAS) with studies measuring the effects of genetic variants on transcriptional and epigenetic molecular phenotypes (QTL studies). We demonstrate for the first time that many risk variants may exert their influence on disease by affecting histone mark activity (H3K27ac and H3K4me1), as well as through DNA methylation, gene expression, and splicing, and show that increasing the variety of molecular phenotypes included in a mediation analysis improves our ability to identify mechanisms underlying disease. We saw that most genes were linked to disease primarily through one molecular phenotype, but we did identify examples of genes that are linked to disease through both expression and histone mark activity, through both DNA methylation and histone marks, and other combinations of mechanisms. This work suggests that epigenetic and transcriptional mechanisms play complementary roles in disease pathogenesis and may enable research into new therapeutic targets.

5.2 Limitations

The contributions I have described each have limitations that restrict their applicability.

5.2.1 Chromatin state signatures of differentiation

The model we formulated to determine if epigenetic features are associated with differentiation is not well suited to jointly studying changes that occur at both early and late stages of differentiation. In our model, we specify the cell lineage of each sample and used a random effect intercept to capture heterogeneity across these lineages, but the lineage of early stage samples like embryonic stem cells is not well defined, since they are the progenitors of many lineages of cells. A linear model without random effects would not have this limitation, but it may become more important to control for sampling biases across groups (e.g. the number of samples per group).

In generating our features, we also assumed that differentiation leads to genome-wide changes in abundance of chromatin state. We made this choice because previous studies had observed this type of change occurring in certain cell differentiation contexts, and because this feature is easily calculated across cell types. However, differentiation might also lead to more localized and/or lineage-specific changes, such as changes in epigenetic marks on regulatory regions controlled by celltype-specific transcription factors. Future work may be able to test these types of features as well.

Furthermore, while the Roadmap Epigenomics Project was the most comprehensive source of epigenomic data at the time, recent work applying imputation has greatly expanded the number of epigenomes available and the number of epigenetic assays available per epigenome[4]. Applying our discovery method to this expanded set of epigenomes would allow us to evaluate the generalizability of our initial findings, add to our understanding of the role of the "ZNF/Rpts" chromatin state, and potentially discover additional epigenetic signatures of stem cell differentiation.

Finally, since we did not find that the loss of "ZNF/Rpts" is strongly linked to changes in gene expression levels, the functional role of "ZNF/Rpts" is still unclear. Additional studies on the other roles epigenetic marks may play in genome regulation, such as regulation of splicing and 3-D organization, may help us establish a stronger functional link to the loss of "ZNF/Rpts". Experimental data targeting

"ZNF/Rpts" specifically may also help us establish whether "ZNF/Rpts" is required for multipotency or simply correlated with a multipotent state.

5.2.2 Impact of chromatin accessibility on transcriptional reprogramming

To more comprehensively understand the *Xenopus* reprogramming system, it would be valuable to study the cell-to-cell variability in chromatin accessibility changes and gene expression activation. We could also analyze data from epigenetic assays other than chromatin accessibility, and compare the epigenetic changes that occur during nuclear transfer in mouse embryonic fibroblasts to those that occur in other somatic cell types. Further, while we observed that the accessibility of a promoter in donor cells is predictive of successful transcriptional activation, we could try to develop a more formal model of this process, in order to learn additional factors that affect reprogramming success.

Finally, this study focused on one specific experimental system, namely transcriptional reprogramming induced by nuclear transfer of somatic nuclei to *Xenopus* oocytes, meaning that caution must be taken in generalizing our conclusions. Our results should be considered in the context of other studies investigating how epigenetic state of a cell influences its response to external stimuli.

5.2.3 Multi-omic mediation analysis

In our multi-omic mediation analysis, we explored the relationship between molecular phenotypes by considering how often genes are associated with different combinations of molecular mediators. This is a sub-optimal way of describing the relationship between molecular phenotypes, though, for several reasons. First, current methods of linking regulatory regions to genes are imperfect, leading to some cases where we erroneously conclude that a given gene is associated with certain type mediator, and to other cases where we erroneously conclude that a gene is not associated with a certain type of mediator. Secondly, even if two different types of mediators are linked

to the same gene, it doesn't necessarily mean that the same variants cause both mediation signals – the risk from the different mediators may be independent. Another computational approach must be developed to clarify the statistical dependence of the mediation signals in each locus.

We also observed that in many cases, the interpretation of the mediation signals is complicated by the sharing of QTLs across cell types. It may be possible to disentangle these effects with computational methods if QTL studies for these molecular phenotypes across many more cell types were available. This data is resource-intensive to collect, but it would also make multi-omic mediation analysis applicable to a wider range of traits. In our study, we observed that the relative amount of information that each QTL type contributed to explaining the heritability of a trait seemed to vary with how immune-related the trait was.

It would also be nice to investigate how the evidence from the molecular phenotypes considered here relates to evidence for other mechanisms, as QTL studies focusing on new molecular phenotypes, such as m⁶A methylation, are being produced every year, and coding sequence mutation is a very important factor in disease risk.

Finally, while the computational evidence suggests that we have identified novel mechanisms of disease pathogenesis, experimental validation will be required to confirm that these are causal mechanisms of disease.

5.3 Looking to the future

Additional basic research into epigenetics is still needed. Currently, the paradigm is that that epigenetic marks regulate transcription by either regulating chromatin accessibility or regulating the binding of transcription factors and other proteins. However, the specific effects of various epigenetic marks are still not clear, and emerging evidence suggests that epigenetic marks may have effects that do not fit neatly into this paradigm, e.g. effects on 3D chromatin conformation [38, 14]. In addition, it is unclear in many cases if making a certain epigenetic change will influence the phenotype of interest, or in other words, if an observed epigenetic mark has a causal effect

or not. A better understanding of the range of effects that epigenetic marks can have, as well as the overlap between the functions of various marks, would help guide the design of studies looking for epigenetic differences between cell types or conditions, and improve the interpretation of existing data.

The computational analysis of epigenetic data also relies heavily on techniques for inferring the function of regulatory regions. These include techniques for linking regulatory regions to genes, biological pathway enrichment, and motif enrichment, among others. Improvement in these methods would have a broad impact on the field of epigenetics.

Addressing these questions will help enable important translational research. Understanding epigenetics is crucial to our ability to solve problems such as how to engineer cellular behavior, how environmental signals may contribute to phenotypic traits, and predicting how different individuals may respond to a medical treatment.

Appendix A

Supplementary information related to Chapter 2

Table A.1: Epigenomes analyzed for differentiation signatures

ID	Epigenome name	Diff. stage	Cell lineage
E006	H1 Derived Mesenchymal Stem Cells	Multipotent	Mesenchymal
E007	H1 Derived Neuronal Progenitor Cultured Cells	Multipotent	Neural
E009	H9 Derived Neuronal Progenitor Cultured Cells	Multipotent	Neural
E010	H9 Derived Neuron Cultured Cells	Terminal	Neural
E023	Mesenchymal Stem Cell Derived Adipocyte Cultured Cells	Terminal	Mesenchymal
E025	Adipose Derived Mesenchymal Stem Cell Cultured Cells	Multipotent	Mesenchymal
E026	Bone Marrow Derived Cultured Mesenchymal Stem Cells	Multipotent	Mesenchymal
E029	Primary monocytes from peripheral blood	Terminal	Hematopoietic
E030	Primary neutrophils from peripheral blood	Terminal	Hematopoietic
E031	Primary B cells from cord blood	Terminal	Hematopoietic
E032	Primary B cells from peripheral blood	Terminal	Hematopoietic
E033	Primary T cells from cord blood	Terminal	Hematopoietic
E034	Primary T cells from peripheral blood	Terminal	Hematopoietic
E035	Primary hematopoietic stem cells	Multipotent	Hematopoietic
E036	Primary hematopoietic stem cells short term culture	Multipotent	Hematopoietic
E037	Primary T helper memory cells from peripheral blood 2	Terminal	Hematopoietic
E038	Primary T helper naive cells from peripheral blood	Terminal	Hematopoietic
E039	Primary T helper naive cells from peripheral blood	Terminal	Hematopoietic
E040	Primary T helper memory cells from peripheral blood 1	Terminal	Hematopoietic
E041	Primary T helper cells PMA-I stimulated	Terminal	Hematopoietic
E042	Primary T helper 17 cells PMA-I stimulated	Terminal	Hematopoietic
E043	Primary T helper cells from peripheral blood	Terminal	Hematopoietic
E044	Primary T regulatory cells from peripheral blood	Terminal	Hematopoietic
E045	Primary T cells effector/memory enriched from peripheral blood	Terminal	Hematopoietic
E046	Primary Natural Killer cells from peripheral blood	Terminal	Hematopoietic
E047	Primary T CD8+ naive cells from peripheral blood	Terminal	Hematopoietic
E048	Primary T CD8+ memory cells from peripheral blood	Terminal	Hematopoietic
E049	Mesenchymal Stem Cell Derived Chondrocyte Cultured Cells	Terminal	Mesenchymal
E050	Primary hematopoietic stem cells G-CSF-mobilized Female	Multipotent	Hematopoietic
E051	Primary hematopoietic stem cells G-CSF-mobilized Male	Multipotent	Hematopoietic
E053	Cortex derived primary cultured neurospheres	Multipotent	Neural
E054	Ganglion Eminence derived primary cultured neurospheres	Multipotent	Neural
E062	Primary mononuclear cells from peripheral blood	Terminal	Hematopoietic
E063	Adipose Nuclei	Terminal	Mesenchymal
E067	Brain Angular Gyrus	Terminal	Neural
E068	Brain Anterior Caudate	Terminal	Neural
E069	Brain Cingulate Gyrus	Terminal	Neural
E070	Brain Germinal Matrix	Multipotent	Neural
E071	Brain Hippocampus Middle	Terminal	Neural
E072	Brain Inferior Temporal Lobe	Terminal	Neural
E073	Brain Dorsolateral Prefrontal Cortex	Terminal	Neural
E074	Brain Substantia Nigra	Terminal	Neural
E081	Fetal Brain Male	Multipotent	Neural
E082	Fetal Brain Female	Multipotent	Neural
E129	Osteoblast Primary Cells	Terminal	Mesenchymal

Appendix B

Supplementary information related to Chapter 4

Table B.1: Genome-wide association studies analyzed

Trait	Reference	ncase	ncontrol	Description
Crohn's disease	[54]	5, 956	14, 927	European CD 1KGP imputed GWAS meta-analysis
Ulcerative colitis	[54]	6, 968	20, 464	European UC 1KGP imputed GWAS meta-analysis
Cardiovascular disease	[71]	60, 801	123, 504	CARDIoGRAMplusC4D 1000 Genomes-based GWAS, additive model
Systemic lupus erythematosus	[9]	7, 219	15, 991	
Multiple sclerosis	[39]	9, 772	17, 376	
Alzheimer's disease	[50]	17, 008	37, 154	Stage 1

Table B.2: Molecular mediators associated with Crohn's disease (top 10 per QTL type, ranked by log odds)

QTL type	Mediator	Location (hg19)	GWAS p-value	lodds (mono)	lodds (neut)	lodds (tcel)	Closest coding gene
DNAme	cg11257741	chr5:40438831-40438831	6.967e-35	13.1	14.3	5.17	PTGER4
DNAme	cg14541801	chr5:40437301-40437301	6.967e-35	4.11	-0.0891	8.4	PTGER4
DNAme	cg06885583	chr11:10715445-10715445	4.516e-05	7.39	-4.52	-4.52	MRVI1
DNAme	cg03902329	chr6:33291306-33291306	7.712e-15	-3.1	7.33	-4.19	DAXX
DNAme	cg23553762	chr14:69735749-69735749	4.339e-08	-4.26	6.93	-0.294	GALNTL1
DNAme	cg10116893	chr17:9019124-9019124	2.225e-05	3.39	6.88	-2.52	NTN1
DNAme	cg08290212	chr11:62473659-62473659	6.768e-07	-4.32	-4.06	6.86	BSCL2, GNG3
DNAme	cg22024760	chr1:40942625-40942625	1.530e-05	-2.85	6.76	-4.11	ZNF642
DNAme	cg21286967	chr6:31696710-31696710	4.519e-15	6.72	-3.18	-3.99	DDAH2
DNAme	cg17309904	chr7:1960073-1960073	9.329e-06	-3.83	NA	6.57	MAD1L1
K27AC	chr5:40395403-40414751	chr5:40395403-40414751	6.967e-35	13.4	-3.6	NA	PTGER4
K27AC	chr6:167319628-167323536	chr6:167319628-167323536	1.025e-12	-4.36	-1.27	8.35	RP11-514O12.4
K27AC	chr5:14430913-14432702	chr5:14430913-14432702	5.709e-06	-4.57	NA	6.46	FAM105A
K27AC	chr5:40427379-40430778	chr5:40427379-40430778	6.967e-35	NA	5.75	NA	PTGER4
K27AC	chr8:126538085-126543685	chr8:126538085-126543685	6.396e-12	5.35	-2.73	NA	TRIB1
K27AC	chr7:33131012-33150504	chr7:33131012-33150504	2.748e-05	-4.54	5.24	-4.42	RP9
K27AC	chr2:25107703-25111521	chr2:25107703-25111521	2.649e-07	5.06	-4.34	-3.82	ADCY3
K27AC	chr3:161027465-161031857	chr3:161027465-161031857	3.180e-06	4.81	-4.28	-2.88	SPTSSB
K27AC	chr10:14653943-14662409	chr10:14653943-14662409	1.026e-05	-4.58	-4.21	4.75	FRMD4A
K27AC	chr3:149687095-149690371	chr3:149687095-149690371	8.079e-05	NA	4.66	NA	ANKUB1
K4ME1	chr5:40370892-40416227	chr5:40370892-40416227	6.967e-35	13.9	NA	-3.93	PTGER4
K4ME1	chr5:40427379-40444505	chr5:40427379-40444505	6.967e-35	-4.52	11.3	12.3	PTGER4
K4ME1	chr11:11156249-11181699	chr11:11156249-11181699	4.516e-05	7.59	NA	NA	ZBED5
K4ME1	chr6:32664121-32670135	chr6:32664121-32670135	7.712e-15	6.43	-2.77	-4.19	HLA-DQB1
K4ME1	chr5:150222268-150235431	chr5:150222268-150235431	8.887e-18	5.79	-2.02	-2.39	IRGM
K4ME1	chr6:32061779-32066621	chr6:32061779-32066621	4.519e-15	-2.51	5.38	NA	TNXB
K4ME1	chr3:160934879-160963756	chr3:160934879-160963756	3.180e-06	NA	5.21	-0.539	B3GALNT1
K4ME1	chr12:13177121-13182213	chr12:13177121-13182213	1.870e-07	-0.801	4.86	-2.81	KIAA1467
K4ME1	chr4:101769956-101774062	chr4:101769956-101774062	1.210e-08	NA	4.81	NA	EMCN
K4ME1	chr18:74688455-74882842	chr18:74688455-74882842	5.274e-05	NA	4.72	-4.55	MBP
PSI	7:87910395:87911955:2:2:1	chr7:87910395-87911955	3.950e-05	NA	8.45	NA	STEAP4
PSI	7:87910395:87912367:2:2:0						
PSI	7:87910395:87913128:2:2:1						
PSI	18:11897620:11905736:2:2:1	chr18:11897620-11905736	1.082e-15	-4.5	-1.01	6.73	MPPE1
PSI	11:1491715:1492524:2:2:0	chr11:1491715-1492524	9.505e-05	NA	NA	4.99	MOB2
PSI	20:62202838:62202888:2:2:0	chr20:62202838-62202888	5.441e-11	NA	4.84	NA	PRIC285
PSI	6:128040969:128128360:2:2:1	chr6:128040969-128128360	3.511e-05	NA	NA	4.34	THEMIS
PSI	7:87910395:87913128:2:2:1	chr7:87910395-87913128	3.950e-05	NA	4.02	NA	STEAP4
PSI	19:11152237:11168930:1:1:1	chr19:11152237-11168930	1.628e-10	3.81	NA	-4.5	SMARCA4
PSI	19:53352467:53393493:2:2:0	chr19:53352467-53393493	3.288e-05	NA	NA	3.6	ZNF28, ZNF320, ZNF468
PSI	16:86582184:86585638:2:2:1	chr16:86582184-86585638	2.472e-05	3.57	NA	-3.88	MTHFSD
PSI	2:153512913:153514402:2:2:1	chr2:153512913-153514402	7.823e-05	NA	3.43	-4.58	PRPF40A
gene	ENSG00000168917	chr3:136537489-136574734	1.202e-05	11.1	NA	-4.48	TMEM22
gene	ENSG0000013583	chr12:13127798-13153207	1.870e-07	-2.21	-1.17	9.04	HEBP1
gene	ENSG00000177112	chr11:10562819-10621479	4.516e-05	-3.22	6.95	NA	
gene	ENSG00000167208	chr16:50700211-50715264	3.245e-61	6.93	-3.74	-4.08	SNX20
gene	ENSG00000026297	chr6:167342992-167370679	1.025e-12	5.84	-3.97	-2.28	RNASSET2
gene	ENSG00000179630	chr13:44453420-44468067	6.021e-09	5.03	-4.53	-2.98	LACC1
gene	ENSG00000253772	chr5:159674524-159675224	9.061e-17	4.95	NA	-4.44	
gene	ENSG00000265735	chr9:9442060-9442347	4.617e-05	NA	4.91	-3.51	
gene	ENSG00000256427	chr12:9399093-9410556	3.648e-05	-4.48	NA	4.67	
gene	ENSG00000235008	chr6:167114582-167116327	1.025e-12	0.41	4.61	-0.481	

Table B.3: Molecular mediators associated with ulcerative colitis (top 10 per QTL type, ranked by log odds)

QTL type	Mediator	Location (hg19)	GWAS p-value	lodds (mono)	lodds (neut)	lodds (tcel)	Closest coding gene
DNAme	cg16492584	chr5:96139282-96139282	1.884e-05	9.79	6.88	-1.97	ERAP1
DNAme	cg04131431	chr1:202183387-202183387	4.581e-13	9.6	NA	-3.79	LGR6
DNAme	cg11986643	chr6:32634316-32634316	2.187e-53	-2.28	-3.36	9.34	HLA-DQB1
DNAme	cg15956162	chr14:39605657-39605657	1.929e-05	NA	7.84	NA	SIP1
DNAme	cg24634333	chr6:111887834-111887834	4.148e-07	7.78	3.28	-1.93	TRAF3IP2
DNAme	cg11821200	chr5:40685989-40685989	4.949e-08	-4.34	7.36	-4.16	PTGER4
DNAme	cg14747132	chr12:74561645-74561645	3.416e-05	NA	NA	7.31	ATXN7L3B
DNAme	cg05617030	chr2:62824045-62824045	5.194e-13	7.22	NA	NA	EHBP1
DNAme	cg00272484	chr8:143400520-143400520	6.381e-05	7.18	NA	-1.16	TSNARE1
DNAme	cg23475955	chr11:1985762-1985762	5.723e-07	7.18	-4.4	-4.18	TNNT3
K27AC	chr6:32439081-32443374	chr6:32439081-32443374	2.187e-53	7.41	-4.11	-3.27	HLA-DRA
K27AC	chr6:149982764-149987876	chr6:149982764-149987876	6.085e-05	-3.59	-4.31	6.24	KATNA1
K27AC	chr2:40548419-40550683	chr2:40548419-40550683	1.024e-05	NA	NA	5.71	SLC8A1
K27AC	chr6:32545108-32547703	chr6:32545108-32547703	2.187e-53	-0.9	5.66	1.13	HLA-DRB1
K27AC	chr17:4421763-4424262	chr17:4421763-4424262	4.578e-05	5.51	NA	-4.02	SPNS2
K27AC	chr2:102305027-102307670	chr2:102305027-102307670	2.120e-07	5.29	-4.05	NA	MAP4K4
K27AC	chr1:181055127-181062607	chr1:181055127-181062607	3.183e-06	-4.5	-4.12	5.16	IER5
K27AC	chr12:62974744-62978095	chr12:62974744-62978095	2.208e-05	-1.68	-3.03	5.03	C12orf61
K27AC	chr2:43810816-43812939	chr2:43810816-43812939	1.024e-05	4.96	-4.46	-4.48	THADA
K27AC	chr6:10642295-10645371	chr6:10642295-10645371	5.921e-05	-2.88	-4.39	4.9	GCNT6
K4ME1	chr6:32619385-32643950	chr6:32619385-32643950	2.187e-53	-0.345	12.1	1.01	HLA-DQB1
K4ME1	chr5:158733049-158741846	chr5:158733049-158741846	3.994e-11	-4.53	7.15	-4.19	IL12B
K4ME1	chr1:162399796-162402861	chr1:162399796-162402861	3.778e-17	6.91	-4.42	-2.13	SH2D1B
K4ME1	chr6:32678465-32684359	chr6:32678465-32684359	2.187e-53	6.67	-3.37	NA	HLA-DQA2
K4ME1	chr6:150397168-150398741	chr6:150397168-150398741	6.085e-05	-3.63	NA	6.28	ULBP3
K4ME1	chr7:2896402-2925110	chr7:2896402-2925110	1.213e-08	-4.32	-4.37	5.65	GNA12
K4ME1	chr3:52729495-52761947	chr3:52729495-52761947	4.762e-07	NA	-4.55	5.33	GLT8D1, SPCS1
K4ME1	chr14:45595089-45614836	chr14:45595089-45614836	5.650e-05	NA	4.03	NA	FANCM, FKBP3
K4ME1	chr22:38872985-38913195	chr22:38872985-38913195	2.986e-06	-4.45	-4.05	3.92	DDX17
K4ME1	chr18:55221460-55224296	chr18:55221460-55224296	3.612e-05	3.89	NA	NA	FECH
PSI	6:32485530:32548023:2:2:0 _	chr6:32485530-32548023	2.187e-53	8.03	6.33	NA	HLA-DRB1, HLA-DRB5
PSI	6:32546882:32548023:2:2:1 _						CTSC
PSI	11:88059568:88061279:2:2:1 _	chr11:88059568-88061279	9.546e-05	NA	6.9	NA	
PSI	11:88059612:88061279:2:2:1 _						
PSI	1:8021796:8022822:1:1:1 _	chr1:8021796-8022822	2.288e-09	6.49	NA	-2.88	PARK7
PSI	1:8021854:8022822:1:1:1 _						
PSI	2:42552695:42553184:1:1:0 _	chr2:42552695-42553184	1.024e-05	NA	5.87	-3.83	EML4
PSI	2:42552695:42553293:1:1:1 _						
PSI	16:70502872:70504216:1:1:1 _	chr16:70502872-70504216	3.189e-06	NA	NA	5.73	FUK
PSI	16:70503227:70504216:1:1:1 _						
PSI	16:19528508:19533025:2:2:1 _	chr16:19528508-19533025	4.968e-06	5.25	-4.52	-4.56	GDE1
PSI	16:19528512:19533025:2:2:1 _						
PSI	1:161594430:161595480:2:2:0 _	chr1:161594430-161595480	3.778e-17	NA	4.9	NA	FCGR2B, FCGR3A, FCGR3B
PSI	1:161594430:161595767:2:2:0 _						
PSI	1:161594430:161595934:2:2:1 _						
PSI	1:161594430:161599500:2:2:0 _						
PSI	1:161594430:161599567:2:2:0 _						
PSI	12:57869399:57869558:2:2:1 _	chr12:57869399-57869558	4.874e-06	-0.867	1.79	4.51	ARHGAP9, MARS
PSI	12:57869399:57869615:2:2:1 _						
PSI	19:32955691:32958379:1:1:1 _	chr19:32955691-32958379	5.839e-08	-4.61	4.45	NA	DPY19L3
PSI	19:32955691:32959636:1:1:1 _						
PSI	6:20480189:20481436:1:1:1 _	chr6:20480189-20481436	3.049e-05	NA	4.28	NA	E2F3
PSI	6:20480189:20481454:1:1:1 _						
gene	ENSG00000164307	chr5:96096521-96143803	1.884e-05	11.3	-4.49	-4.3	ERAP1
gene	ENSG00000255389	chr6:111921078-111923498	4.148e-07	6.11	6.17	10.3	
gene	ENSG00000227036	chr17:70319264-70636611	1.235e-05	7.54	-2.3	-3.66	
gene	ENSG00000155542	chr5:56205087-56221359	4.342e-05	-4.57	6.08	-3.05	C5orf35
gene	ENSG00000225364	chr5:132284185-132284410	7.628e-09	-4.51	-3.47	5.75	
gene	ENSG00000232629	chr6:32723875-32731311	2.187e-53	5.52	5.7	-0.265	HLA-DQB2
gene	ENSG00000166986	chr12:57869228-57911352	4.874e-06	-4.07	NA	5.24	MARS
gene	ENSG00000236468	chr1:199356269-199357603	4.581e-13	-4.56	-4.39	4.87	
gene	ENSG00000212769	chr10:85841185-85841751	1.427e-05	-4.46	NA	4.83	
gene	ENSG0000025772	chr20:43570771-43589127	1.385e-15	3.42	4.57	-3.44	TOMM34

Table B.4: Molecular mediators associated with coronary artery disease (top 10 per QTL type, ranked by log odds)

QTL type	Mediator	Location (hg19)	GWAS p-value	lodds (mono)	lodds (neut)	lodds (tcel)	Closest coding gene
DNAme	cg09874482	chr15:86098694-86098694	4.770e-05	-4.12	NA	7.97	AKAP13
DNAme	cg22464447	chr15:91054211-91054211	5.010e-08	7.26	7.68	4.67	CRTC3
DNAme	cg09797463	chr10:90995297-90995297	5.150e-12	6.78	-4.02	-4.53	LIPA
DNAme	cg04963535	chr15:79603582-79603582	4.440e-16	-4.53	-3.64	5.71	TMED3
DNAme	cg18338521	chr10:69995036-69995036	9.120e-05	0.66	5.69	-1.7	ATOH7
DNAme	cg18164784	chr5:140762229-140762229	9.560e-05	2.85	5.66	-3.66	PCDHGA1, PCDHGA2, PCDHGA3, PCDHGA4, PCDHGA5, PCDHGA6, PCDHGA7, PCDHGB1, PCDHGB2, PCDHGB3
DNAme	cg01152302	chr15:84323770-84323770	1.180e-05	5.47	-0.512	-0.766	ADAMTSL3
DNAme	cg07586081	chr14:23772235-23772235	9.210e-05	-4.44	5.46	-4.51	PPPIR3E
DNAme	cg23750338	chr8:142222470-142222470	3.460e-06	5.24	0.685	-4.49	SLC45A4
DNAme	cg20999565	chr14:68265668-68265668	7.560e-05	-4.21	5.07	NA	ZFYVE26
K27AC	chr7:106164152-106166936	chr7:106164152-106166936	6.260e-06	-4.06	6.21	-4.47	C7orf74
K27AC	chr10:32247181-32254209	chr10:32247181-32254209	1.020e-05	4.28	-2.83	-1.53	ARHGAP12
K27AC	chr7:12104757-12108382	chr7:12104757-12108382	5.280e-05	-1.87	4.25	-4.33	TMEM106B
K27AC	chr10:43274376-43275919	chr10:43274376-43275919	5.550e-15	3.92	-3.76	NA	BMS1
K27AC	chr11:10140425-10142269	chr11:10140425-10142269	1.280e-08	3.32	NA	NA	SBF2
K27AC	chr7:150275384-150286353	chr7:150275384-150286353	1.690e-09	2.92	-3.86	-4.54	GIMAP4
K27AC	chr19:44266361-44290433	chr19:44266361-44290433	7.070e-11	-4.27	-4.08	2.49	KCNN4
K27AC	chr10:32235003-32239213	chr10:32235003-32239213	1.020e-05	-2.59	2.47	-4.28	ARHGAP12
K27AC	chr15:86032442-86037002	chr15:86032442-86037002	4.770e-05	2.39	2.43	-4.36	AKAP13
K27AC	chr10:76153257-76154852	chr10:76153257-76154852	2.670e-05	2.42	-4.1	NA	ADK
K4ME1	chr7:12246689-12275830	chr7:12246689-12275830	5.280e-05	-0.005	6.12	-4.57	TMEM106B
K4ME1	chr6:12897440-12910698	chr6:12897440-12910698	1.810e-42	4.94	1.76	-4.53	PHACTR1
K4ME1	chr8:97375574-97379389	chr8:97375574-97379389	1.820e-05	-4.08	-4.39	4.91	KB-1448A5.1
K4ME1	chr10:26039796-26048920	chr10:26039796-26048920	1.390e-05	NA	4.46	NA	MYO3A
K4ME1	chr2:160547720-160679069	chr2:160547720-160679069	1.790e-05	3.58	-3.95	-4.49	CD302, MARCH7
K4ME1	chr7:20833636-20841301	chr7:20833636-20841301	6.440e-06	3.52	-4.23	-4.51	SP8
K4ME1	chr15:83496424-83504647	chr15:83496424-83504647	1.180e-05	3.51	-4.32	-3.63	WHAMM
K4ME1	chr10:90987967-91024823	chr10:90987967-91024823	5.150e-12	3.08	-3.48	-3.88	CH25H
K4ME1	chr1:109779241-109861456	chr1:109779241-109861456	1.970e-23	2.89	-1.11	NA	CELSR2, MYBPHL, PSRC1
K4ME1	chr5:122235191-122240028	chr5:122235191-122240028	6.000e-05	NA	-3.84	2.84	SNX24
PSI	5:96235894:96237209:1:1:1 5:96235950:96237209:1:1:1	chr5:96235894-96237209	3.510e-05	5.22	NA	5.06	CTD- 2260A17.2, ERAP2
PSI	15:99467914:99472786:1:1:1 15:99467914:99472789:1:1:1 1:154422457:154437609:1:1:1 1:154427058:154437609:1:1:1	chr15:99467914-99472786	4.520e-05	NA	4.43	NA	IGF1R
PSI	1:154422457:154437609:1:1:1 1:154427058:154437609:1:1:1	chr1:154422457-154437609	2.600e-09	-1.55	-0.0853	3.98	IL6R
PSI	21:35304047:35310528:1:1:1 21:35308169:35310528:1:1:0	chr21:35304047-35310528	1.330e-15	NA	NA	3.79	
PSI	10:43656044:43657647:1:1:1 10:43656044:43659313:1:1:1 10:43656044:43662451:1:1:0	chr10:43656044-43657647	5.550e-15	NA	3.76	NA	CSGALNACT2
PSI	16:84853740:84872027:1:1:1 16:84870042:84872027:1:1:1	chr16:84853740-84872027	8.280e-06	NA	3.53	NA	CRISPLD2
PSI	2:228416834:228418419:1:1:1 2:228416834:228418425:1:1:1	chr2:228416834-228418419	1.510e-06	-4.51	-3.31	2.53	AGFG1
PSI	6:31237303:31237742:2:2:1 6:31237303:31237986:2:2:0 6:31237303:3122883:2:2:0 6:31237303:31323093:2:2:0	chr6:31237303-31237742	9.600e-07	-1.08	2.28	-1.79	HLA-C
PSI	15:91009669:91010681:1:1:1 15:91009669:91016071:1:1:0	chr15:91009669-91010681	5.010e-08	NA	2.23	NA	IQGAP1
PSI	8:53586830:53588928:2:2:0 8:53586835:53588928:2:2:1	chr8:53586830-53588928	7.980e-05	NA	2.08	NA	RB1CC1
gene	ENSG00000234297	chr9:15055055-15056050	8.020e-05	6.76	-4.59	-4.57	
gene	ENSG00000106460	chr7:12250867-12276886	5.280e-05	-4.06	6.43	-2.51	TMEM106B
gene	ENSG00000148175	chr9:124101355-124132531	7.230e-07	-2.84	4.94	-4.56	STOM
gene	ENSG00000234335	chr10:32391450-32392221	1.020e-05	NA	4.43	-4.04	
gene	ENSG00000006747	chr7:12610203-12693228	5.280e-05	-4.35	-3.47	4.02	SCIN
gene	ENSG00000164975	chr9:15422702-15465951	8.020e-05	3.26	-3.9	NA	SNAPC3
gene	ENSG00000225556	chr1:151810295-151813033	1.030e-06	-4.47	2.45	NA	C2CD4D
gene	ENSG00000211818	chr14:22771939-22772438	9.210e-05	2.27	NA	-4.39	
gene	ENSG00000115486	chr2:85774743-85788670	3.620e-10	1.32	-3.78	2.23	GGCX
gene	ENSG00000175893	chr9:14615530-14693469	8.020e-05	2.22	-4.6	-4.61	ZDHHC21

Table B.5: Molecular mediators associated with Alzheimer’s disease (top 10 per QTL type, ranked by log odds)

QTL type	Mediator	Location (hg19)	GWAS p-value	lodds (mono)	lodds (neut)	lodds (fcel)	Closest coding gene
DNAme	cg05908241	chr7:143109367-143109367	1.417e-11	-1.34	8.44	-2.92	EPHA1
DNAme	cg01965533	chr14:75348459-75348459	1.004e-07	6.21	NA	-4.05	DLST
DNAme	cg07661177	chr2:127864026-127864026	1.001e-26	6.02	-1.5	-1.82	BIN1
DNAme	cg20406878	chr1:199996237-199996237	5.456e-05	6.02	NA	-4.37	NR5A2
DNAme	cg03290677	chr10:135122746-135122746	5.561e-05	5.68	0.357	-3.53	TUBGCP2, ZNF511
DNAme	cg16821175	chr3:128080818-128080818	2.735e-05	NA	-4.35	5.44	EEFSEC, MIR1280
DNAme	cg14961031	chr3:4344921-4344921	9.051e-05	NA	NA	5.41	SETMAR
DNAme	cg27056437	chr15:85177026-85177026	2.819e-06	-4.57	5.39	-4.45	SCAND2
DNAme	cg21446981	chr7:37534909-37534909	4.557e-06	1.94	5.3	-3.85	ELMO1
DNAme	cg09435798	chr10:135122752-135122752	5.561e-05	-4.56	-0.962	5.06	TUBGCP2, ZNF511
K27AC	chr11:85865916-85876777	chr11:85865916-85876777	6.534e-16	2.69	6.22	NA	EED
K27AC	chr8:103412080-103426942	chr8:103412080-103426942	7.073e-05	-3.66	5.54	-3.25	UBR5
K27AC	chr10:28747640-28750322	chr10:28747640-28750322	9.189e-05	-4.07	NA	5.2	WAC
K27AC	chr14:51629098-51632478	chr14:51629098-51632478	1.010e-05	4.71	-3.2	-4.43	TRIM9
K27AC	chr4:36289780-36291481	chr4:36289780-36291481	9.801e-05	NA	4.69	-2.64	DTHD1
K27AC	chr19:19667333-19674959	chr19:19667333-19674959	2.006e-05	NA	NA	4.26	CILP2
K27AC	chr2:127861936-127866241	chr2:127861936-127866241	1.001e-26	4.07	1.45	-3.69	BIN1
K27AC	chr15:85936196-85938677	chr15:85936196-85938677	1.684e-05	-1.86	-4.56	4	AKAP13
K27AC	chr5:178121304-178124171	chr5:178121304-178124171	2.815e-05	3.73	-2.1	-3.73	ZNF354A
K27AC	chr4:111084020-111089235	chr4:111084020-111089235	3.632e-05	NA	3.41	-4.43	ELOVL6
K4ME1	chr2:127814061-127870567	chr2:127814061-127870567	1.001e-26	1.09	5.07	-4.04	BIN1
K4ME1	chr1:200902206-200912056	chr1:200902206-200912056	5.456e-05	-3.41	4.88	NA	C1orf106
K4ME1	chr2:127884897-127896433	chr2:127884897-127896433	1.001e-26	4.76	-1.34	-3.8	BIN1
K4ME1	chr7:143052447-143144656	chr7:143052447-143144656	1.417e-11	4.75	-1.28	-4.44	EPHA1, FAM131B, TAS2R60, ZYX
K4ME1	chr16:11369336-11375392	chr16:11369336-11375392	5.752e-06	3.79	-3.25	-4.23	PRM1, PRM2
K4ME1	chr9:137531286-137537745	chr9:137531286-137537745	1.502e-05	-3.64	3.75	-4.5	COL5A1
K4ME1	chr5:110094501-110097674	chr5:110094501-110097674	9.804e-05	NA	3.52	NA	TMEM232
K4ME1	chr10:4811776-4815589	chr10:4811776-4815589	6.500e-05	NA	-4.04	3.45	AKR1E2
K4ME1	chr20:57261983-57285054	chr20:57261983-57285054	1.054e-06	-4.29	3.36	NA	NPEPL1
K4ME1	chr15:60262398-60278237	chr15:60262398-60278237	5.262e-06	-2.93	3.29	-3.28	FOXBP1
PSI	19:51728412:51729058:1:1:1_19:51728855:51729058:1:1:1_	chr19:51728412-51729058	5.120e-08	3.35	NA	NA	CD33
PSI	6:79679635:79679765:2:2:1_6:79679635:79680497:2:2:0_	chr6:79679635-79679765	5.608e-05	NA	3.15	NA	PHIP
PSI	4:109541757:109543106:1:1:1_4:109541869:109543106:1:1:1_	chr4:109541757-109543106	3.632e-05	NA	2.48	NA	
PSI	15:91428489:91428641:1:1:1_15:91428489:91430190:1:1:1_	chr15:91428489-91428641	2.731e-05	2.47	-4.21	NA	FES
PSI	1:228871752:228873419:1:1:1_1:228871752:228879031:1:1:0_	chr1:228871752-228873419	3.403e-05	2.42	NA	NA	RHOA
PSI	11:85726007:85742510:2:2:0_11:85733513:85742510:2:2:1_11:85737410:85742510:2:2:1_	chr11:85726007-85742510	6.534e-16	NA	2.15	NA	PICALM
PSI	13:108922773:108939150:1:1:1_13:108922773:108955600:1:1:1_	chr13:108922773-108939150	2.138e-05	2	-4.56	NA	TNFSF13B
PSI	7:138758777:138763297:2:2:0_7:138761156:138763297:2:2:1_	chr7:138758777-138763297	4.854e-05	-4.33	1.94	NA	ZC3HAV1
PSI	14:53184861:53185694:1:1:0_14:53185071:53185694:1:1:1_	chr14:53184861-53185694	1.010e-05	-4.6	1.88	NA	PSMC6
PSI	6:79655978:79656254:2:2:0_6:79655978:79656427:2:2:1_	chr6:79655978-79656254	5.608e-05	-4.91	1.55	-4.63	PHIP
gene	ENSG00000184967	chr12:132628993-132637013	2.456e-05	10.1	NA	NA	NOC4L
gene	ENSG00000136717	chr2:127805603-127864931	1.001e-26	8.51	-3.13	-3.47	BIN1
gene	ENSG00000251881	chr15:80434400-80434469	2.819e-06	5.52	-4.53	-3.72	
gene	ENSG00000122779	chr7:138144953-138274738	4.854e-05	4.33	NA	-0.718	TRIM24
gene	ENSG00000229153	chr7:143104906-143220542	1.417e-11	3.91	-2.08	-4.22	
gene	ENSG00000234066	chr7:143134127-143135066	1.417e-11	-3.63	3.19	-4.34	
gene	ENSG00000213558	chr9:110220319-110220587	6.884e-05	-4.62	NA	3.15	
gene	ENSG00000236152	chr3:6814724-6815030	8.535e-05	2.81	NA	-4.64	
gene	ENSG00000223040	chr3:2130617-2130895	2.932e-05	NA	2.7	-4.58	
gene	ENSG00000100485	chr14:50583847-50698276	1.010e-05	NA	-4.54	2.64	SOS2

Table B.6: Molecular mediators associated with systemic lupus erythematosus (top 10 per QTL type, ranked by log odds)

QTL type	Mediator	Location (hg19)	GWAS p-value	lodds (mono)	lodds (neut)	lodds (tcel)	Closest coding gene
DNAme	cg16287740	chr7:128595046-128595046	1.861e-45	14.3	-3.61	-4.13	TNPO3
DNAme	cg02846316	chr16:31340340-31340340	5.027e-48	13.1	13	8.26	ITGAM
DNAme	cg20802826	chr6:30782080-30782080	6.700e-103	13.1	-3.36	-2.81	DDR1
DNAme	cg07389699	chr6:32728786-32728786	6.700e-103	12.6	-1.59	0.548	HLA-DQB2
DNAme	cg00381684	chr6:31913814-31913814	6.700e-103	-2.96	-1.58	11.8	CFB
DNAme	cg16505946	chr6:31870305-31870305	6.700e-103	-1.08	-1.08	0.157	ZBTB12
DNAme	cg13914531	chr7:128579876-128579876	1.861e-45	-1.96	-4.35	11.7	IRF5
DNAme	cg03796735	chr6:32632643-32632643	6.700e-103	11	-0.333	-2.88	HLA-DQB1
DNAme	cg21612666	chr6:26537797-26537797	1.115e-34	-2.43	-1.34	10.9	HMGNA4
DNAme	cg04014328	chr17:46653615-46653615	1.223e-06	10.8	-2.64	-4.43	HOXB4
K27AC	chr6:31797205-31818044	chr6:31797205-31818044	6.700e-103	12.6	-3.81	-3.11	C6orf48
K27AC	chr6:32650849-32655154	chr6:32650849-32655154	6.700e-103	12.2	-3.84	-4.45	HLA-DQB1
K27AC	chr6:31457099-31467946	chr6:31457099-31467946	6.700e-103	-3.97	0.29	11.5	MICB
K27AC	chr6:31690633-31713468	chr6:31690633-31713468	6.700e-103	1.77	10.1	7.99	CLIC1, DDAH2, MSH5, MSH5-SAPCD1
K27AC	chr6:30225646-30229311	chr6:30225646-30229311	1.546e-96	-0.715	9.83	0.628	TRIM26
K27AC	chr6:30739568-30806629	chr6:30739568-30806629	6.700e-103	8.99	0.821	0.451	IER3
K27AC	chr6:31468114-31471249	chr6:31468114-31471249	6.700e-103	8.35	-1.18	-2.78	MICB
K27AC	chr6:26358262-26369051	chr6:26358262-26369051	1.115e-34	-2.21	-4.53	7.82	BTN3A2
K27AC	chr6:29973310-29975990	chr6:29973310-29975990	1.546e-96	-1.06	-1.39	7.63	ZNRD1
K27AC	chr3:187430320-187432323	chr3:187430320-187432323	4.267e-05	-4.53	NA	7.48	RTP2
K4ME1	chr6:31225358-31246870	chr6:31225358-31246870	6.700e-103	12.8	1.9	-1.27	HLA-C
K4ME1	chr6:30373698-30375735	chr6:30373698-30375735	1.546e-96	NA	NA	11.6	RPP21
K4ME1	chr6:31311305-31338443	chr6:31311305-31338443	6.700e-103	-2.67	-3.25	11.6	HLA-B
K4ME1	chr6:31073716-31086136	chr6:31073716-31086136	6.700e-103	11.3	NA	0.137	C6orf15, PSORS1C1
K4ME1	chr6:29745177-29749509	chr6:29745177-29749509	2.530e-65	-3.73	-4.13	10.9	HLA-G
K4ME1	chr6:32130018-32174407	chr6:32130018-32174407	6.700e-103	9.28	-4.24	-4.33	AGER, PAT1, EGFL8, GPSM3, PBX2, RNF5
K4ME1	chr6:32894459-32958593	chr6:32894459-32958593	6.700e-103	1.91	9.12	4.7	BRD2, HLA-DMA, HLA-DMB
K4ME1	chr6:32218325-32225919	chr6:32218325-32225919	6.700e-103	-2.69	-3.21	8.29	NOTCH4
K4ME1	chr6:32072640-32101980	chr6:32072640-32101980	6.700e-103	8.08	-3.1	-2.02	ATF6B, FKBPL, TNXB
K4ME1	chr6:31152700-31211658	chr6:31152700-31211658	6.700e-103	7.98	0.313	2.81	HCG27
PSI	6:31239646:31324862:2:2:0_6:31324735:31324862:2:2:1	chr6:31239646-31324862	6.700e-103	14.3	14.4	13.8	HLA-B, HLA-C
PSI	6:31238263:31238408:2:2:1_6:31238263:31238849:2:2:1	chr6:31238263-31238408	6.700e-103	NA	13.3	NA	HLA-C
PSI	6:29693084:29857483:1:1:0_6:29856002:29857483:1:1:0_6:29857382:29857483:1:1:1	chr6:29693084-29857483	1.546e-96	NA	13.1	NA	HLA-F, HLA-G
PSI	6:30708576:30708966:2:2:1_6:30708576:30709390:2:2:1	chr6:30708576-30708966	6.700e-103	NA	12.6	NA	FLOT1
PSI	6:32486444:32522466:2:2:0_6:32521744:32522466:2:2:1	chr6:32486444-32522466	6.700e-103	12.3	NA	NA	HLA-DRB5
PSI	12:120662181:120703419:2:2:1_12:120664921:120703419:2:2:1	chr12:120662181-120703419	5.867e-05	11.1	NA	-3.38	PXN
PSI	2:191514709:191523883:1:1:1_2:191520880:191523883:1:1:1	chr2:191514709-191523883	9.733e-66	8.06	NA	NA	NAB1
PSI	6:31237321:31237742:2:2:1_6:31237321:31237742:2:2:0	chr6:31237321-31237742	6.700e-103	-2.5	NA	6.94	HLA-C
PSI	16:31371984:31372100:1:1:0_16:31371984:31372383:1:1:0	chr16:31371984-31372100	5.027e-48	NA	6.85	NA	ITGAX
PSI	6:29857382:30459313:1:1:0_6:30229257:30459313:1:1:0_6:30459190:30459313:1:1:1	chr6:29857382-30459313	1.546e-96	2.83	6.21	NA	HLA-A, HLA-E, PPP1R11, RNF39, RPP21, TRIM10, TRIM15, TRIM26, TRIM31, TRIM39, TRIM39-RPP21, TRIM40, ZNRD1
gene	ENSG00000204338	chr6:31973413-31976228	6.700e-103	14	-0.434	0.457	
gene	ENSG00000137411	chr6:30876019-30894236	6.700e-103	0.832	-1.08	12.1	VAR52
gene	ENSG00000206341	chr6:29855350-29858259	1.546e-96	8.8	11.4	-2.28	
gene	ENSG00000224389	chr6:31982539-32003195	6.700e-103	-4.44	-1.9	11	C4B
gene	ENSG00000204421	chr6:31686425-31689622	6.700e-103	-4.48	10.3	-4.27	LY6G6C
gene	ENSG00000204257	chr6:32916390-32938493	6.700e-103	-2.6	0.063	9.84	HLA-DMA
gene	ENSG00000110921	chr12:110011060-110035067	3.507e-08	9.11	0.925	-4.46	MVK
gene	ENSG00000214894	chr6:30766431-30798436	6.700e-103	8.03	4.59	-0.316	
gene	ENSG00000206503	chr6:29909037-29913661	1.546e-96	-2.81	7.56	-2.48	HLA-A
gene	ENSG00000230174	chr6:31409444-31414750	6.700e-103	-4.5	7.44	-0.0125	

Bibliography

- [1] The Gene Ontology Resource: 20 years and still GOing strong. *Nucleic Acids Research*, 47(D1):D330–D338, January 2019. Publisher: Oxford Academic.
- [2] Picard toolkit. <http://broadinstitute.github.io/picard/>, 2019.
- [3] Karen Adelman and John T. Lis. Promoter-proximal pausing of RNA polymerase II: emerging roles in metazoans. *Nature reviews. Genetics*, 13(10):720–731, October 2012.
- [4] Carles B. Adsera, Yongjin P. Park, Wouter Meuleman, and Manolis Kellis. Integrative analysis of 10,000 epigenomic maps across 800 samples for regulatory genomics and disease dissection. *bioRxiv*, page 810291, October 2019. Publisher: Cold Spring Harbor Laboratory Section: New Results.
- [5] Michael Ashburner, Catherine A. Ball, Judith A. Blake, David Botstein, Heather Butler, J. Michael Cherry, Allan P. Davis, Kara Dolinski, Selina S. Dwight, Janan T. Eppig, Midori A. Harris, David P. Hill, Laurie Issel-Tarver, Andrew Kasarskis, Suzanna Lewis, John C. Matese, Joel E. Richardson, Martin Ringwald, Gerald M. Rubin, and Gavin Sherlock. Gene Ontology: tool for the unification of biology. *Nature Genetics*, 25(1):25–29, May 2000. Number: 1 Publisher: Nature Publishing Group.
- [6] Patrik Asp, Roy Blum, Vasupradha Vethantham, Fabio Parisi, Mariann Micsinai, Jemmie Cheng, Christopher Bowman, Yuval Kluger, and Brian David Dynlacht. Genome-wide remodeling of the epigenetic landscape during myogenic differentiation. *Proceedings of the National Academy of Sciences*, 108(22):E149–E158, May 2011. Publisher: National Academy of Sciences Section: PNAS Plus.
- [7] Christian Baarlink, Matthias Plessner, Alice Sherrard, Kohtaro Morita, Shinji Misu, David Virant, Eva-Maria Kleinschnitz, Robert Harniman, Dominic Alibhai, Stefan Baumeister, Kei Miyamoto, Ulrike Endesfelder, Abderrahmane Kaidi, and Robert Grosse. A transient pool of nuclear F-actin at mitotic exit controls chromatin organization. *Nature Cell Biology*, 19(12):1389–1399, December 2017. Number: 12 Publisher: Nature Publishing Group.
- [8] Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1):1–48, 2015.

- [9] James Bentham, David L. Morris, Deborah S. Cunninghame Graham, Christopher L. Pinder, Philip Tomblason, Timothy W. Behrens, Javier Martínez, Benjamin P. Fairfax, Julian C. Knight, Lingyan Chen, Joseph Replogle, Ann-Christine Syvärdén, Lars Rönnblom, Robert R. Graham, Joan E. Wither, John D. Rioux, Marta E. Alarcón-Riquelme, and Timothy J. Vyse. Genetic association analyses implicate aberrant regulation of innate and adaptive immunity genes in the pathogenesis of systemic lupus erythematosus. *Nature Genetics*, 47(12):1457–1464, December 2015. Number: 12 Publisher: Nature Publishing Group.
- [10] Tomaz Berisa and Joseph K. Pickrell. Approximately independent linkage disequilibrium blocks in human populations. *Bioinformatics*, 32(2):283–285, January 2016.
- [11] Bradley E. Bernstein, Alexander Meissner, and Eric S. Lander. The Mammalian Epigenome. *Cell*, 128(4):669–681, February 2007.
- [12] Bradley E. Bernstein, Tarjei S. Mikkelsen, Xiaohui Xie, Michael Kamal, Dana J. Huebert, James Cuff, Ben Fry, Alex Meissner, Marius Wernig, Kathrin Plath, Rudolf Jaenisch, Alexandre Wagschal, Robert Feil, Stuart L. Schreiber, and Eric S. Lander. A Bivalent Chromatin Structure Marks Key Developmental Genes in Embryonic Stem Cells. *Cell*, 125(2):315–326, April 2006.
- [13] Christoph Bock, Isabel Beerman, Wen-Hui Lien, Zachary D. Smith, Hongcang Gu, Patrick Boyle, Andreas Gnirke, Elaine Fuchs, Derrick J. Rossi, and Alexander Meissner. DNA Methylation Dynamics during In Vivo Differentiation of Blood and Skin Stem Cells. *Molecular Cell*, 47(4):633–647, August 2012.
- [14] Alistair N. Boettiger, Bogdan Bintu, Jeffrey R. Moffitt, Siyuan Wang, Brian J. Believeau, Geoffrey Fudenberg, Maxim Imakaev, Leonid A. Mirny, Chao-ting Wu, and Xiaowei Zhuang. Super-resolution imaging reveals distinct chromatin folding for different epigenetic states. *Nature*, 529(7586):418–422, January 2016. Number: 7586 Publisher: Nature Publishing Group.
- [15] Alan P. Boyle, Sean Davis, Hennady P. Shulha, Paul Meltzer, Elliott H. Margulies, Zhiping Weng, Terrence S. Furey, and Gregory E. Crawford. High-Resolution Mapping and Characterization of Open Chromatin across the Genome. *Cell*, 132(2):311–322, January 2008.
- [16] Jason D. Buenrostro, Paul G. Giresi, Lisa C. Zaba, Howard Y. Chang, and William J. Greenleaf. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nature Methods*, 10(12):1213–1218, December 2013.
- [17] Hong-Thuy Bui, Sayaka Wakayama, Satoshi Kishigami, Keun-Kyu Park, Jin-Hoi Kim, Nguyen Van Thuan, and Teruhiko Wakayama. Effect of Trichostatin

A on Chromatin Remodeling, Histone Modifications, DNA Replication, and Transcriptional Activity in Cloned Mouse Embryos. *Biology of Reproduction*, 83(3):454–463, September 2010. Publisher: Oxford Academic.

- [18] Peter Carbonetto and Matthew Stephens. Scalable Variational Inference for Bayesian Variable Selection in Regression, and Its Accuracy in Genetic Association Studies. *Bayesian Analysis*, 7(1):73–108, March 2012. Publisher: International Society for Bayesian Analysis.
- [19] Sihem Cheloufi, Ulrich Elling, Barbara Hopfgartner, Youngsook L. Jung, Jernej Murn, Maria Ninova, Maria Hubmann, Aimee I. Badeaux, Cheen Euong Ang, Danielle Tenen, Daniel J. Wesche, Nadezhda Abazova, Max Hogue, Nilgun Tasdemir, Justin Brumbaugh, Philipp Rathert, Julian Jude, Francesco Ferrari, Andres Blanco, Michaela Fellner, Daniel Wenzel, Marietta Zinner, Simon E. Vidal, Oliver Bell, Matthias Stadtfeld, Howard Y. Chang, Genevieve Almouzni, Scott W. Lowe, John Rinn, Marius Wernig, Alexei Aravin, Yang Shi, Peter J. Park, Josef M. Penninger, Johannes Zuber, and Konrad Hochedlinger. The histone chaperone CAF-1 safeguards somatic cell identity. *Nature*, 528(7581):218–224, December 2015.
- [20] Lu Chen, Bing Ge, Francesco Paolo Casale, Louella Vasquez, Tony Kwan, Diego Garrido-Martn, Stephen Watt, Ying Yan, Kousik Kundu, Simone Ecker, Avik Datta, David Richardson, Frances Burden, Daniel Mead, Alice L. Mann, Jose Maria Fernandez, Sophia Rowston, Steven P. Wilder, Samantha Farrow, Xiaojian Shao, John J. Lambourne, Adriana Redensek, Cornelis A. Albers, Vyacheslav Amstislavskiy, Sofie Ashford, Kim Berentsen, Lorenzo Bomba, Guillaume Bourque, David Bujold, Stephan Busche, Maxime Caron, Shu-Huang Chen, Warren Cheung, Oliver Delaneau, Emmanouil T. Dermitzakis, Heather Elding, Irina Colgiu, Frederik O. Bagger, Paul Flicek, Ehsan Habibi, Valentina Iotchkova, Eva Janssen-Megens, Bowon Kim, Hans Lehrach, Ernesto Lowy, Amit Mandoli, Filomena Matarese, Matthew T. Maurano, John A. Morris, Vera Pancaldi, Farzin Pourfarzad, Karola Rehnstrom, Augusto Rendon, Thomas Risch, Nilofar Sharifi, Marie-Michelle Simon, Marc Sultan, Alfonso Valencia, Klaudia Walter, Shuang-Yin Wang, Mattia Frontini, Stylianos E. Antonarakis, Laura Clarke, Marie-Laure Yaspo, Stephan Beck, Roderic Guigo, Daniel Rico, Joost H. A. Martens, Willem H. Ouwehand, Taco W. Kuijpers, Dirk S. Paul, Hendrik G. Stunnenberg, Oliver Stegle, Kate Downes, Tomi Pastinen, and Nicole Soranzo. Genetic Drivers of Epigenetic and Transcriptional Variation in Human Immune Cells. *Cell*, 167(5):1398–1414.e24, November 2016.
- [21] The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature*, 489(7414):57–74, September 2012.
- [22] M. Ryan Corces, Jason D. Buenrostro, Beijing Wu, Peyton G. Greenside, Steven M. Chan, Julie L. Koenig, Michael P. Snyder, Jonathan K. Pritchard, Anshul Kundaje, William J. Greenleaf, Ravindra Majeti, and Howard Y.

- Chang. Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. *Nature Genetics*, 48(10):1193–1203, October 2016. Number: 10 Publisher: Nature Publishing Group.
- [23] George Davey Smith and Shah Ebrahim. Mendelian randomization: can genetic epidemiology contribute to understanding environmental determinants of disease? *International Journal of Epidemiology*, 32(1):1–22, February 2003. Publisher: Oxford Academic.
- [24] Laurence Delacroix, Emmanuel Moutier, Gioia Altobelli, Stephanie Legras, Olivier Poch, Mohamed-Amin Choukrallah, Isabelle Bertin, Bernard Jost, and Irwin Davidson. Cell-Specific Interaction of Retinoic Acid Receptors with Target Genes in Mouse Embryonic Fibroblasts and Embryonic Stem Cells. *Molecular and Cellular Biology*, 30(1):231–244, January 2010. Publisher: American Society for Microbiology Journals Section: ARTICLES.
- [25] Jason Ernst and Manolis Kellis. Discovery and characterization of chromatin states for systematic annotation of the human genome. *Nature Biotechnology*, 28(8):817–825, August 2010. Number: 8 Publisher: Nature Publishing Group.
- [26] Eric R. Gamazon, Heather E. Wheeler, Kanaan P. Shah, Sahar V. Mozaffari, Keston Aquino-Michaels, Robert J. Carroll, Anne E. Eyler, Joshua C. Denny, GTEx Consortium, Dan L. Nicolae, Nancy J. Cox, and Hae Kyung Im. A gene-based association method for mapping traits using reference transcriptome data. *Nature Genetics*, 47(9):1091–1098, September 2015.
- [27] Claudia Giambartolomei, Damjan Vukcevic, Eric E. Schadt, Lude Franke, Aroon D. Hingorani, Chris Wallace, and Vincent Plagnol. Bayesian Test for Colocalisation between Pairs of Genetic Association Studies Using Summary Statistics. *PLoS Genetics*, 10(5), May 2014.
- [28] Matthew G. Guenther, Stuart S. Levine, Laurie A. Boyer, Rudolf Jaenisch, and Richard A. Young. A Chromatin Landmark and Transcription Initiation at Most Promoters in Human Cells. *Cell*, 130(1):77–88, July 2007.
- [29] Alexander Gusev, Arthur Ko, Huwenbo Shi, Gaurav Bhatia, Wonil Chung, Brenda W. J. H. Penninx, Rick Jansen, Eco J. C. de Geus, Dorret I. Boomsma, Fred A. Wright, Patrick F. Sullivan, Elina Nikkola, Marcus Alvarez, Mete Civelek, Aldons J. Lusi, Terho Lehtimäki, Emma Raitoharju, Mika Kähönen, Ilkka Seppälä, Olli T. Raitakari, Johanna Kuusisto, Markku Laakso, Alkes L. Price, Pami Pajukanta, and Bogdan Pasaniuc. Integrative approaches for large-scale transcriptome-wide association studies. *Nature Genetics*, 48(3):245–252, March 2016.
- [30] Eilis Hannon, Mike Weedon, Nicholas Bray, Michael O'Donovan, and Jonathan Mill. Pleiotropic Effects of Trait-Associated Genetic Variation on DNA Methylation: Utility for Refining GWAS Loci. *The American Journal of Human Genetics*, 100(6):954–959, June 2017.

- [31] R. David Hawkins, Gary C. Hon, Leonard K. Lee, QueMinh Ngo, Ryan Lister, Mattia Pelizzola, Lee E. Edsall, Samantha Kuan, Ying Luu, Sarit Klugman, Jessica Antosiewicz-Bourget, Zhen Ye, Celso Espinoza, Saurabh Agarwahl, Li Shen, Victor Ruotti, Wei Wang, Ron Stewart, James A. Thomson, Joseph R. Ecker, and Bing Ren. Distinct Epigenomic Landscapes of Pluripotent and Lineage-Committed Human Cells. *Cell Stem Cell*, 6(5):479–491, May 2010.
- [32] Nathaniel D. Heintzman, Rhona K. Stuart, Gary Hon, Yutao Fu, Christina W. Ching, R. David Hawkins, Leah O. Barrera, Sara Van Calcar, Chunxu Qu, Keith A. Ching, Wei Wang, Zhiping Weng, Roland D. Green, Gregory E. Crawford, and Bing Ren. Distinct and predictive chromatin signatures of transcriptional promoters and enhancers in the human genome. *Nature Genetics*, 39(3):311–318, March 2007.
- [33] Sven Heinz, Christopher Benner, Nathanael Spann, Eric Bertolino, Yin C. Lin, Peter Laslo, Jason X. Cheng, Cornelis Murre, Harinder Singh, and Christopher K. Glass. Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. *Molecular Cell*, 38(4):576–589, May 2010.
- [34] Konrad Hochedlinger and Rudolf Jaenisch. Nuclear reprogramming and pluripotency. *Nature*, 441(7097):1061–1067, June 2006. Number: 7097 Publisher: Nature Publishing Group.
- [35] Kyung-Won Hong, Ji-Eun Lim, and Bermseok Oh. A regulatory SNP in AKAP13 is associated with blood pressure in Koreans. *Journal of Human Genetics*, 56(3):205–210, March 2011. Number: 3 Publisher: Nature Publishing Group.
- [36] Hailiang Huang, Ming Fang, Luke Jostins, MaÅaa UmiÄĖviÄĖ Mirkov, Gabrielle Boucher, Carl A. Anderson, Vibeke Andersen, Isabelle Cleynen, Adrian Cortes, FranÄĖois Crins, Mauro DÄĖAmato, ValÄĖrie Deffontaine, Julia Dmitrieva, Elisa Docampo, Mahmoud Elansary, Kyle Kai-How Farh, Andre Franke, Ann-Stephan Gori, Philippe Goyette, Jonas Halfvarson, Talin Haritunians, Jo Knight, Ian C. Lawrance, Charlie W. Lees, Edouard Louis, Rob Mariman, Theo Meuwissen, Myriam Mni, Yukihide Momozawa, Miles Parkes, Sarah L. Spain, Emilie ThÄĖÄĖtre, Gosia Trynka, Jack Satsangi, Suzanne van Sommeren, Severine Vermeire, Ramnik J. Xavier, Rinse K. Weersma, Richard H. Duerr, Christopher G. Mathew, John D. Rioux, Dermot P. B. McGovern, Judy H. Cho, Michel Georges, Mark J. Daly, and Jeffrey C. Barrett. Fine-mapping inflammatory bowel disease loci to single-variant resolution. *Nature*, 547(7662):173–178, July 2017. Number: 7662 Publisher: Nature Publishing Group.
- [37] Eva HÄĖrmanseder, Angela Simeone, George E. Allen, Charles R. Bradshaw, Magdalena FiglmÄĖjller, John Gurdon, and Jerome Jullien. H3K4 Methylation-

Dependent Memory of Somatic Cell Identity Inhibits Reprogramming and Development of Nuclear Transfer Embryos. *Cell Stem Cell*, 21(1):135–143.e6, July 2017.

- [38] Daishi Inoue, Hitoshi Aihara, Tatsuharu Sato, Hirofumi Mizusaki, Masamichi Doiguchi, Miki Higashi, Yuko Imamura, Mitsuhiro Yoneda, Takayuki Miyanishi, Satoshi Fujii, Akihiko Okuda, Takeya Nakagawa, and Takashi Ito. Dzip3 regulates developmental genes in mouse embryonic stem cells by reorganizing 3D chromatin conformation. *Scientific Reports*, 5(1):1–11, November 2015. Number: 1 Publisher: Nature Publishing Group.
- [39] International Multiple Sclerosis Genetics Consortium, Wellcome Trust Case Control Consortium 2, Stephen Sawcer, Garrett Hellenthal, Matti Pirinen, Chris C. A. Spencer, Nikolaos A. Patsopoulos, Loukas Moutsianas, Alexander Dilthey, Zhan Su, Colin Freeman, Sarah E. Hunt, Sarah Edkins, Emma Gray, David R. Booth, Simon C. Potter, An Goris, Gavin Band, Annette Bang Oturai, Amy Strange, Janna Saarela, CÃline Bellenguez, Bertrand Fontaine, Matthew Gillman, Bernhard Hemmer, Rhian Gwilliam, Frauke Zipp, Alagurevathi Jayakumar, Roland Martin, Stephen Leslie, Stanley Hawkins, Eleni Giannoulatou, Sandra D’alfonso, Hannah Blackburn, Filippo Martinelli Boneschi, Jennifer Liddle, Hanne F. Harbo, Marc L. Perez, Anne Spurkland, Matthew J. Waller, Marcin P. Mycko, Michelle Ricketts, Manuel Comabella, Naomi Hammond, Ingrid Kockum, Owen T. McCann, Maria Ban, Pamela Whittaker, Anu Kempainen, Paul Weston, Clive Hawkins, Sara Widaa, John Zajicek, Serge Dronov, Neil Robertson, Suzannah J. Bumpstead, Lisa F. Barcellos, Rathi Ravindrarajah, Roby Abraham, Lars Alfredsson, Kristin Ardlie, Cristin Aubin, Amie Baker, Katharine Baker, Sergio E. Baranzini, Laura Bergamaschi, Roberto Bergamaschi, Allan Bernstein, Achim Berthele, Mike Boggild, Jonathan P. Bradfield, David Brassat, Simon A. Broadley, Dorothea Buck, Helmut Butzkueven, Ruggero Capra, William M. Carroll, Paola Cavalla, Elisabeth G. Celius, Sabine Cepok, Rosetta Chiavacci, FranÃgoise Clerget-Darpoux, Kathleen Clysters, Giancarlo Comi, Mark Cossburn, Isabelle Cournu-Rebeix, Mathew B. Cox, Wendy Cozen, Bruce A. C. Cree, Anne H. Cross, Daniele Cusi, Mark J. Daly, Emma Davis, Paul I. W. de Bakker, Marc Debouverie, Marie Beatrice D’hooghe, Katherine Dixon, Rita Dobosi, BÃnÃdicte Dubois, David Ellinghaus, Irina Elovaara, Federica Esposito, Claire Fontenille, Simon Foote, Andre Franke, Daniela Galimberti, Angelo Ghezzi, Joseph Glessner, Refujia Gomez, Olivier Gout, Colin Graham, Struan F. A. Grant, Franca Rosa Guerini, Hakon Hakonarson, Per Hall, Anders Hamsten, Hans-Peter Hartung, Rob N. Heard, Simon Heath, Jeremy Hobart, Muna Hoshi, Carmen Infante-Duarte, Gillian Ingram, Wendy Ingram, Talat Islam, Maja Jagodic, Michael Kabesch, Allan G. Kermod, Trevor J. Kilpatrick, Cecilia Kim, Norman Klopp, Keijo Koivisto, Malin Larsson, Mark Lathrop, Jeannette S. Lechner-Scott, Maurizio A. Leone, Virpi LeppÃd, Ulrika Liljedahl, Izaura Lima Bomfim, Robin R. Lincoln, Jenny Link, Jianjun Liu, Aslaug R. Lorentzen, Sara Lupoli, Fabio

Macciardi, Thomas Mack, Mark Marriott, Vittorio Martinelli, Deborah Mason, Jacob L. McCauley, Frank Mentch, Inger-Lise Mero, Tania Mihalova, Xavier Montalban, John Mottershead, Kjell-Morten Myhr, Paola Naldi, William Ollier, Alison Page, Aarno Palotie, Jean Pelletier, Laura Piccio, Trevor Pickersgill, Fredrik Piehl, Susan Pobywajlo, Hong L. Quach, Patricia P. Ramsay, Mauri Reunanen, Richard Reynolds, John D. Rioux, Mariaemma Rodegher, Sabine Roesner, Justin P. Rubio, Ina-Maria RÃijckert, Marco Salvetti, Erika Salvi, Adam Santaniello, Catherine A. Schaefer, Stefan Schreiber, Christian Schulze, Rodney J. Scott, Finn Sellebjerg, Krzysztof W. Selmaj, David Sexton, Ling Shen, Brigid Simms-Acuna, Sheila Skidmore, Patrick M. A. Sleiman, Cathrine Smestad, Per Soelberg SÃyrensen, Helle Bach SÃyndergaard, Jim Stankovich, Richard C. Strange, Anna-Maija Sulonen, Emilie Sundqvist, Ann-Christine SyvÃdnen, Francesca Taddeo, Bruce Taylor, Jenefer M. Blackwell, Pentti Tienari, Elvira Bramon, Ayman Tourbah, Matthew A. Brown, Ewa Tronczynska, Juan P. Casas, Niall Tubridy, Aiden Corvin, Jane Vickery, Janusz Jankowski, Pablo Villoslada, Hugh S. Markus, Kai Wang, Christopher G. Mathew, James Wason, Colin N. A. Palmer, H.-Erich Wichmann, Robert Plomin, Ernest Willoughby, Anna Rautanen, Juliane Winkelmann, Michael Wittig, Richard C. Trembath, Jacqueline Yaouanq, Ananth C. Viswanathan, Haitao Zhang, Nicholas W. Wood, Rebecca Zuvich, Panos Deloukas, Cordelia Langford, Audrey Duncanson, Jorge R. Oksenberg, Margaret A. Pericak-Vance, Jonathan L. Haines, Tomas Olsson, Jan Hillert, Adrian J. Ivinson, Philip L. De Jager, Leena Peltonen, Graeme J. Stewart, David A. Hafler, Stephen L. Hauser, Gil McVean, Peter Donnelly, and Alastair Compston. Genetic risk and a primary role for cell-mediated immune mechanisms in multiple sclerosis. *Nature*, 476(7359):214–219, August 2011.

- [40] Thomas Jenuwein and C. David Allis. Translating the Histone Code. *Science*, 293(5532):1074–1080, August 2001. Publisher: American Association for the Advancement of Science Section: Special Reviews.
- [41] Keven R. Johnson, Jessie Nicodemus-Johnson, Mathew J. Spindler, and Graeme K. Carnegie. Genome-Wide Gene Expression Analysis Shows AKAP13-Mediated PKD1 Signaling Regulates the Transcriptional Response to Cardiac Hypertrophy. *PLoS ONE*, 10(7), July 2015.
- [42] Jerome Jullien, Kei Miyamoto, Vincent Pasque, GeorgeÃE. Allen, CharlesÃR. Bradshaw, NigelÃJ. Garrett, RichardÃP. Halley-Stott, Hiroshi Kimura, Keita Ohsumi, and JohnÃB. Gurdon. Hierarchical Molecular Events Driven by Oocyte-Specific Factors Lead to Rapid and Extensive Reprogramming. *Molecular Cell*, 55(4):524–536, August 2014.
- [43] Jerome Jullien, Vincent Pasque, Richard P. Halley-Stott, Kei Miyamoto, and J. B. Gurdon. Mechanisms of nuclear reprogramming by eggs and oocytes: a deterministic process? *Nature reviews. Molecular cell biology*, 12(7):453–459, June 2011.

- [44] Jerome Jullien, Munender Vodnala, Vincent Pasque, Mami Oikawa, Kei Miyamoto, George Allen, Sarah Anne David, Vincent Brochard, Stan Wang, Charles Bradshaw, Haruhiko Koseki, Vittorio Sartorelli, Nathalie Beaujean, and John Gurdon. Gene Resistance to Transcriptional Reprogramming following Nuclear Transfer Is Directly Mediated by Multiple Chromatin-Repressive Pathways. *Molecular Cell*, 65(5):873–884.e8, March 2017.
- [45] Minoru Kanehisa and Susumu Goto. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Research*, 28(1):27–30, January 2000. Publisher: Oxford Academic.
- [46] Donna Karolchik, Angela S. Hinrichs, Terrence S. Furey, Krishna M. Roskin, Charles W. Sugnet, David Haussler, and W. James Kent. The UCSC Table Browser data retrieval tool. *Nucleic Acids Research*, 32(suppl_1):D493–D496, January 2004.
- [47] W. James Kent, Charles W. Sugnet, Terrence S. Furey, Krishna M. Roskin, Tom H. Pringle, Alan M. Zahler, and David Haussler. The Human Genome Browser at UCSC. *Genome Research*, 12(6):996–1006, June 2002.
- [48] Sandy L. Klemm, Zohar Shipony, and William J. Greenleaf. Chromatin accessibility and the regulatory epigenome. *Nature Reviews Genetics*, 20(4):207–220, April 2019. Number: 4 Publisher: Nature Publishing Group.
- [49] Qing-Ran Kong, Bing-Teng Xie, Heng Zhang, Jing-Yu Li, Tian-Qing Huang, Ren-Yue Wei, and Zhong-Hua Liu. RE1-silencing Transcription Factor (REST) Is Required for Nuclear Reprogramming by Inhibiting Transforming Growth Factor β Signaling Pathway. *Journal of Biological Chemistry*, 291(53):27334–27342, December 2016. Publisher: American Society for Biochemistry and Molecular Biology.
- [50] J. C. Lambert, C. A. Ibrahim-Verbaas, D. Harold, A. C. Naj, R. Sims, C. Bellenguez, A. L. DeStafano, J. C. Bis, G. W. Beecham, B. Grenier-Boley, G. Russo, T. A. Thornton-Wells, N. Jones, A. V. Smith, V. Chouraki, C. Thomas, M. A. Ikram, D. Zelenika, B. N. Vardarajan, Y. Kamatani, C. F. Lin, A. Gerish, H. Schmidt, B. Kunkle, M. L. Dunstan, A. Ruiz, M. T. Bihoreau, S. H. Choi, C. Reitz, F. Pasquier, C. Cruchaga, D. Craig, N. Amin, C. Berr, O. L. Lopez, P. L. De Jager, V. Deramecourt, J. A. Johnston, D. Evans, S. Lovestone, L. Letenneur, F. J. Morais, D. C. Rubinsztein, G. Eiriksdottir, K. Sleegers, A. M. Goate, N. F. Aliev, M. W. Huentelman, M. Gill, K. Brown, M. I. Kamboh, L. Keller, P. Barberger-Gateau, B. McGuinness, E. B. Larson, R. Green, A. J. Myers, C. Dufouil, S. Todd, D. Wallon, S. Love, E. Rogaeva, J. Gallacher, P. St George-Hyslop, J. Clarimon, A. Lleo, A. Bayer, D. W. Tsuang, L. Yu, M. Tsolaki, P. Bossa, G. Spalletta, P. Proitsi, J. Collinge, S. Sorbi, F. Sanchez-Garcia, N. C. Fox, J. Hardy, M. C. Deniz Naranjo, P. Bosco, R. Clarke, C. Brayne, D. Galimberti, M. Mancuso, F. Matthews, European Alzheimer’s

- Disease Initiative (EADI), Genetic and Environmental Risk in Alzheimer’s Disease, Alzheimer’s Disease Genetic Consortium, Cohorts for Heart and Aging Research in Genomic Epidemiology, S. Moebus, P. Mecocci, M. Del Zompo, W. Maier, H. Hampel, A. Pilotto, M. Bullido, F. Panza, P. Caffarra, B. Nacmias, J. R. Gilbert, M. Mayhaus, L. Lannefelt, H. Hakonarson, S. Pichler, M. M. Carrasquillo, M. Ingelsson, D. Beekly, V. Alvarez, F. Zou, O. Valladares, S. G. Younkin, E. Coto, K. L. Hamilton-Nelson, W. Gu, C. Razquin, P. Pastor, I. Mateo, M. J. Owen, K. M. Faber, P. V. Jonsson, O. Combarros, M. C. O’Donovan, L. B. Cantwell, H. Soininen, D. Blacker, S. Mead, T. H. Mosley, D. A. Bennett, T. B. Harris, L. Fratiglioni, C. Holmes, R. F. de Bruijn, P. Passmore, T. J. Montine, K. Bettens, J. I. Rotter, A. Brice, K. Morgan, T. M. Foroud, W. A. Kukull, D. Hannequin, J. F. Powell, M. A. Nalls, K. Ritchie, K. L. Lunetta, J. S. Kauwe, E. Boerwinkle, M. Riemenschneider, M. Boada, M. Hiltunen, E. R. Martin, R. Schmidt, D. Rujescu, L. S. Wang, J. F. Dartigues, R. Mayeux, C. Tzourio, A. Hofman, M. M. NÃthen, C. Graff, B. M. Psaty, L. Jones, J. L. Haines, P. A. Holmans, M. Lathrop, M. A. Pericak-Vance, L. J. Launer, L. A. Farrer, C. M. van Duijn, C. Van Broeckhoven, V. Moskvina, S. Seshadri, J. Williams, G. D. Schellenberg, and P. Amouyel. Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer’s disease. *Nature Genetics*, 45(12):1452–1458, December 2013.
- [51] Ben Langmead and Steven L. Salzberg. Fast gapped-read alignment with Bowtie 2. *Nature Methods*, 9(4):357–359, April 2012.
- [52] Dongwei Li, Jing Liu, Xuejie Yang, Chunhua Zhou, Jing Guo, Chuman Wu, Yue Qin, Lin Guo, Jiangping He, Shenyong Yu, He Liu, Xiaoshan Wang, Fang Wu, Junqi Kuang, Andrew P. Hutchins, Jiekai Chen, and Duanqing Pei. Chromatin Accessibility Dynamics during iPSC Reprogramming. *Cell Stem Cell*, 21(6):819–833.e6, December 2017.
- [53] Arthur Liberzon, Chet Birger, Helga ThorvaldsdÃsttir, Mahmoud Ghandi, Jil-lÃP. Mesirov, and Pablo Tamayo. The Molecular Signatures Database Hallmark Gene Set Collection. *Cell Systems*, 1(6):417–425, December 2015.
- [54] Jimmy Z. Liu, Suzanne van Sommeren, Hailiang Huang, Siew C. Ng, Rudi Alberts, Atsushi Takahashi, Stephan Ripke, James C. Lee, Luke Jostins, Tejas Shah, Shifteh Abedian, Jae Hee Cheon, Judy Cho, Naser E. Dayani, Lude Franke, Yuta Fuyuno, Ailsa Hart, Ramesh C. Juyal, Garima Juyal, Won Ho Kim, Andrew P. Morris, Hossein Poustchi, William G. Newman, Vandana Midha, Timothy R. Orchard, Homayon Vahedi, Ajit Sood, Joseph Y. Sung, Reza Malekzadeh, Harm-Jan Westra, Keiko Yamazaki, Suk-Kyun Yang, International Multiple Sclerosis Genetics Consortium, International IBD Genetics Consortium, Jeffrey C. Barrett, Behrooz Z. Alizadeh, Miles Parkes, Thelma Bk, Mark J. Daly, Michiaki Kubo, Carl A. Anderson, and Rinse K. Weersma. Association analyses identify 38 susceptibility loci for inflammatory bowel dis-

- ease and highlight shared genetic risk across populations. *Nature Genetics*, 47(9):979–986, September 2015.
- [55] Ruize Liu, Juha Karjalainen, Andrea Byrnes, Beryl B. Cummings, Padhraig Gormley, Yunfeng Ruan, Yongjin Park, Lei Hou, Khoi T. Nguyen, Ramnik Xavier, Mark J. Daly, and Hailiang Huang. Junction-skipping regulation in complex disease. *bioRxiv*, page 804039, October 2019. Publisher: Cold Spring Harbor Laboratory Section: New Results.
- [56] Falong Lu, Yuting Liu, Azusa Inoue, Tsukasa Suzuki, Keji Zhao, and Yi Zhang. Establishing Chromatin Regulatory Landscape during Mouse Preimplantation Development. *Cell*, 165(6):1375–1388, June 2016.
- [57] Liang Ma, Peilin Jia, and Zhongming Zhao. Splicing QTL of human adipose-related traits. *Scientific Reports*, 8(1):1–5, January 2018. Number: 1 Publisher: Nature Publishing Group.
- [58] Walid E. Maalouf, Zichuan Liu, Vincent Brochard, Jean-Paul Renard, Pascale Debey, Nathalie Beaujean, and Daniele Zink. Trichostatin A treatment of cloned mouse embryos improves constitutive heterochromatin remodeling as well as developmental potential to term. *BMC Developmental Biology*, 9(1):11, February 2009.
- [59] Chantal M. Mayers, Jennifer Wadell, Kate McLean, Monica Venere, Minnie Malik, Takahisa Shibata, Paul H. Driggers, Tomoshige Kino, X. Catherine Guo, Hisashi Koide, Marat Gorivodsky, Alex Grinberg, Mahua Mukhopadhyay, Mones Abu-Asab, Heiner Westphal, and James H. Segars. The Rho Guanine Nucleotide Exchange Factor AKAP13 (BRX) Is Essential for Cardiac Development in Mice. *Journal of Biological Chemistry*, 285(16):12344–12354, April 2010. Publisher: American Society for Biochemistry and Molecular Biology.
- [60] Cory Y. McLean, Dave Bristor, Michael Hiller, Shoa L. Clarke, Bruce T. Schaar, Craig B. Lowe, Aaron M. Wenger, and Gill Bejerano. GREAT improves functional interpretation of cis-regulatory regions. *Nature Biotechnology*, 28(5):495–501, May 2010.
- [61] Gil A. McVean, David M. Altshuler (Co-Chair), Richard M. Durbin (Co-Chair), Gonçalo R. Abecasis, David R. Bentley, Aravinda Chakravarti, Andrew G. Clark, Peter Donnelly, Evan E. Eichler, Paul Flicek, Stacey B. Gabriel, Richard A. Gibbs, Eric D. Green, Matthew E. Hurles, Bartha M. Knoppers, Jan O. Korbel, Eric S. Lander, Charles Lee, Hans Lehrach, Elaine R. Mardis, Gabor T. Marth, Gil A. McVean, Deborah A. Nickerson, Jeanette P. Schmidt, Stephen T. Sherry, Jun Wang, Richard K. Wilson, Richard A. Gibbs (Principal Investigator), Huyen Dinh, Christie Kovar, Sandra Lee, Lora Lewis, Donna Muzny, Jeff Reid, Min Wang, Jun Wang (Principal Investigator), Xiaodong Fang, Xiaosen Guo, Min Jian, Hui Jiang, Xin Jin, Guoqing Li, Jingxiang Li, Yingrui Li, Zhuo Li, Xiao Liu, Yao Lu, Xuedi Ma,

Zhe Su, Shuaishuai Tai, Meifang Tang, Bo Wang, Guangbiao Wang, Honglong Wu, Renhua Wu, Ye Yin, Wenwei Zhang, Jiao Zhao, Meiru Zhao, Xiaole Zheng, Yan Zhou, Eric S. Lander (Principal Investigator), David M. Altshuler, Stacey B. Gabriel (Co-Chair), Namrata Gupta, Paul Flicek (Principal Investigator), Laura Clarke, Rasko Leinonen, Richard E. Smith, Xiangqun Zheng-Bradley, David R. Bentley (Principal Investigator), Russell Grocock, Sean Humphray, Terena James, Zoya Kingsbury, Hans Lehrach (Principal Investigator), Ralf Sudbrak (Project Leader), Marcus W. Albrecht, Vyacheslav S. Amstislavskiy, Tatiana A. Borodina, Matthias Lienhard, Florian Mertes, Marc Sultan, Bernd Timmermann, Marie-Laure Yaspo, Stephen T. Sherry (Principal Investigator), Gil A. McVean (Principal Investigator), Elaine R. Mardis (Co-Principal Investigator) (Co-Chair), Richard K. Wilson (Co-Principal Investigator), Lucinda Fulton, Robert Fulton, George M. Weinstock, Richard M. Durbin (Principal Investigator), Senduran Balasubramaniam, John Burton, Petr Danecek, Thomas M. Keane, Anja Kolb-Kokocinski, Shane McCarthy, James Stalker, Michael Quail, Jeanette P. Schmidt (Principal Investigator), Christopher J. Davies, Jeremy Gollub, Teresa Webster, Brant Wong, Yiping Zhan, Adam Auton (Principal Investigator), Richard A. Gibbs (Principal Investigator), Fuli Yu (Project Leader), Matthew Bainbridge, Danny Challis, Uday S. Evani, James Lu, Donna Muzny, Uma Nagaswamy, Jeff Reid, Aniko Sabo, Yi Wang, Jin Yu, Jun Wang (Principal Investigator), Lachlan J. M. Coin, Lin Fang, Xiaosen Guo, Xin Jin, Guoqing Li, Qibin Li, Yingrui Li, Zhenyu Li, Haoxiang Lin, Binghang Liu, Ruibang Luo, Nan Qin, Haojing Shao, Bingqiang Wang, Yinlong Xie, Chen Ye, Chang Yu, Fan Zhang, Hancheng Zheng, Hongmei Zhu, Gabor T. Marth (Principal Investigator), Erik P. Garrison, Deniz Kural, Wan-Ping Lee, Wen Fung Leong, Alistair N. Ward, Jiantao Wu, Mengyao Zhang, Charles Lee (Principal Investigator), Lauren Griffin, Chih-Heng Hsieh, Ryan E. Mills, Xinghua Shi, Marcin von Grotthuss, Chengsheng Zhang, Mark J. Daly (Principal Investigator), Mark A. DePristo (Project Leader), David M. Altshuler, Eric Banks, Gaurav Bhatia, Mauricio O. Carneiro, Guillermo del Angel, Stacey B. Gabriel, Giulio Genovese, Namrata Gupta, Robert E. Handaker, Chris Hartl, Eric S. Lander, Steven A. McCarroll, James C. Nemes, Ryan E. Poplin, Stephen F. Schaffner, Khalid Shakir, Seungtae C. Yoon (Principal Investigator), Jayon Lihm, Vladimir Makarov, Hanjun Jin (Principal Investigator), Wook Kim, Ki Cheol Kim, Jan O. Korbel (Principal Investigator), Tobias Rausch, Paul Flicek (Principal Investigator), Kathryn Beal, Laura Clarke, Fiona Cunningham, Javier Herrero, William M. McLaren, Graham R. S. Ritchie, Richard E. Smith, Xiangqun Zheng-Bradley, Andrew G. Clark (Principal Investigator), Srikanth Gottipati, Alon Keinan, Juan L. Rodriguez-Flores, Pardis C. Sabeti (Principal Investigator), Sharon R. Grossman, Shervin Tabrizi, Ridhi Tariyal, David N. Cooper (Principal Investigator), Edward V. Ball, Peter D. Stenson, David R. Bentley (Principal Investigator), Bret Barnes, Markus Bauer, R. Keira Cheetham, Tony Cox, Michael Eberle, Sean Humphray, Scott Kahn, Lisa Murray, John Peden, Richard Shaw, Kai Ye (Principal Investigator), Mark A. Batzer (Principal Investigator), Miriam K. Konkel, Jerilyn A.

Walker, Daniel G. MacArthur (Principal Investigator), Monkol Lek, Sudbrak (Project Leader), Vyacheslav S. Amstislavskiy, Ralf Herwig, Mark D. Shriver (Principal Investigator), Carlos D. Bustamante (Principal Investigator), Jake K. Byrnes, Francisco M. De La Vega, Simon Gravel, Eimear E. Kenny, Jeffrey M. Kidd, Phil Lacroute, Brian K. Maples, Andres Moreno-Estrada, Fouad Zakharia, Eran Halperin (Principal Investigator), Yael Baran, David W. Craig (Principal Investigator), Alexis Christoforides, Nils Homer, Tyler Izatt, Ahmet A. Kurdoglu, Shripad A. Sinari, Kevin Squire, Stephen T. Sherry (Principal Investigator), Chunlin Xiao, Jonathan Sebat (Principal Investigator), Vineet Bafna, Kenny Ye, Esteban G. Burchard (Principal Investigator), Ryan D. Hernandez (Principal Investigator), Christopher R. Gignoux, David Haussler (Principal Investigator), Sol J. Katzman, W. James Kent, Bryan Howie, Andres Ruiz-Linares (Principal Investigator), The 1000 Genomes Project Consortium, Corresponding Author, Steering committee, Production group:, Baylor College of Medicine, BGI-Shenzhen, Broad Institute of MIT and Harvard, European Bioinformatics Institute, Illumina, Max Planck Institute for Molecular Genetics, US National Institutes of Health, University of Oxford, Washington University in St Louis, Wellcome Trust Sanger Institute, Analysis group:, Affymetrix, Albert Einstein College of Medicine, Boston College, Brigham and Women's Hospital, Cold Spring Harbor Laboratory, Dankook University, European Molecular Biology Laboratory, Cornell University, Harvard University, Human Gene Mutation Database, Leiden University Medical Center, Louisiana State University, Massachusetts General Hospital, Pennsylvania State University, Stanford University, Tel-Aviv University, Translational Genomics Research Institute, San Diego University of California, San Francisco University of California, Santa Cruz University of California, University of Chicago, University College London, and University of Geneva. An integrated map of genetic variation from 1,092 human genomes. *Nature*, 491(7422):56–65, November 2012. Number: 7422 Publisher: Nature Publishing Group.

- [62] Eran Meshorer and Tom Misteli. Chromatin in pluripotent embryonic stem cells and differentiation. *Nature Reviews Molecular Cell Biology*, 7(7):540–546, July 2006.
- [63] Tarjei S. Mikkelsen, Manching Ku, David B. Jaffe, Biju Issac, Erez Lieberman, Georgia Giannoukos, Pablo Alvarez, William Brockman, Tae-Kyung Kim, Richard P. Koche, William Lee, Eric Mendenhall, Aisling O'Donovan, Aviva Presser, Carsten Russ, Xiaohui Xie, Alexander Meissner, Marius Wernig, Rudolf Jaenisch, Chad Nusbaum, Eric S. Lander, and Bradley E. Bernstein. Genome-wide maps of chromatin state in pluripotent and lineage-committed cells. *Nature*, 448(7153):553–560, August 2007.
- [64] Saverio Minucci, Jiemin Wong, Jorge C. G. Blanco, Yun-Bo Shi, Alan P. Wolffe, and Keiko Ozato. Retinoid Receptor-Induced Alteration of the Chromatin Assembled on a Ligand-Responsive Promoter in *Xenopus* Oocytes. *Molecular Endocrinology*, 12(3):315–324, March 1998. Publisher: Oxford Academic.

- [65] T. J. Mitchell and J. J. Beauchamp. Bayesian Variable Selection in Linear Regression. *Journal of the American Statistical Association*, 83(404):1023–1032, December 1988. Publisher: Taylor & Francis.
- [66] Kei Miyamoto, Khoi T. Nguyen, George E. Allen, Jerome Jullien, Dinesh Kumar, Tomoki Otani, Charles R. Bradshaw, Frederick J. Livesey, Manolis Kellis, and John B. Gurdon. Chromatin Accessibility Impacts Transcriptional Reprogramming in Oocytes. *Cell Reports*, 24(2):304–311, July 2018.
- [67] Kei Miyamoto, Vincent Pasque, Jerome Jullien, and John B. Gurdon. Nuclear actin polymerization is required for transcriptional reprogramming of Oct4 by oocytes. *Genes & Development*, 25(9):946–958, May 2011. Company: Cold Spring Harbor Laboratory Press Distributor: Cold Spring Harbor Laboratory Press Institution: Cold Spring Harbor Laboratory Press Label: Cold Spring Harbor Laboratory Press Publisher: Cold Spring Harbor Lab.
- [68] Kei Miyamoto, Yosuke Tajima, Koki Yoshida, Mami Oikawa, Rika Azuma, George E. Allen, Tomomi Tsujikawa, Tomomasa Tsukaguchi, Charles R. Bradshaw, Jerome Jullien, Kazuo Yamagata, Kazuya Matsumoto, Masayuki Anzai, Hiroshi Imai, John B. Gurdon, and Masayasu Yamada. Reprogramming towards totipotency is greatly facilitated by synergistic effects of small molecules. *Biology Open*, 6(4):415–424, April 2017. Publisher: The Company of Biologists Ltd Section: Research Article.
- [69] Kei Miyamoto, Marta Teperek, Kosuke Yusa, George E. Allen, Charles R. Bradshaw, and J. B. Gurdon. Nuclear Wave1 Is Required for Reprogramming Transcription in Oocytes and for Normal Development. *Science*, 341(6149):1002–1005, August 2013. Publisher: American Association for the Advancement of Science Section: Report.
- [70] Vamsi K. Mootha, Cecilia M. Lindgren, Karl-Fredrik Eriksson, Aravind Subramanian, Smita Sihag, Joseph Lehar, Pere Puigserver, Emma Carlsson, Martin Ridderstråle, Esa Laurila, Nicholas Houstis, Mark J. Daly, Nick Patterson, Jill P. Mesirov, Todd R. Golub, Pablo Tamayo, Bruce Spiegelman, Eric S. Lander, Joel N. Hirschhorn, David Altshuler, and Leif C. Groop. PGC-1 β -responsive genes involved in oxidative phosphorylation are coordinately down-regulated in human diabetes. *Nature Genetics*, 34(3):267–273, July 2003. Number: 3 Publisher: Nature Publishing Group.
- [71] Majid Nikpay, Anuj Goel, Hong-Hee Won, Leanne M. Hall, Christina Wilenborg, Stavroula Kanoni, Danish Saleheen, Theodosios Kyriakou, Christopher P. Nelson, Jemma C. Hopewell, Thomas R. Webb, Lingyao Zeng, Abbas Dehghan, Maris Alver, Sebastian M. Armasu, Kirsi Auro, Andrew Bjornes, Daniel I. Chasman, Shufeng Chen, Ian Ford, Nora Franceschini, Christian Gieger, Christopher Grace, Stefan Gustafsson, Jie Huang, Shih-Jen Hwang, Yun Kyoung Kim, Marcus E. Kleber, King Wai Lau, Xiangfeng Lu, Yingchang Lu, Leo-Pekka Lyytikäinen, Evelin Mihailov, Alanna C. Morrison, Natalia

Pervjakova, Liming Qu, Lynda M. Rose, Elias Salfati, Richa Saxena, Markus Scholz, Albert V. Smith, Emmi Tikkanen, Andre Uitterlinden, Xueli Yang, Weihua Zhang, Wei Zhao, Mariza de Andrade, Paul S. de Vries, Natalie R. van Zuydam, Sonia S. Anand, Lars Bertram, Frank Beutner, George Dedoussis, Philippe Frossard, Dominique Gauguier, Alison H. Goodall, Omri Gottesman, Marc Haber, Bok-Ghee Han, Jianfeng Huang, Shapour Jalilzadeh, Thorsten Kessler, Inke R. KÃ¼nig, Lars Lannfelt, Wolfgang Lieb, Lars Lind, Cecilia M. Lindgren, Marja-Liisa Lokki, Patrik K. Magnusson, Nadeem H. Mallick, Narinder Mehra, Thomas Meitinger, Fazal-Ur-Rehman Memon, Andrew P. Morris, Markku S. Nieminen, Nancy L. Pedersen, Annette Peters, Loukianos S. Rallidis, Asif Rasheed, Maria Samuel, Svati H. Shah, Juha Sinisalo, Kathleen E. Stirrups, Stella Trompet, Laiyuan Wang, Khan S. Zaman, Diego Ardissino, Eric Boerwinkle, Ingrid B. Borecki, Erwin P. Bottinger, Julie E. Buring, John C. Chambers, Rory Collins, L. Adrienne Cupples, John Danesh, Ilja Demuth, Roberto Elosua, Stephen E. Epstein, TÃ¤nu Esko, Mary F. Feitosa, Oscar H. Franco, Maria Grazia Franzosi, Christopher B. Granger, Dongfeng Gu, Vilmundur Gudnason, Alistair S. Hall, Anders Hamsten, Tamara B. Harris, Stanley L. Hazen, Christian Hengstenberg, Albert Hofman, Erik Ingelsson, Carlos Iribarren, J. Wouter Jukema, Pekka J. Karhunen, Bong-Jo Kim, Jaspal S. Kooner, Iftikhar J. Kullo, Terho LehtimÃ¤ki, Ruth J. F. Loos, Olle Melander, Andres Metspalu, Winfried MÃ¤rz, Colin N. Palmer, Markus Perola, Thomas Quertermous, Daniel J. Rader, Paul M. Ridker, Samuli Ripatti, Robert Roberts, Veikko Salomaa, Dharambir K. Sanghera, Stephen M. Schwartz, Udo Seedorf, Alexandre F. Stewart, David J. Stott, Joachim Thiery, Pierre A. Zalloua, Christopher J. O'Donnell, Muredach P. Reilly, Themistocles L. Assimes, John R. Thompson, Jeanette Erdmann, Robert Clarke, Hugh Watkins, Sekar Kathiresan, Ruth McPherson, Panos Deloukas, Heribert Schunkert, Nilesh J. Samani, and Martin Farrall. A comprehensive 1,000 Genomes-based genome-wide association meta-analysis of coronary artery disease. *Nature Genetics*, 47(10):1121–1130, October 2015.

- [72] H. Okazawa, K. Okamoto, F. Ishino, T. Ishino-Kaneko, S. Takeda, Y. Toyoda, M. Muramatsu, and H. Hamada. The oct3 gene, a gene for an embryonic transcription factor, is controlled by a retinoic acid repressible enhancer. *The EMBO Journal*, 10(10):2997–3005, October 1991. Publisher: John Wiley & Sons, Ltd.
- [73] Nuala A. O'Leary, Mathew W. Wright, J. Rodney Brister, Stacy Ciufu, Diana Haddad, Rich McVeigh, Bhanu Rajput, Barbara Robbertse, Brian Smith-White, Danso Ako-Adjei, Alexander Astashyn, Azat Badretdin, Yiming Bao, Olga Blinkova, Vyacheslav Brover, Vyacheslav Chetvernin, Jinna Choi, Eric Cox, Olga Ermolaeva, Catherine M. Farrell, Tamara Goldfarb, Tripti Gupta, Daniel Haft, Eneida Hatcher, Wratkano Hlavina, Vinita S. Joardar, Vamsi K. Kodali, Wenjun Li, Donna Maglott, Patrick Masterson, Kelly M. McGarvey, Michael R. Murphy, Kathleen O'Neill, Shashikant Pujar, Sanjida H. Rangwala,

- Daniel Rausch, Lillian D. Riddick, Conrad Schoch, Andrei Shkeda, Susan S. Storz, Hanzhen Sun, Francoise Thibaud-Nissen, Igor Tolstoy, Raymond E. Tully, Anjana R. Vatsan, Craig Wallin, David Webb, Wendy Wu, Melissa J. Landrum, Avi Kimchi, Tatiana Tatusova, Michael DiCuccio, Paul Kitts, Terence D. Murphy, and Kim D. Pruitt. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Research*, 44(D1):D733–745, January 2016.
- [74] Yongjin Park, Abhishek K. Sarkar, Liang He, Jose Davila-Velderrain, Philip L. De Jager, and Manolis Kellis. A Bayesian approach to mediation analysis predicts 206 causal target genes in Alzheimer’s disease. *bioRxiv*, page 219428, December 2017. Publisher: Cold Spring Harbor Laboratory Section: New Results.
- [75] Aaron R. Quinlan and Ira M. Hall. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*, 26(6):841–842, March 2010. Publisher: Oxford Academic.
- [76] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2017.
- [77] Roadmap Epigenomics Consortium, Anshul Kundaje, Wouter Meuleman, Jason Ernst, Misha Bilenky, Angela Yen, Alireza Heravi-Moussavi, Pouya Kheradpour, Zhizhuo Zhang, Jianrong Wang, Michael J. Ziller, Viren Amin, John W. Whitaker, Matthew D. Schultz, Lucas D. Ward, Abhishek Sarkar, Gerald Quon, Richard S. Sandstrom, Matthew L. Eaton, Yi-Chieh Wu, Andreas R. Pfening, Xinchun Wang, Melina Claussnitzer, Yaping Liu, Cristian Coarfa, R. Alan Harris, Noam Shores, Charles B. Epstein, Elizabeta Gjoneska, Danny Leung, Wei Xie, R. David Hawkins, Ryan Lister, Chibo Hong, Philippe Gascard, Andrew J. Mungall, Richard Moore, Eric Chuah, Angela Tam, Theresa K. Canfield, R. Scott Hansen, Rajinder Kaul, Peter J. Sabo, Mukul S. Bansal, Annaick Carles, Jesse R. Dixon, Kai-How Farh, Soheil Feizi, Rosa Karlic, Ah-Ram Kim, Ashwinikumar Kulkarni, Daofeng Li, Rebecca Lowdon, GiNell Elliott, Tim R. Mercer, Shane J. Neph, Vitor Onuchic, Paz Polak, Nisha Rajagopal, Pradipta Ray, Richard C. Sallari, Kyle T. Siebenthal, Nicholas A. Sinnott-Armstrong, Michael Stevens, Robert E. Thurman, Jie Wu, Bo Zhang, Xin Zhou, Arthur E. Beaudet, Laurie A. Boyer, Philip L. De Jager, Peggy J. Farnham, Susan J. Fisher, David Haussler, Steven J. M. Jones, Wei Li, Marco A. Marra, Michael T. McManus, Shamil Sunyaev, James A. Thomson, Thea D. Tlsty, Li-Huei Tsai, Wei Wang, Robert A. Waterland, Michael Q. Zhang, Lisa H. Chadwick, Bradley E. Bernstein, Joseph F. Costello, Joseph R. Ecker, Martin Hirst, Alexander Meissner, Aleksandar Milosavljevic, Bing Ren, John A. Stamatoyannopoulos, Ting Wang, and Manolis Kellis. Integrative analysis of 111 reference human epigenomes. *Nature*, 518(7539):317–330, February 2015.
- [78] Dustin E. Schones and Keji Zhao. Genome-wide approaches to studying chromatin modifications. *Nature Reviews Genetics*, 9(3):179–191, March 2008.

- [79] Andrey A. Shabalín. Matrix eQTL: ultra fast eQTL analysis via large matrix operations. *Bioinformatics*, 28(10):1353–1358, May 2012. Publisher: Oxford Academic.
- [80] Yin Shen, Feng Yue, David F. McCleary, Zhen Ye, Lee Edsall, Samantha Kuan, Ulrich Wagner, Jesse Dixon, Leonard Lee, Victor V. Lobanenko, and Bing Ren. A map of the cis -regulatory sequences in the mouse genome. *Nature*, 488(7409):116–120, August 2012. Number: 7409 Publisher: Nature Publishing Group.
- [81] Huwenbo Shi, Gleb Kichaev, and Bogdan Pasaniuc. Contrasting the Genetic Architecture of 30 Complex Traits from Summary Association Data. *The American Journal of Human Genetics*, 99(1):139–153, July 2016.
- [82] Abdenour Soufi, Greg Donahue, and Kenneth S. Zaret. Facilitators and Impediments of the Pluripotency Reprogramming Factors’ Initial Engagement with the Genome. *Cell*, 151(5):994–1004, November 2012.
- [83] Sarah L. Spain and Jeffrey C. Barrett. Strategies for fine-mapping complex traits. *Human Molecular Genetics*, 24(R1):R111–R119, October 2015. Publisher: Oxford Academic.
- [84] Jurg Stalder, Alf Larsen, James D. Engel, Maureen Dolan, Mark Groudine, and Harold Weintraub. Tissue-specific DNA cleavages in the globin chromatin domain introduced by DNAase I. *Cell*, 20(2):451–460, June 1980. Publisher: Elsevier.
- [85] Oliver Stegle, Leopold Parts, Matias Piipari, John Winn, and Richard Durbin. Using probabilistic estimation of expression residuals (PEER) to obtain increased power and interpretability of gene expression analyses. *Nature Protocols*, 7(3):500–507, March 2012. Number: 3 Publisher: Nature Publishing Group.
- [86] Przemyslaw Stempor and Julie Ahringer. SeqPlots - Interactive software for exploratory data analyses, pattern discovery and visualization in genomics. *Wellcome Open Research*, 1, November 2016.
- [87] Aravind Subramanian, Pablo Tamayo, Vamsi K. Mootha, Sayan Mukherjee, Benjamin L. Ebert, Michael A. Gillette, Amanda Paulovich, Scott L. Pomeroy, Todd R. Golub, Eric S. Lander, and Jill P. Mesirov. Gene Set Enrichment Analysis: A Knowledge-Based Approach for Interpreting Genome-Wide Expression Profiles. *Proceedings of the National Academy of Sciences of the United States of America*, 102(43):15545–15550, 2005. Publisher: National Academy of Sciences.
- [88] Kazutoshi Takahashi, Koji Tanabe, Mari Ohnuki, Megumi Narita, Tomoko Ichisaka, Kiichiro Tomoda, and Shinya Yamanaka. Induction of Pluripotent Stem Cells from Adult Human Fibroblasts by Defined Factors. *Cell*, 131(5):861–872, November 2007.

- [89] Atsushi Takata, Naomichi Matsumoto, and Tadafumi Kato. Genome-wide identification of splicing QTLs in the human brain and their enrichment among schizophrenia-associated loci. *Nature Communications*, 8(1):1–11, February 2017. Number: 1 Publisher: Nature Publishing Group.
- [90] Robert E. Thurman, Eric Rynes, Richard Humbert, Jeff Vierstra, Matthew T. Maurano, Eric Haugen, Nathan C. Sheffield, Andrew B. Stergachis, Hao Wang, Benjamin Vernot, Kavita Garg, Sam John, Richard Sandstrom, Daniel Bates, Lisa Boatman, Theresa K. Canfield, Morgan Diegel, Douglas Dunn, Abigail K. Ebersol, Tristan Frum, Erika Giste, Audra K. Johnson, Ericka M. Johnson, Tanya Kutuyavin, Bryan Lajoie, Bum-Kyu Lee, Kristen Lee, Darin London, Dimitra Lotakis, Shane Neph, Fidencio Neri, Eric D. Nguyen, Hongzhu Qu, Alex P. Reynolds, Vaughn Roach, Alexias Safi, Minerva E. Sanchez, Amartya Sanyal, Anthony Shafer, Jeremy M. Simon, Lingyun Song, Shinny Vong, Molly Weaver, Yongqi Yan, Zhancheng Zhang, Zhuzhu Zhang, Boris Lenhard, Muneesh Tewari, Michael O. Dorschner, R. Scott Hansen, Patrick A. Navas, George Stamatoyannopoulos, Vishwanath R. Iyer, Jason D. Lieb, Shamil R. Sunyaev, Joshua M. Akey, Peter J. Sabo, Rajinder Kaul, Terrence S. Furey, Job Dekker, Gregory E. Crawford, and John A. Stamatoyannopoulos. The accessible chromatin landscape of the human genome. *Nature*, 489(7414):75–82, September 2012.
- [91] Peter M. Visscher, Naomi R. Wray, Qian Zhang, Pamela Sklar, Mark I. McCarthy, Matthew A. Brown, and Jian Yang. 10 Years of GWAS Discovery: Biology, Function, and Translation. *American Journal of Human Genetics*, 101(1):5–22, July 2017.
- [92] Philipp Voigt, Wee-Wei Tee, and Danny Reinberg. A double take on bivalent promoters. *Genes & Development*, 27(12):1318–1338, June 2013. Company: Cold Spring Harbor Laboratory Press Distributor: Cold Spring Harbor Laboratory Press Institution: Cold Spring Harbor Laboratory Press Label: Cold Spring Harbor Laboratory Press Publisher: Cold Spring Harbor Lab.
- [93] Klaudia Walter, Josine L. Min, Jie Huang, Lucy Crooks, Yasin Memari, Shane McCarthy, John R. B. Perry, ChangJiang Xu, Marta Futema, Daniel Lawson, Valentina Iotchkova, Stephan Schiffels, Audrey E. Hendricks, Petr Danecek, Rui Li, James Floyd, Louise V. Wain, InÃs Barroso, Steve E. Humphries, Matthew E. Hurles, Eleftheria Zeggini, Jeffrey C. Barrett, Vincent Plagnol, J. Brent Richards, Celia M. T. Greenwood, Nicholas J. Timpson, Richard Durbin, Nicole Soranzo, Senduran Bala, Peter Clapham, Guy Coates, Tony Cox, Allan Daly, Petr Danecek, Yuanping Du, Richard Durbin, Sarah Edkins, Peter Ellis, Paul Flicek, Xiaosen Guo, Xueqin Guo, Liren Huang, David K. Jackson, Chris Joyce, Thomas Keane, Anja Kolb-Kokocinski, Cordelia Langford, Yingrui Li, Jieqin Liang, Hong Lin, Ryan Liu, John Maslen, Shane McCarthy, Dawn Muddyman, Michael A. Quail, Jim Stalker, Jianping Sun, Jing Tian, Guangbiao Wang, Jun Wang, Yu Wang, Kim Wong, Pingbo Zhang, InÃs Barroso, Ewan Birney, Chris Boustred, Lu Chen, Gail Clement, Massimiliano Cocca,

Petr Danecek, George Davey Smith, Ian N. M. Day, Aaron Day-Williams, Thomas Down, Ian Dunham, Richard Durbin, David M. Evans, Tom R. Gaunt, Matthias Geihs, Celia M. T. Greenwood, Deborah Hart, Audrey E. Hendricks, Bryan Howie, Jie Huang, Tim Hubbard, Pirro Hysi, Valentina Iotchkova, Yalda Jamshidi, Konrad J. Karczewski, John P. Kemp, Genevieve Lachance, Daniel Lawson, Monkol Lek, Margarida Lopes, Daniel G. MacArthur, Jonathan Marchini, Massimo Mangino, Iain Mathieson, Shane McCarthy, Yasin Memari, Sarah Metrustry, Josine L. Min, Alireza Moayyeri, Dawn Muddyman, Kate Northstone, Kalliope Panoutsopoulou, Lavinia Paternoster, John R. B. Perry, Lydia Quaye, J. Brent Richards, Susan Ring, Graham R. S. Ritchie, Stephan Schiffels, Hashem A. Shihab, So-Youn Shin, Kerrin S. Small, Mar~Aja Soler Artigas, Nicole Soranzo, Lorraine Southam, Timothy D. Spector, Beate St Pourcain, Gabriela Surdulescu, Ioanna Tachmazidou, Nicholas J. Timpson, Martin D. Tobin, Ana M. Valdes, Peter M. Visscher, Louise V. Wain, Klaudia Walter, Kirsten Ward, Scott G. Wilson, Kim Wong, Jian Yang, Eleftheria Zeggini, Feng Zhang, Hou-Feng Zheng, Richard Anney, Muhammad Ayub, Jeffrey C. Barrett, Douglas Blackwood, Patrick F. Bolton, Gerome Breen, David A. Collier, Nick Craddock, Lucy Crooks, Sarah Curran, David Curtis, Richard Durbin, Louise Gallagher, Daniel Geschwind, Hugh Gurling, Peter Holmans, Irene Lee, Jouko L~Åunnqvist, Shane McCarthy, Peter McGuffin, Andrew M. McIntosh, Andrew G. McKechnie, Andrew McQuillin, James Morris, Dawn Muddyman, Michael C. O'Donovan, Michael J. Owen, Aarno Palotie, Jeremy R. Parr, Tiina Paunio, Olli Pietilainen, Karola Rehnstr~Åm, Sally I. Sharp, David Skuse, David St Clair, Jaana Suvisaari, James T. R. Walters, Hywel J. Williams, In~Åls Barroso, Elena Bochukova, Rebecca Bounds, Anna Dominiczak, Richard Durbin, I. Sadaf Farooqi, Audrey E. Hendricks, Julia Keogh, Ga~Ånille Marenne, Shane McCarthy, Andrew Morris, Dawn Muddyman, Stephen O'Rahilly, David J. Porteous, Blair H. Smith, Ioanna Tachmazidou, Eleanor Wheeler, Eleftheria Zeggini, Saeed Al Turki, Carl A. Anderson, Dinu Antony, In~Åls Barroso, Phil Beales, Jamie Bentham, Shoumo Bhattacharya, Mattia Calissano, Keren Carss, Krishna Chatterjee, Sebahattin Cirak, Catherine Cosgrove, Richard Durbin, David R. Fitzpatrick, James Floyd, A. Reghan Foley, Christopher S. Franklin, Marta Futema, Detelina Grozeva, Steve E. Humphries, Matthew E. Hurles, Shane McCarthy, Hannah M. Mitchison, Dawn Muddyman, Francesco Muntoni, Stephen O'Rahilly, Alexandros Onoufriadis, Victoria Parker, Felicity Payne, Vincent Plagnol, F. Lucy Raymond, Nicola Roberts, David B. Savage, Peter Scambler, Miriam Schmidts, Nadia Schoenmakers, Robert K. Semple, Eva Serra, Olivera Spasic-Boskovic, Elizabeth Stevens, Margriet van Kogelenberg, Parthiban Vijayarangakannan, Klaudia Walter, Kathleen A. Williamson, Crispian Wilson, Tamieka Whyte, Antonio Ciampi, Celia M. T. Greenwood, Audrey E. Hendricks, Rui Li, Sarah Metrustry, Karim Oualkacha, Ioanna Tachmazidou, ChangJiang Xu, Eleftheria Zeggini, Martin Bobrow, Patrick F. Bolton, Richard Durbin, David R. Fitzpatrick, Heather Griffin, Matthew E. Hurles, Jane Kaye, Karen Kennedy, Alastair Kent, Dawn Muddyman, Francesco Muntoni, F. Lucy Raymond, Robert K.

- Semple, Carol Smee, Timothy D. Spector, Nicholas J. Timpson, Ruth Charlton, Rosemary Ekong, Marta Futema, Steve E. Humphries, Farrah Khawaja, Luis R. Lopes, Nicola Migone, Stewart J. Payne, Vincent Plagnol, Rebecca C. Pollitt, Sue Povey, Cheryl K. Ridout, Rachel L. Robinson, Richard H. Scott, Adam Shaw, Petros Syrris, Rohan Taylor, Anthony M. Vandersteen, Jeffrey C. Barrett, InÃs Barroso, George Davey Smith, Richard Durbin, I. Sadaf Farooqi, David R. Fitzpatrick, Matthew E. Hurles, Jane Kaye, The UK10K Consortium, Writing group, Production group, Cohorts group, Neurodevelopmental disorders group, Obesity group, Rare disease group, Statistics group, Ethics group, Incidental findings group, and Management committee. The UK10K project identifies rare variants in health and disease. *Nature*, 526(7571):82–90, October 2015. Number: 7571 Publisher: Nature Publishing Group.
- [94] Dong Wang, Ivan Garcia-Bassets, Chris Benner, Wenbo Li, Xue Su, Yiming Zhou, Jinsong Qiu, Wen Liu, Minna U. Kaikkonen, Kenneth A. Ohgi, Christopher K. Glass, Michael G. Rosenfeld, and Xiang-Dong Fu. Reprogramming transcription by distinct classes of enhancers functionally defined by eRNA. *Nature*, 474(7351):390–394, June 2011. Number: 7351 Publisher: Nature Publishing Group.
- [95] Wei Wang, Jian Yang, Hui Liu, Dong Lu, Xiongfeng Chen, Zenon Zenonos, Lia S. Campos, Roland Rad, Ge Guo, Shujun Zhang, Allan Bradley, and Pentao Liu. Rapid and efficient reprogramming of somatic cells to induced pluripotent stem cells by retinoic acid receptor gamma and liver receptor homolog 1. *Proceedings of the National Academy of Sciences*, 108(45):18283–18288, November 2011. Publisher: National Academy of Sciences Section: Biological Sciences.
- [96] Danielle Welter, Jacqueline MacArthur, Joannella Morales, Tony Burdett, Peggy Hall, Heather Junkins, Alan Klemm, Paul Flicek, Teri Manolio, Lucia Hindorff, and Helen Parkinson. The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Research*, 42(D1):D1001–D1006, January 2014. Publisher: Oxford Academic.
- [97] Jason A. West, April Cook, Burak H. Alver, Matthias Stadtfeld, Aimee M. Deaton, Konrad Hochedlinger, Peter J. Park, Michael Y. Tolstorukov, and Robert E. Kingston. Nucleosomal occupancy changes locally over key regulatory regions during cell differentiation and reprogramming. *Nature Communications*, 5(1):1–12, August 2014. Number: 1 Publisher: Nature Publishing Group.
- [98] WarrenÃA. Whyte, DavidÃA. Orlando, Denes Hnisz, BrianÃA. Abraham, CharlesÃA. Lin, MichaelÃA. Kagey, PeterÃA. Rahl, TongÃA. Ihn Lee, and RichardÃA. Young. Master Transcription Factors and Mediator Establish Super-Enhancers at Key Cell Identity Genes. *Cell*, 153(2):307–319, April 2013.
- [99] Jingyi Wu, Bo Huang, He Chen, Qiangzong Yin, Yang Liu, Yunlong Xiang, Bingjie Zhang, Bofeng Liu, Qiujun Wang, Weikun Xia, Wenzhi Li, Yuanyuan

- Li, Jing Ma, Xu Peng, Hui Zheng, Jia Ming, Wenhao Zhang, Jing Zhang, Geng Tian, Feng Xu, Zai Chang, Jie Na, Xuerui Yang, and Wei Xie. The landscape of accessible chromatin in mammalian preimplantation embryos. *Nature*, 534(7609):652–657, June 2016.
- [100] Wei Xie, Matthew D. Schultz, Ryan Lister, Zhonggang Hou, Nisha Rajagopal, Pradipta Ray, John W. Whitaker, Shulan Tian, R. David Hawkins, Danny Leung, Hongbo Yang, Tao Wang, Ah Young Lee, Scott A. Swanson, Jiuchun Zhang, Yun Zhu, Audrey Kim, Joseph R. Nery, Mark A. Urich, Samantha Kuan, Chia-an Yen, Sarit Klugman, Pengzhi Yu, Kran Suknuntha, Nicholas E. Propson, Huaming Chen, Lee E. Edsall, Ulrich Wagner, Yan Li, Zhen Ye, Ashwinikumar Kulkarni, Zhenyu Xuan, Wen-Yu Chung, Neil C. Chi, Jessica E. Antosiewicz-Bourget, Igor Slukvin, Ron Stewart, Michael Q. Zhang, Wei Wang, James A. Thomson, Joseph R. Ecker, and Bing Ren. Epigenomic Analysis of Multilineage Differentiation of Human Embryonic Stem Cells. *Cell*, 153(5):1134–1148, May 2013.
- [101] Jingli A. Zhang, Ali Mortazavi, Brian A. Williams, Barbara J. Wold, and Ellen V. Rothenberg. Dynamic Transformations of Genome-wide Epigenetic Marking and Transcriptional Control Establish T Cell Identity. *Cell*, 149(2):467–482, April 2012.
- [102] Yong Zhang, Tao Liu, Clifford A. Meyer, Jirme Eeckhoutte, David S. Johnson, Bradley E. Bernstein, Chad Nusbaum, Richard M. Myers, Myles Brown, Wei Li, and X. Shirley Liu. Model-based analysis of ChIP-Seq (MACS). *Genome Biology*, 9(9):R137, 2008.
- [103] Vicky W. Zhou, Alon Goren, and Bradley E. Bernstein. Charting histone modifications and the functional organization of mammalian genomes. *Nature Reviews Genetics*, 12(1):7–18, January 2011.
- [104] Jiang Zhu, Mazhar Adli, James Y. Zou, Griet Verstappen, Michael Coyne, Xiaolan Zhang, Timothy Durham, Mohammad Miri, Vikram Deshpande, Philip L. DeJager, David A. Bennett, Joseph A. Houmard, Deborah M. Muoio, Tamer T. Onder, Ray Camahort, Chad A. Cowan, Alexander Meissner, Charles B. Epstein, Noam Shores, and Bradley E. Bernstein. Genome-wide Chromatin State Transitions Associated with Developmental and Environmental Cues. *Cell*, 152(3):642–654, January 2013.