

Dissecting the gene-regulatory circuitry of disease-associated genetic variants

By Taylor Herr

Submitted to the Department of Electrical Engineering and Computer Science

In partial fulfillment of the requirements for the degree of

Masters of Engineering in Computer Science and Electrical Engineering

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2019

© Massachusetts Institute of Technology 2017. All rights reserved.

Author
Taylor Herr

Thesis Advisor.....
Manolis Kellis

Accepted By.....
Katrina LaCurts
Chair, Master of Engineering Thesis Committee

Dissecting the gene-regulatory circuitry of disease-associated genetic variants

By Taylor Herr

Submitted to the Department of Electrical Engineering and Computer Science

May 24, 2019

In partial fulfillment of the requirements for the degree of

Masters of Engineering in Computer Science and Electrical Engineering

Abstract: Disease-associated nucleotides lie primarily in non-coding regions, increasing the urgency of understanding how gene-regulatory circuitry impacts human disease. Here, we use the increasing availability of functional genomics datasets and models elucidating how regulatory proteins control genes, to evaluate the impact of genetic variants on the activity of diverse regulators. First, we generate a comprehensive compendium of predicted binding intensities across the entire genome for over 500 transcription factors. Second, we create a novel dataset to connect how these binding intensities change in the context of disease datasets. Third, we develop a statistical framework to integrate these two datasets using dimensionality reduction, latent cluster discovery, and topic modeling. We use these techniques to show that regulatory proteins with analogous biological functions share similar global changes in binding due to genome-wide genetic variation. We also use our framework to discover a latent set of topics behind all genomic locations in chromosome 1, to link the locations in each of the topic clusters with a class of related diseases, and to show that relevant biological processes are statistically enriched in the genomic locations most related to each cluster.

Table of Contents

Introduction	8
Chapter 1 - Background Research:	11
Kipoi: Model Zoo for Genomics	11
<i>DeepBind</i>	12
Latent Dirichlet Allocation:	14
Variational Autoencoder and TSNE	17
GWAS Summary Statistics	20
Historical Transcription Factor Results	22
Chapter 2 – Building the Datasets:	24
Building and Analyzing the Variant Effect Database	24
<i>Selection of Kipoi Model:</i>	24
<i>Choosing the Variants and Constructing the vcfs:</i>	26
Calculating Variant Effect Sizes:	26
<i>Analyzing the Variant Effect Sizes:</i>	28
Building the GWAS, Transcription Factor, Genomic Location Tensor	29
<i>Selection of the GWAS Traits</i>	29
<i>Creating a Framework to Build the Tensor in Parallel</i>	31
<i>Creating a Framework to Build the Tensor in Parallel</i>	32

<i>The LD-corrected strategy</i>	34
<i>Runtime Considerations for Building the Dataset</i>	35
Chapter 3 – Dataset Analysis Results	37
Variant Effect Database Results:	37
<i>GWAS Correlation Analysis:</i>	37
<i>Autoencoder and TSNE Results:</i>	40
GWAS, Location, TF Tensor Analysis:	45
<i>Latent Dirichlet Allocation on the Tensor</i>	46
<i>Flattening the Tensor along both the TF and GWAS Axis</i>	66
Conclusion:	72
Table of Figures	76
Appendix	81
<i>Gwas Names:</i>	81
<i>Transcription Factors</i>	84
References	89

Introduction

The field of computational genetics is currently experiencing a revolution powered by the emergence of many important technologies. Low-cost sequencing has allowed researchers to determine the exact DNA sequences for massive cohorts of people. Cheap data storage has given researchers the ability to access and replicate these massive data sets. Modern machine learning has made it incredibly easy for non-statisticians to spin up state-of-the-art models to effectively make predictions from data. Huge GPU mega-clusters have made it quick and easy to train these models. Thanks to all of these advances, many of the questions that seemed outlandish to even ask can now be easily answered.

One of these questions involves determining the relationship between transcription factors and diseases. A transcription factor is a protein which binds to DNA to regulate the expression of a nearby gene. There are hundreds of these proteins, each of which has a diverse function and might help regulate the expression of numerous genes. At each genomic location, it could be possible that a single base pair change or other genetic variant might affect how one of these transcription factor functions in the body. These small perturbations can have massive downstream impacts on the diseases that a person experiences on an everyday basis. But there are billions of base pairs in the human genome. How can we even begin to understand the relationship between these transcription factors and the traits that actually matter to a person?

This is where we combine all of the aforementioned tools together. First, thanks to new initiatives and innovations in the machine learning community, it is very easy to grab already pre-trained models for genetics. For a given transcription factor, we can use a unique machine-learning model to predict the binding affinity of that transcription factor on a short segment of

DNA using only the base pair sequence. For any genetic variant, we can then find the delta score for the transcription factor on the DNA sequence with and without the variant. If we do this across all transcription factors and the whole genome, we generate a matrix relating how important each genetic variant is to every transcription factor. Once built using a large computing cluster, we analyze this variant-effect dataset to explore the relationships between different transcription factors.

Next, we bring in a dataset linking the same list of genetic variants with a list of 100s of observable human traits. Thanks to the UK Biobank study, hundreds of thousands of people have had their biological and genetic data tested by researchers. By linking this biological and genetic data, researchers at the Broad Institute have generated a statistical z-score for each pair of tested variant and observed trait. These z-scores simply score how associated a particular genetic variant is with an observed trait. With these statistics in hand, we construct a new dataset: a tensor relating disease, genomic location, and transcription factor activity. For a given choice of genomic window, human trait, and transcription factor, we take the z-scores and the transcription factor scores for the variants in our genomic region and merge them to create a single score. This score serves as a proxy for identifying how relevant a particular transcription factor is to a disease in a small genomic window. When analyzed, this tensor paints a picture of the meta relationships between regions of DNA, transcription factors, and observable human traits.

We perform four different analyses on the two datasets. First, we apply a neural network-based dimensionality reduction technique on the variant effect dataset, and then study this data with a very powerful visualization tool. Second, we cluster the correlation coefficients linking the z-scores of each disease with the variant effect sizes for each transcription factor. Third, we run a topic model on the tensor in order to understand where in the DNA are the

different transcription factors relevant, and what diseases are important to these locations. Finally, we perform a flatten analysis on the tensor to explore latent structures along both the transcription factor and human trait dimensions.

The work of this thesis has been split into three different chapters. The first is background research, a section which sets the intellectual groundwork for the rest of the project. The second describes the construction of two novel datasets: the variant effect dataset; and the GWAS, Transcription Factor, Genomic Location tensor. The third chapter discusses the results of the four different follow-up studies that we perform on the datasets.

Chapter 1 - Background Research:

Deep Learning has slowly infiltrated many different fields, as its state of the art performance and growing ease of implementation make it a very attractive modern choice for performing novel research. In a paper by Yongjin Park and Manolis Kellis, my two main research advisors, they discussed some of the details of how to use deep learning in genomics.ⁱ In general, using machine learning techniques in a genetics context has gotten even easier than it was when they initially wrote their paper, thanks to new initiatives within the genetics community and the ML community at large to consolidate machine learning models and make them more available. This fundamental belief, that deep learning is easy, ready to use, and already commoditized led to the development of the ideas that became this project.

Kipoi: Model Zoo for Genomics

In the recent past, if you wanted to use a particular machine learning model from a genetics paper, you had to hunt down the authors or attempt to make poorly written code on Github work on your machine. This is no longer the case, now that a standardized repository of machine learning models exists for performing genomics. Kipoi is this model zoo. Authors of papers can easily convert their existing Pytorch models into Kipoi models, and then upload them to the cloud for use by any researcher.ⁱⁱ

For example, the model zoo already contains machine learning models built to impute RNA expression, histone modification, DNA accessibility, RNA splicing, RNA binding, DNA binding, and much more. Now, we can think of computing any imputed intermediate phenotype as commoditized and trivial, as we can easily spin up any one of these Kipoi models to compute

our intended intermediate phenotype for any given genomic location. This opens a world of possible use cases.

One major benefit of the Kipoi suite is that it is more than just a repository of models, it also is a software framework which allows the user to easily extend any machine learning model to a number of other tasks. The two major plugins which exist for Kipoi include one which predicts the effects of variants, as well as one which will perform model interpretation using feature importance scores. The first of these plugins turned out to be vitally important for our project, as we were able to quickly and easily calculate the predicted difference of a variant on the intermediate phenotype score for all SNPs in our dataset.

For this thesis, the main set of models we used were those published by the authors of the DeepBind paper. The DeepBind models allow the user to impute the binding affinity for any of a large set of transcription factors on any genomic region. Thanks to the aforementioned benefits of the Kipoi model, this means we can easily compute variant scores across the entire genome.

DeepBind

The main set of models which will power our analysis is DeepBind, a neural net designed to impute transcription factor binding sites from DNA sequences alone.ⁱⁱⁱ DeepBind achieves state of the art performance for both in vitro and in vivo detection of transcription factors found in ChIP-seq experiments, as well as strong predictive accuracy of transcriptional binding affinity. It achieves this state of the art performance using a simple neural network architecture of a convolutional layer, pooling layer, then finally a neural network layer. For each of 538 transcription factors, an individual DeepBind model was trained, which means for each genetic variant, we can test its effect on a very large number of transcription factors. DeepBind simply

requires a FASTA file and a list of segments of DNA in order to run, and then it spits out transcription factor binding affinity scores. Interestingly, DeepBind also contains information about RNA binding protein affinities. Although this is an avenue that we can explore, the paper describes that their model had far worse predictive power on RBP affinities. Because a weak predictive power would strongly affect the signal, we decided not to analyze these proteins in this work.

The training data used to build each DeepBind model is specific to a particular protein and to a specific protein-binding-site detection technique. It is fine if the binding scores are real-valued measurements or binary class-labels. Because of this, the Kipoi suite contains pre-trained models for predicting experimental scores from PBM, CHIP-seq and HT-SELEX experiments across a number of transcription factors. Hyper-parameters for the models were chosen in a very systematic way. Thirty random sets of hyper-parameters were tested on a threefold cross-validation experiment, and the one with the best validation accuracy was trained using all the data. This ensures that the generation of the final models was unbiased as well as highly-performant across the massive set of models the authors trained.

When the authors tested DeepBind, they found that it significantly outperformed all of the state-of-the-art TF-detection algorithms that participated in the DREAM5 TF-DNA motif recognition challenge. They tested their algorithms under the same conditions as the challenge, and found they had the highest correlation of predicted binding sites with ground truth as well as the highest AUC score out of all models tested. Another way that they validated that their methodology was accurate was they checked that the convolutional layer had a similar structure to known PWMs (position weight matrices) for various transcription factors. They found that the

convolutional weights that the machine-learning model learned were very similar to known PWMs.

Another useful concept that the DeepBind paper suggests is that of a mutation map, which specifies how important each base pair is to the overall predicted transcription factor binding affinity. Knowing that individual point mutations can have a large impact on affinity is very important, as we want to investigate how much individual SNPs change the TF binding affinities.

Other models that could be useful in the Kipoi suite include BassNet for imputing DNA accessibility, CpGenie for imputing DNA methylation, and extended_coda for imputing histone modifications. For example, CpGenie, a model built by MIT's own Zeng et al., can predict the methylation status of all CpGs in a 500 base pair window around a particular variant.^{iv} It does this using a CNN, much like DeepBind. CpGenie had state of the art performance when predicting methylation levels, and also had solid predictive accuracy when trained to detect known meQTL's, verifying that it can predict differences in methylation even when applied to single base pair changes (this is very important for our research).

Latent Dirichlet Allocation:

One of the main motivating works for this research is Latent Dirichlet Allocation, a topic modeling technique which has risen to prominence in the field of natural language processing. LDA is a probabilistic generative model which asserts that each document is a mixture of topics, and that words are generated from topics based on a unique probability distribution for each topic. The algorithm outputs the mixture of topics for each document, as well as the probability of each word for each topic.^v

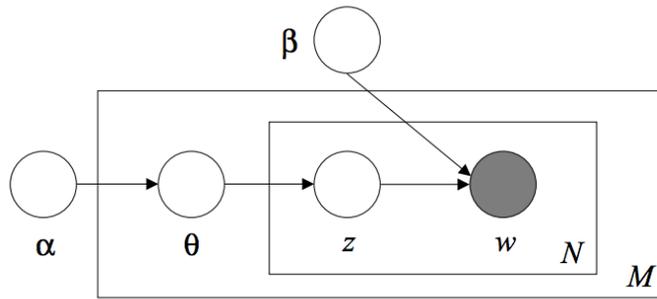


Figure 1: Plate Notation for Latent Dirichlet Allocation. The outer plates represent documents, while the inner plates represent the repeated choice of topics and words.

In LDA, we assume that for each of M documents, an N representing the number of words is chosen based on a Poisson process. Then, a θ is generated based on a Dirichlet random variable with parameter α . This θ represents the probability distribution of topics that make up the document. Using this θ in a multinomial distribution, we can then randomly generate a topic, and then choose words from this topic based on the matrix β , which stores the probability of each word for each topic. The learning process consists of determining the matrix β , as well as the θ for every document using variational inference.

LDA has already been applied to the field of genetics in a number of different ways. Some of the earliest work in this space was performed by Rogers et al., when they built a software implementation of a modified LDA algorithm for analyzing cDNA microarray data sets.^{vi} Their approach is very relevant to our research, as they represented each of the protein expression levels in a continuous fashion rather than a binary fashion. In their context, each word was an expression level, and they modeled these expression levels as normal distributions. This paper contains exact formulations for how to solve for α and β in this context, and thus serves as a good framework for developing a continuous LDA implementation. Building off of this work, Bicego et al. tried a bunch of different dimensionality reduction techniques on the

classification of leukemia based on cell expression data.^{vii} The way they tested the different techniques was to compare the classification accuracy of the very simple K-Nearest Neighbor model applied on the low-dimensional spaces generated by each of the various topic modeling/low dimensionality reduction algorithms. LDA performed the best out of these techniques, supporting the idea that it's a very effective algorithm for discovering latent classes behind sets of "documents" even in a genetics context.

In a later paper, Pinoli et al. used LDA to predict gene function annotations in an older gene annotation dataset, and then found that many of the gene functions they predicted showed up in a later gene annotation dataset.^{viii} This study in particular was interesting because it discovered that one could obtain similar results with a modified, faster version of LDA, known as LDA with Collapsed Gibbs Sampling. In this algorithm, instead of finding values for the priors beta and alpha in a maximum likelihood fashion, they use a Gibbs Sampling approach. They iteratively guess values for these hyper-parameters until they converge to the ones with the best predictive power. Their predictive accuracy for new gene ontology terms was 30-45 percent, which is pretty impressive considering many of the terms they proposed might be accurate if we wait longer for more ontology terms to be generated with time and research.

Liu et al. used LDA to discover modules of miRNA, and without using information about the corresponding mRNA binding sites was able to detect predicted miRNA / mRNA interactions.^{ix} They used a modified version of LDA for handling pairs of data sources, where one acts as an annotation for the other. In the original paper by Blei et al. describing this correspondence-LDA algorithm, the example of annotating images with text was used.^x This generative model works by first breaking the image into N individual regions, then sampling a region description for each based on an LDA model. Then, for each of the M words in a caption,

the word is generated by sampling a region, then sampling a word based on the corresponding region's description. Back to the example of miRNA and mRNA, we now assume that these are the two disparate data sources, and that they're linked with some underlying set of topics. Mapping to the original paper, we let the miRNA be images, and the mRNA be descriptions. Then, the latent classes that Liu discovered represent functional micro-RNA regulatory modules (FMRM). This paper was able to discover some known FMRMs using this approach, and validated the model by showing that without using target binding information, the predicted FMRMs captured hypothesized miRNA / mRNA relationships.

Variational Autoencoder and TSNE

At its core, an autoencoder is a dimensionality reduction technique. The general idea is that you feed a neural network with a very specific architecture some training data, and ask the neural network to predict the same output as the input you feed it. This architecture will have a bottleneck layer, or layer where the data is forced to go through fewer nodes than were used to generate the data (figure 2). Since the edges and layers after the bottleneck layer can be used to reconstruct the data from the activations of the bottleneck layer, we've essentially shown a way to generate the input data from a vector the size of the bottleneck layer. This reduces the dimension of our data to the same size as that of the bottleneck layer.

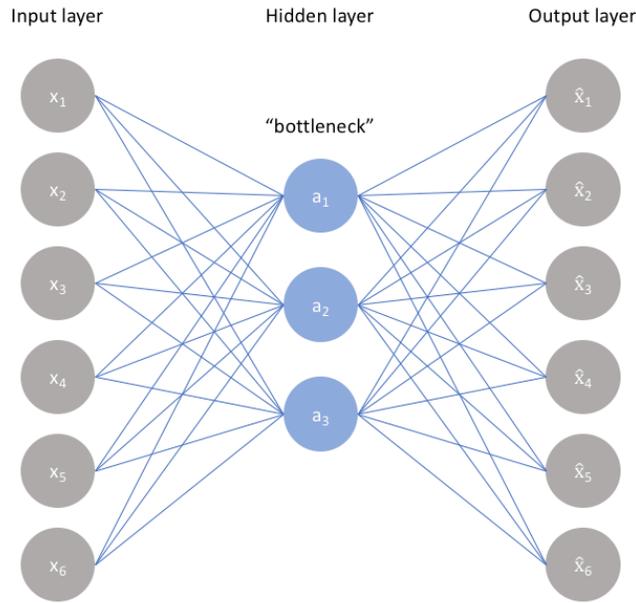


Figure 2: Pictorial Representation of a Sparse Autoencoder. An autoencoder will reduce the dimension of a dataset n to the number of nodes in the bottleneck layer.

Variational Autoencoders are an extremely popular adaptation of the autoencoding dimensionality reduction technique. They've been used for writing captions for images, recommender systems, and text generation.^{xi xii} In the original paper by Kingma et al., the authors laid out an approach to perform variational bayes in a more generalizable format using neural networks.^{xiii} The key to variational bayes is the idea that we can sample approximate posterior distributions from a class of distributions. Then, by choosing the distribution which minimizes the KL divergence, we find a good estimator of our goal posterior distribution. In variational autoencoders, the possible distributions we use to generate the posterior are the set of neural nets with a given architecture. The reason variational autoencoders are so popular, is that they have the same theoretical guarantees of variational bayes at estimating the posterior distribution, but they come with an implicit dimensionality reduction thanks to their autoencoding properties.

Once we have greatly reduced the dimension of our data using the variational autoencoder, we need to further reduce the dimension in order to visualize it. A good algorithm to reduce the dimension down to two for visualization is the t-SNE algorithm.^{xiv} The t-SNE algorithm works by first converting high-dimensional Euclidean distances into conditional probabilities. Similarity under this new mapping is the probability that one high-dimensional data point x_i would pick another data point x_j as its neighbor. Now, we consider low dimensional datapoints y_i and y_j , and assign them corresponding probabilities. The goal is to find the best set of points Y such that the KL divergence of the distribution on X to the distribution on Y is minimized.

In general, t-SNE works incredibly well by actually pushing dissimilar groups of points far away from one another when they are mapped to a low-dimensional space. This is made incredibly obvious by looking at a comparison of two different dimensionality reduction/visualization techniques. As seen in Figure 3, t-SNE is nice in that it actually forces the data to separate into visual clusters. While a lot of other algorithms will successfully move data from different classes into different parts of a 2-D figure, they will not look nearly as well-clustered as the same data pushed through the t-SNE algorithm.

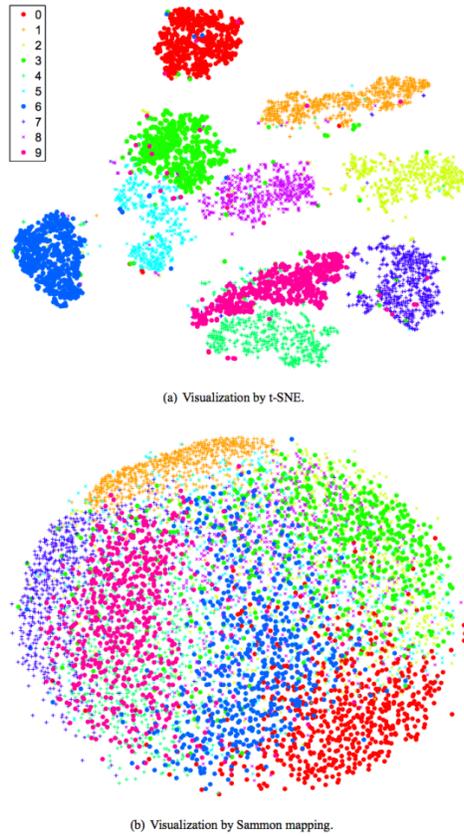


Figure 3: *t-SNE* helps distinguish and separate data clusters that are not clearly separable with other visualization schemes (in this example hand-written digits). (top) *t-SNE* representation of dataset. (bottom) Sammon Mapping of the same data.

GWAS Summary Statistics

In general, we need to use information about how associated a particular SNP is with a trait in order to perform the final mapping from trait to intermediate phenotype. We need knowledge of z-scores and effect sizes for each of the different SNPs across all of the different phenotypes. This is why it is essential that we use a dataset of high quality GWAS summary statistics. Fortunately, Neale et al. have already created an open-source repository of GWAS summary statistics across every trait/SNP pair from the UK Biobank.^{xv} In a broader sense, there's a lot of very interesting things to study in this dataset, as unlocking its underlying structure will yield information about all of the phenotypes that the UK Biobank has collected

data on. As part of our work, we want to explore these GWAS summary statistics to see what interesting insights about human health we learn.

In order to get accurate summary statistics, Neale Lab uses a common GWAS technique known as LD score regression. Initially proposed by Bulik-Sullivan et al., LD score regression uses an elegant relationship between the chi-squared statistic and LD scores to control for population biases.^{xvi} The key equation which describes this clean relationship is:

$$E[\chi^2 | l_j] = \frac{Nh^2 l_j}{M} + Na + 1$$

In this equation, N is the sample size, M is the number of SNPs, h^2 is the heritability, a measures confounding bias, and l_j is the LD score of variant j . The technique is to simply regress the χ^2 statistics from GWAS against the LD scores, and then look at the intercept. The intercept minus one gives us an estimate of the inflation of the test statistics due to confounding bias. If we correct for this confounding bias we get better GWAS summary statistics.

In the original Bulik-Sullivan paper, some very important notes were made about when it is appropriate to use LD score regression. One important aspect is that the trait needs to be highly polygenic. When there are not enough causal variants, the LD score regression estimate is unbiased, but the variance in the intercept can become very large. This is problematic, as we might improperly overcorrect the summary statistics. Another important note is that the reference panel needs to have a similar LD structure to the GWAS panel. If this is not the case, we can again introduce bias into the measurement of our intercept. The final key detail is that LD score regression works best when the variance explained per SNP is uncorrelated with the LD score. All of these details were tested and verified in simulation studies by the authors of the original LD score regression paper.

Historical Transcription Factor Results

After we've built a dataset connecting phenotype to transcription factor, one way we can validate our results is by connecting them back to schizophrenia. There's already a few known transcription factors associated with schizophrenia. Xia et al. discovered that Transcription Factor 4 (TCF4) is highly associated with schizophrenia, after they found that the binding sites for TCF4 were modestly enriched in schizophrenia patients. They additionally found that the genes related to TCF4 were generally associated with brain development.^{xvii} Additionally, Georgieva et al. discovered that variation in the transcription factor OLIG2 is related to an increased risk in schizophrenia.^{xviii} It will be important to verify that these transcription factors emerge as relevant to schizophrenia from our LDA methods.

There's already a lot of precedent for assuming that there exist classes of different transcription factors. Many of these "families" of DNA binding proteins can be grouped together based on shared-evolution, structure, motif-recognition, and function.^{xix} One easy way to validate whether or not the latent topics of motifs we discover are accurate is to compare them to known TF protein families. One example of a family of transcription factors that will be worth investigating is the Nuclear Factor of Activated T-Cells (NFAT) family.^{xx} These particular proteins are interesting biologically because they are the targets of immunosuppressive drugs. Many of these proteins share characteristics of their motifs, which means mutations will likely affect binding affinity scores in similar ways. Thus, we should expect them to have correlated scores with respect to binding affinities, and since they have similar physiological effects, we would expect them to be related to similar phenotypes. This property makes them a good candidate for validation with our LDA methods.

Chapter 2 – Building the Datasets:

Building and Analyzing the Variant Effect Database

The end goal of this portion of the project was to build a database of variant-effect scores. Say we start with a machine learning model that can predict an intermediate phenotype score for a specific sequence of DNA. Now, consider a variant on this piece of DNA, for example a single nucleotide polymorphism. This variant will change the score that the machine learning model predicts for this intermediate phenotype. We want to capture these deltas across all different genomic locations for a number of different intermediate phenotypes and analyze this broader dataset. To begin with, we needed to identify first what intermediate phenotypes we thought were most interesting.

Selection of Kipoi Model:

Our goal was to determine the relationship between an imputed intermediate-phenotype and a set of UK Biobank phenotypes. The first step was to choose this intermediate-phenotype. Some of the possible options included RNA splicing, DNA methylation, DNA accessibility, histone modification, and DNA binding. While all of these intermediate phenotypes would've been interesting to work on, most of them were not highly dimensional. What we really wanted to investigate was a set of intermediate phenotypes with the property that for any given genetic variant, we would have hundreds of different score differentials. With this in mind, evidence was pointing towards transcription factor binding affinity as a good choice for our analysis.

Moreover, transcription factors are a very interesting and not well-understood topic. We know of hundreds of transcription factors that bind to millions of different binding sites in DNA,

however their relationships to diseases and larger, macro-phenotypes is not well-understood. By taking a holistic approach towards studying them in parallel across the entire genome, we were confident that we'd be able to find some interesting results. In addition, without Kipoi there was no way we would be able to do this analysis. With all this in mind, we thought transcription factors would make a good intermediate-phenotype for us to investigate.

The Kipoi suite itself does not have a terribly large number of machine learning models on it (yet), but it does have hundreds of DeepBind models already easily accessible in the framework. The authors of DeepBind published hundreds of models in the Kipoi suite, across not only the hundreds of known human transcription factors, but also across many different species. We made an early decision that we would run the analysis on all DeepBind models, as leveraging the power of larger datasets might reveal insights that haven't been discovered before. DeepBind has a combination of models trained on SELEX and CHIP-seq data, we chose to run our analyses on both sets of data. In total, we decided to find scores for about one million variants on each of 23 chromosomes for 515 different machine learning models.

These models have a very simple interface. You can pass them a FASTA file containing a DNA sequence, as well as an intervals file containing a list of DNA intervals, and the model will score each interval using the machine learning model. Our end goal however was to determine the effect of a set of variants on the calculated model scores. The naïve way to do this would be to edit the individual fasta files with each variant, but this posed a lot of problems. Fortunately, an add-on package to Kipoi had a built-in way to perform this with a simpler interface. In order to use this package, we first needed to create a .vcf file containing entries for each of the variants we wanted to analyze.

Choosing the Variants and Constructing the vcfs:

The end goal of this project was to relate GWAS scores with intermediate-phenotype variant effect sizes, so it made a lot of sense that the variants we choose should be the same as the ones that appear in the GWAS summary statistics we use in the next step. At one point, we were considering two different sets of GWAS summary statistics, but in the end, we decided to use the Neale Lab summary statistics. For this reason, it made a lot of sense to use a GWAS file from the Neale Lab summary stats to create the list of variants we were going to analyze.

Once we had the set of GWAS variants, it was very easy to manually convert each variant into a .vcf line. After creating files with the correct header information, it was easy to finish the chromosomal .vcf files.

Calculating Variant Effect Sizes:

Using the Kipoi-veff suite, we were able to avoid messing with individual FASTA files and instead work through a much cleaner interface. First, we started by building very simple .vcf files for each of the different variants that appear in the Neale Lab summary statistics. Once we have a .vcf file built, we can use a reader in the Kipoi-veff suite to read in each variant, score it, and then write it back to a file using a .vcf writer. We used two different approaches to score the data. The first approach was to just take the difference between the binding affinities for the reference and alternative alleles. The second was to record a DeepSEA variant score. The DeepSEA score at a particular genomic location is the product of the geometric mean value of predicted chromatin effects with the geometric mean of feature effect.

In order to create this large of a dataset, we used the Broad computing clusters and ran each model simultaneously. Thankfully, this computing task is an example of an “embarrassingly

parallel” computation, implying we do not need the machines on the cluster to talk to one another to perform all of the jobs. Given a Kipoi model, we start by creating an array of subtasks, one for each chromosome. The goal of each subtask is to generate the variant effect sizes for all variants on that chromosome for the given Kipoi model. For each of these subtasks, we allocate a worker node which starts by creating a working folder for its computation. This worker then reads in a .vcf file containing the list of variants it’s going to score on that chromosome. It then scores all of the variants according to the Kipoi model using the Kipoi-veff suite. Finally, it writes the output to another .vcf file in the working directory, then performs the necessary housekeeping. Wrapping this whole process is a single script which would launch all of these jobs and sub processes. In order to run this sub-process, we simply run a single Python file after passing it the model name and chromosome.

Running jobs on the cluster came with a host of challenges. For a while, we were unable to figure out why my jobs did not always consistently run as expected. The main issue was that in order to use TensorFlow and run the Kipoi models, we needed to be on a RedHat7 machine. Additionally, for a while all of our jobs were set to output their stdout’s to the same file. This wasn’t an issue when the first couple jobs fired, but once the file filled up, all of our jobs would mysteriously stop without notifying us that something was wrong. As a result of this, for a while our variant effect size data sets were actually incomplete, with a large portion of the files only containing predicted variant effects for the first 10 million base pairs rather than for the whole chromosome. It wasn’t until we performed some downstream analysis that we realized there were some inconsistencies, and we had to go back in fix it. In general, building this data set takes about 4 days on the cluster using compbio-lab priority. If all jobs fired simultaneously, it would take 2 hours.

Analyzing the Variant Effect Sizes:

Once we had the whole dataset as a set of .vcf files, we found it far easier to work with once it was loaded into a pandas dataframe then saved into a pickle file. One analysis that we ran on the data was to analyze each genomic location using a dimensionality reduction technique. For each genomic location, we essentially have an extremely long vector of scores, one for each different transcription factors. In order to understand this data, we can use a variational autoencoder to reduce the dimension to one that's easier to work with.

It's fairly easy to set up a variational autoencoder in PyTorch, especially working off of existing implementations. To start, we first create the encoding layers of the network. This consists of a small intermediate dense layer, followed by two layers, one to represent the means and one to represent the variances of each of our latent dimensions. Together, we can use these values to sample values from all latent dimensions if we wanted to generate new training data. Then, we construct the decoder network. It takes in all of the latent inputs, and passes them through a dense intermediate layer, then through another dense layer to output vectors the same size as the original inputs.

Now the key detail is that we adjust the loss function so that instead of minimizing binary cross-entropy or some other traditional machine-learning loss function, we're going to minimize the KL divergence of our original data from the outputs our network generates. As discussed in the background section, if we can push this KL divergence towards zero, we will have exactly represented the original data in a lower-dimensional space. This is ideal, as then we have performed the dimensionality reduction we originally wanted to perform. In our case, we set the latent dimension to 32.

Once we have reduced the dimension, we would like to visualize these results. To do this, we used two different approaches. The first was to perform t-SNE and reduce the dimension to exactly two. Because of run-time constraints, we were unable to perform t-SNE on the entire dataset, but we were able to run it on a single chromosome. Additionally, we performed PCA to visualize the data using another traditional machine-learning strategy. Both of these dimensionality reduction techniques were extremely easy to perform out of the box using the sklearn Python package.

Building the GWAS, Transcription Factor, Genomic Location Tensor

The next major portion of this project was creating the tensor which relates genomic location, transcription factor, and GWAS phenotypes. The goal was to get a score which related how much the GWAS z-scores in a particular genomic region matched with the variant effect sizes for each transcription factor. We wanted one score for each 100 kbp region, transcription factor, GWAS phenotype combination. In general, creation of this tensor was an absolute beast of a computational job, as the overall scale of the data was massive. The main components necessary to perform this were selection of the GWAS traits we wanted to analyze, building out a framework for performing all the computations in parallel, and finally determining a scoring mechanism for how to relate 100 kbp regions of z-scores with 100 kbp regions of transcription factor variant scores.

Selection of the GWAS Traits

In general, we did not have an extremely rigorous methodology for determining what GWAS traits were selected, but we did have some goals/selection criteria in mind when we were deciding which ones were worth investigating. In general, we started out by deciding that we

wanted to choose a number of different traits, this way we'd have enough to know that our methodology was actually producing real results. We also chose the traits in a way that many of them were similar to one another. There's a bucket of traits all related to the brain/depression/mental health, some that are related to weight/fitness, some related to addiction, and more.

This is a sampling showing 40 of the 131 GWAS traits that were selected:

100007_raw	Polyunsaturated fat
100008_raw	Total sugars
1001	Duration of strenuous sports
100150	Drinking water intake
100240	Coffee consumed
100580	Alcohol consumed
1011	Frequency of light DIY in last 4 weeks
102_raw	Pulse rate- automated reading
102330	Sweets intake
1031	Frequency of friend/family visits
104400	Fruit consumers
104900	Time spent doing vigorous physical activity
1050	Time spend outdoors in summer
1090	Time spent driving
1100	Drive faster than motorway speed limit
1110	Length of mobile phone use
1160	Sleep duration
1180	Morning/evening person (chronotype)
1190	Nap during day
1200	Sleeplessness / insomnia
1239	Current tobacco smoking
1558	Alcohol intake frequency.
1707_1	Handedness (chirality/laterality): Right-handed
1727	Ease of skin tanning
1717	Skin colour
1757	Facial ageing
1920	Mood swings
1930	Miserableness
1940	Irritability

1950	Sensitivity / hurt feelings
20001_1001	Cancer code- self-reported: lung cancer
20001_1003	Cancer code- self-reported: skin cancer
20001_1006	Cancer code- self-reported: larynx/throat cancer
20001_1020	Cancer code- self-reported: large bowel cancer/colorectal cancer
20001_1002	Cancer code- self-reported: breast cancer
20001_1044	Cancer code- self-reported: prostate cancer
20001_1048	Cancer code- self-reported: leukaemia
20001_1059	Cancer code- self-reported: malignant melanoma
20002_1065	Non-cancer illness code- self-reported: hypertension
20002_1066	Non-cancer illness code- self-reported: heart/cardiac problem

In total, we manually chose 131 different GWAS traits from the nearly 4000 traits that researchers collected data about in the UK biobank. Once we had this list of traits, it was relatively straightforward to pull the GWAS summary statistics from the Neale Lab repository to the cluster. The full list of all GWAS traits are in the appendix.

Creating a Framework to Build the Tensor in Parallel

In general, there were about 1 million genomic locations per chromosome, 23 chromosomes, 500 transcription factors, and 131 GWAS traits. In order to build out the final tensor, we needed to take the Cartesian product of all three, meaning the overall scale of the job we were running was absolutely massive. Again, a parallel solution was needed.

The general framework stayed very similar to the previous big data operation we performed, however we basically needed to do every operation 131 times. In essence, a wrapper script would submit a job for each transcription factor, which contained sub processes for each of the different chromosomes. This would create a worker for each transcription factor/chromosome pair, which would then run a Python script for each of the 131 GWAS traits to compute all of the scores. Each worker needed to first read in the GWAS information pertaining to its assigned

chromosome, read in the variant effect sizes, and then join these two datasets. Once they were matched up, we needed to convert a vector of z-scores and a vector of transcription factor binding-affinity score differentials into a single score. We investigated the use of two different scoring functions, a dot product, and also an LD-score corrected inner product. When we were using the second strategy, we also needed to read in individual level 1K Genomes data.

Running this job was incredibly resource-intensive. In general, it takes 8 gigabytes of ram and 120 hours to run the LD-score corrected strategy on each chromosome for each transcription factor. This means that in order to run the full job on every single chromosome, we needed to run 11500 workers for 120 hours each. Unfortunately, with the priority that we have on the Broad cluster, we cannot get access to that many machines simultaneously. At most we were able to run about 150 machines simultaneously before we started to get throttled, which meant that the task itself was going to take half a year to complete. As a result, we were unable to perform this computation for every single chromosome, and instead ran the LD-score corrected dot product strategy only on chromosome 1.

Creating a Framework to Build the Tensor in Parallel

In order to build the tensor in parallel, we repurposed the same paradigm that we used to build the variant-effect database. Every combination of chromosome, transcription factor, and GWAS trait could be run in parallel with no computational overlap, so again the general framework of the problem was “embarrassingly parallel.” We decided to submit one job per transcription factor binding affinity model and per chromosome, which would then run all of the analysis. Each job consisted of running a single python file, which took as input the model name, gwas file number, and then the chromosome number.

As a reminder, the variant effect sizes are all stored in a database of .vcf files, while the GWAS files are simple .csv files. In order to link the two databases, we simply perform an inner join between the .vcf files and the GWAS files on the tuple of (location, reference allele, alternative allele). Just in case the reference allele and alternative allele are swapped in the two datasets, we also concatenate this dataset with the one created where the inner join between the reference allele and alternative allele are swapped. This gives us a single dataset with a vector of z-scores from the gwas file, and a vector of variant effect scores. For every 100kbp region we want to aggregate these two vectors into a single score. We used two different strategies to perform this aggregation. The first was the dot product strategy, the second was to use a modified inner-product which corrects for LD biases.

The dot product strategy was very simple to run and implement, however its run-time was quite slow. Most of the overall runtime of this strategy was based around reading in the two files and performing the inner-join, but once this was completed it was not at all time-intensive to compute a single dot product for each 100kbp region. The runtime for the LD-corrected strategy was even worse, as for every 100kbp region, we need to perform all of the same steps as in the dot product strategy, and also read in the appropriate data from the 1k genome project. In fact, the first implementation of this strategy was so non-performant that we thought we might need to completely drop the LD-corrected approach.

Our original strategy was to read from the 1k genome .vcf files, as we already understood an easy way to read the data into python. However, this turned out to be incredibly slow because the pyvcf package was very inefficient at parsing each line. Upon googling, we realized this was an issue specifically with reading in the 1k genome .vcf files. We started to develop our own

custom parser to extract just the relevant information we needed, however we soon realized that there was a better way: PLINK.

PLINK is an open-source whole-genome-wide association analysis toolset.^{xxi} Although we did not use it for any of its GWAS/LD software toolsets, we did utilize the fact that the file-formats that it uses to store data are incredibly efficient and easy to parse. On the Broad cluster, we already had all of the requisite files saved in the PLINK file format. Then, using the `pandas_plink` package, we were able to read in the `.bim` and `.bed` file formats into `pandas` dataframes. The `.bim` file contained information about the individual variants, while the `.bed` file contained all of the individuals' genotypes. This switchover to PLINK greatly reduced our runtime, and now it was computationally feasible to build this dataset.

Now that we can (somewhat) quickly read in the 1k-genome data, we can actually build the dataset using the LD-corrected strategy.

The LD-corrected strategy

We need to correct for linkage disequilibrium to ensure that we're not overfitting for a single causal effect in each of our small windows of DNA. To do this, we're going to modify our scoring function so that it is now a dot product adjusted with a coordinate transformation. In order to compute this appropriate coordinate transformation, we're going to construct an R^{-1} matrix for each of the small genomic windows. Assuming the same assumptions made in the LD-regression paper, when we introduce an R^{-1} term into our scoring function, we correct for linkage disequilibrium biases.

First, we load in a small region of the 1k genomes data into the matrix A . Rows are variants, columns are individuals, and entries are the genotype of that individual in a particular position. We start by taking the Singular Value Decomposition of A to get:

$$A = U\Sigma V^T$$

Then, in order to avoid the issue of small entries in Σ being hugely important in the final computation, we filter Σ to create a Σ' by zeroing out all of the entries that explain a negligible amount of the overall variance. Once we've avoided this pitfall, we simply perform the final computation to get R^{-1} :

$$R^{-1} = V^T(\Sigma')^{-2}V$$

Once we have R^{-1} , it is easy to compute our final score. The final score s is computed very easily:

$$score = z^T \cdot R^{-1} \cdot v$$

where z is the vector of z scores

Runtime Considerations for Building the Dataset

For each of the 500+ transcription factors, we needed to run a pretty computationally-intensive script: for each disease we load in a GWAS file, load in a variant scoring file, merge the two datasets, then perform some kind of inner-product computation. If this operation were a simple dot product, then we were done with computationally-expensive operations, but if we wanted to correct for LD-biases, we needed to then repeatedly load in appropriate portions of the 1k-genome data, and then perform a computationally-cheap matrix transformation.

In general, it took about 5 days for one machine to run one chromosome across one transcription factor using the dot product strategy, while it took closer to 8 days for one machine to run one transcription factor using the LD-corrected strategy. Again, getting enough machines to run simultaneously was a large issue because of limited cluster-resources. We could get about 150 machines running simultaneously with sufficient resources, so overall it took close to 4 weeks to build the LD-corrected matrix. The unfortunate consequence of this long build time meant that iteration cycles were very slow.

Chapter 3 – Dataset Analysis Results

Now that the datasets are built, we can begin to analyze them. We have two datasets to explore at this point: the variant effect database and the overall tensor. On the variant effect database, we performed a correlation analysis and an autoencoder/tSNE analysis. On the tensor we explored Latent Dirichlet Allocation as well as performed a flatten analysis.

Variant Effect Database Results:

The two main analyses that we performed on the variant effect database was a correlation analysis to determine the relationship between the transcription factors and various GWAS traits as well as an Auto-encoding/TSNE visualization to determine whether or not there were any latent classes of transcription factors.

GWAS Correlation Analysis:

As discussed in Chapter 2, we ran a simple (but computationally-large) analysis to determine how the variant-effect scores across transcription factors related to the z-scores of various known GWAS traits. In essence, we took the vector of variant effect scores and correlated them with the vector of GWAS z scores across every pair of transcription factor, GWAS pair. Across a dataset linking 500+ transcription factors with 100+ traits, we found that all correlations were extremely small. This is most likely noise, but we can investigate a little further to determine whether any structure exists within the data.

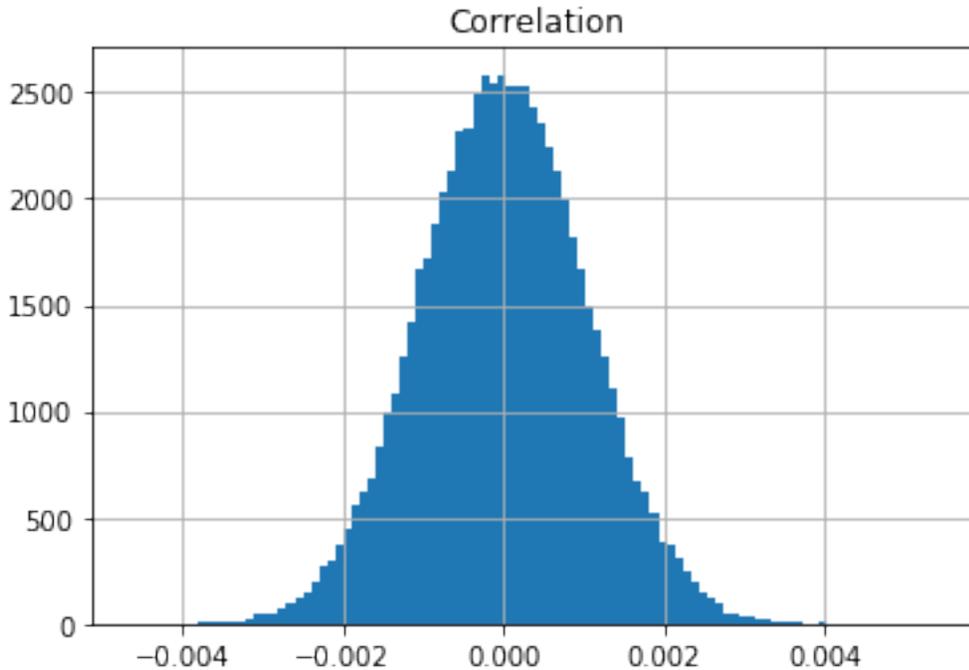


Figure 4: Distribution of Correlation values between GWAS z-scores and transcription factor effect sizes across chromosome 1 shows lack of outlier correlations, indicating that the relationship between TF binding and GWAS is not immediately visible.

It makes sense that the correlation coefficients would be extremely small across the entire genome, as most of the data is in fact noise. The z-scores for a GWAS trait are pretty meaningless except at high thresholds, and the variant effect sizes are mostly random except at key locations. We can then take these correlations across all disease, GWAS pairs and then plot a heatmap of the different correlation values. This is going to look mostly random, but if we start with a clustering step, we might ideally observe a little more structure in the correlation coefficients. Unfortunately much of this still looked like noise, even after using UPGMA to perform clustering on both our GWAS and TF axes.



Figure 5: The pairwise correlation coefficients for each pair of disease and transcription factor does not immediately show any structure, indicating that the relationship between the two is not immediately visible. Hierarchical clustering was performed on both the transcription factors and the GWAS traits using UPGMA to get the ordering along both axes.

As seen in this image, there does not seem to be much structure within this overall heatmap. One way to validate whether or not this clustering approach is actually capturing relationships within the data is to see if the clustering itself relates to our understanding of how the diseases we've been studying relate to one another. In theory, the nearby diseases along the left half of this plot should be similar. There are certain instances where this does actually seem to be true. Some examples include:

- 1 . Saturated fat and sugar consumption are adjacent
- 2 . Ever self-harmed, heard unreal voices, is schizophrenic, experiences panic attacks, and friendship satisfaction are all adjacent
- 3 . Ever taken cannabis and maximum frequency of cannabis use are near one another
- 4 . Basal metabolic rate, body fat percentage, and BMI all appear together
- 5 . Loneliness, guilty feelings, neuroticism, mood swings, miserableness, and irritability all appear together.

Although it isn't as perfect as we might like, a lot of our GWAS traits do seem to appear together in this clustering. In general, the fifth bullet point above was by far the strongest obvious association we noticed within this clustering, as all of the traits clearly relate to one another and there was a block of so many adjacent together.

Autoencoder and TSNE Results:

As discussed in the methods section, there were two different ways we could run the autoencoder and TSNE analysis. The first was to assume that each of the different SNPs was a vector of 500 scores, while the second was to assume that each of the different transcription factors was a massive vector of 100000 scores (one for each snp).

The main hyperparameter that we could tune in our model was the latent dimension of our variational autoencoder. Figure 6 shows the results for a latent dimension of 32.

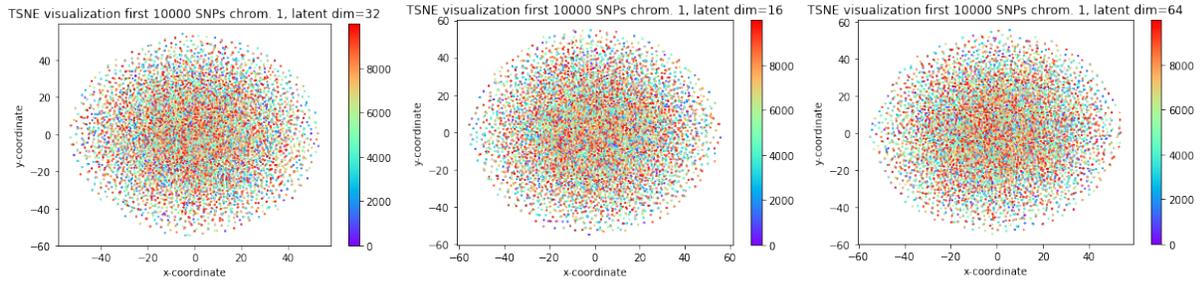


Figure 6: Viewing each SNP as a vector of transcription factor variant effects, we perform autoencoding dimensionality-reduction followed by t-SNE visualization. Across all different settings of the hyper-parameters, we do not observe any observable structure in the plot based on genomic location.

Clearly this is pretty noisy, which is non-ideal. We tried bumping the latent dimension of the variational autoencoder so that it would be able to capture more nuances within the data, but we get a similar plot. Similarly, dropping the dimension so that the variational autoencoder learns more higher-level features still does not yield any perceivable clustering (figure 7).

Now, we adjust the direction of our analysis so that we can start to observe relationships between the transcription factors. Now, we consider each transcription factor to be a massive vector of hundreds of thousands of scores, one for each variant. Since there are only ~ 500 transcription factors in our dataset, we run into the $n \ll d$ problem, where the dimensionality of our dataset is far higher than the number of data points we use. Unfortunately, this means that our variational autoencoder will no longer work, as the KL divergences will blow up when the learned variances get too high. However, we can still perform a similar analysis as before, just by using a regular autoencoder this time. Running an autoencoder with a latent dimension of 64 yields a more meaningful plot this time. We notice that the transcription factor data does actually tend to cluster into some clearly observable areas (figure 8).

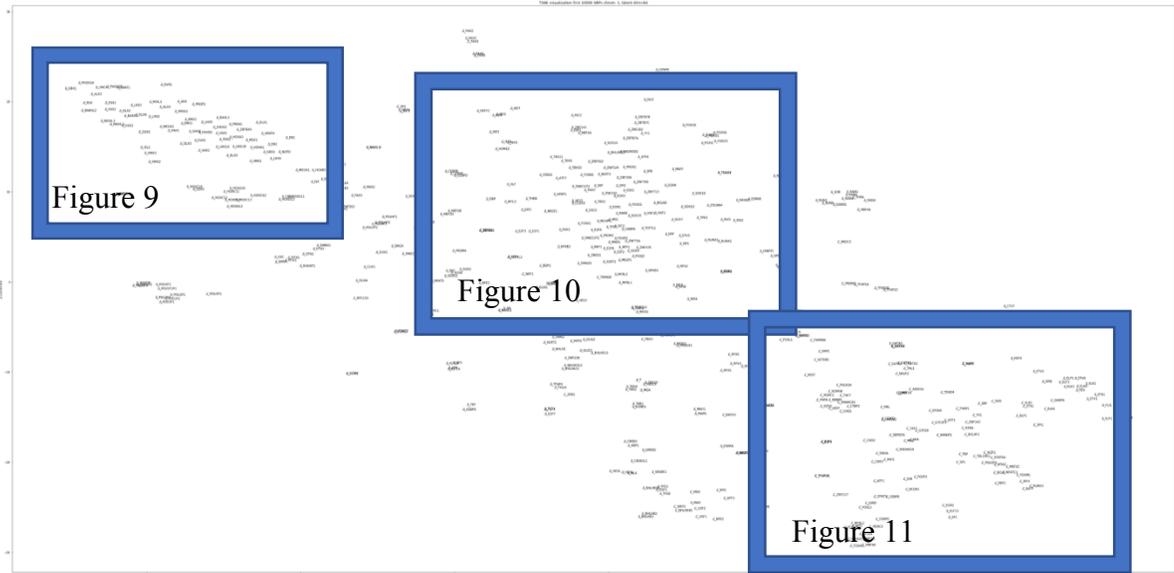


Figure 7: We now consider each transcription factor to be a vector of variant effect-sizes across all different SNPs, and perform an autoencoding dimensionality reduction followed by *t*-SNE visualization. Each labeled point is a transcription factor, with the first letter representing whether or not the model was built using a SELEX (S) experiment or a ChIP-seq (C) experiment. In the next couple figures we zoom into three noticeable large clusters.

The first thing we notice when we actually dive into this plot is that the models built using ChIP-seq data cluster together in the bottom right corner of this figure. This reveals that there are actually noticeable differences between the scores outputted by machine-learning models built on the two different datasets. Zooming in on the cluster at the upper left corner allows us to construct the plot in figure 9. Most of these transcription factors are homeobox transcription factors, an important subgroup of transcription factors which relate to the early formation of many body structures during embryonic development. This is exactly the result we would hope would appear when we cluster transcription factors together using *t*-SNE dimensionality reduction. The C_ prefix refers to ChIP-seq experiments, the S_ prefix refers to SELEX experiments.

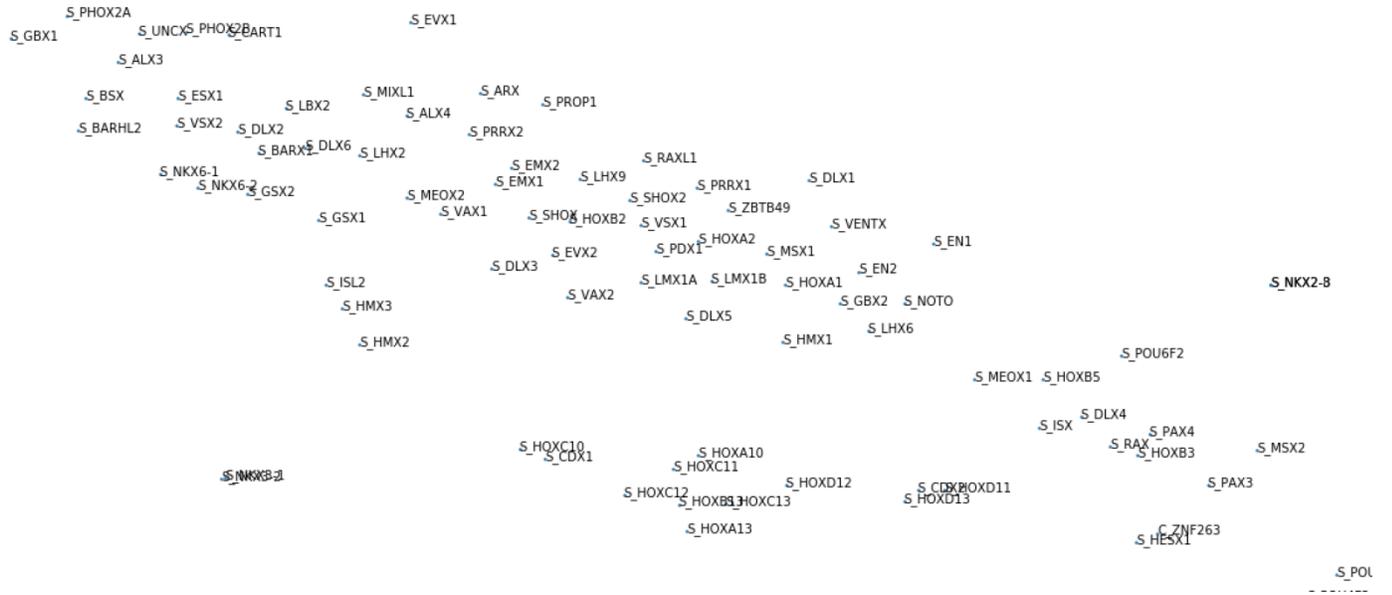


Figure 8: Zoomed in view of the upper left corner of figure 8 (t-SNE plot of transcription factors). This cluster is well structured as we see that many of the transcription factors have similar names, and also that most TFs are homeobox transcription factors, a class of transcription factor related to embryonic growth.

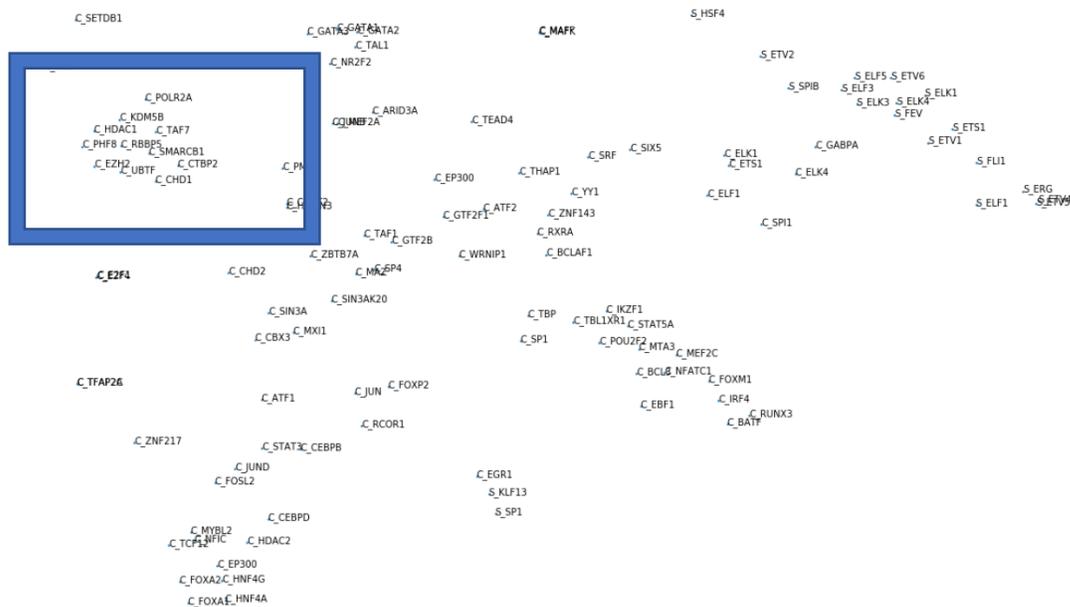


Figure 9: Zoomed in view of the bottom right corner of figure 8. This cluster is also well-structured: we see that all of the transcription factor models built using ChIP-seq data cluster together. The upper left corner of this picture shows a dense cluster of transcription repressors.

In figure 10, we notice that many of the FOX class transcription factors cluster together. The right side of this plot shows that the ELK and ELF classes of transcription factors also pop up near one another. The super dense cluster in the upper left corner of this plot consists of a group of transcription repressors. It's pretty interesting to see that the small subset of transcription factors which repress gene expression actually cluster together using this method.

The center of our original figure features a very large cluster:

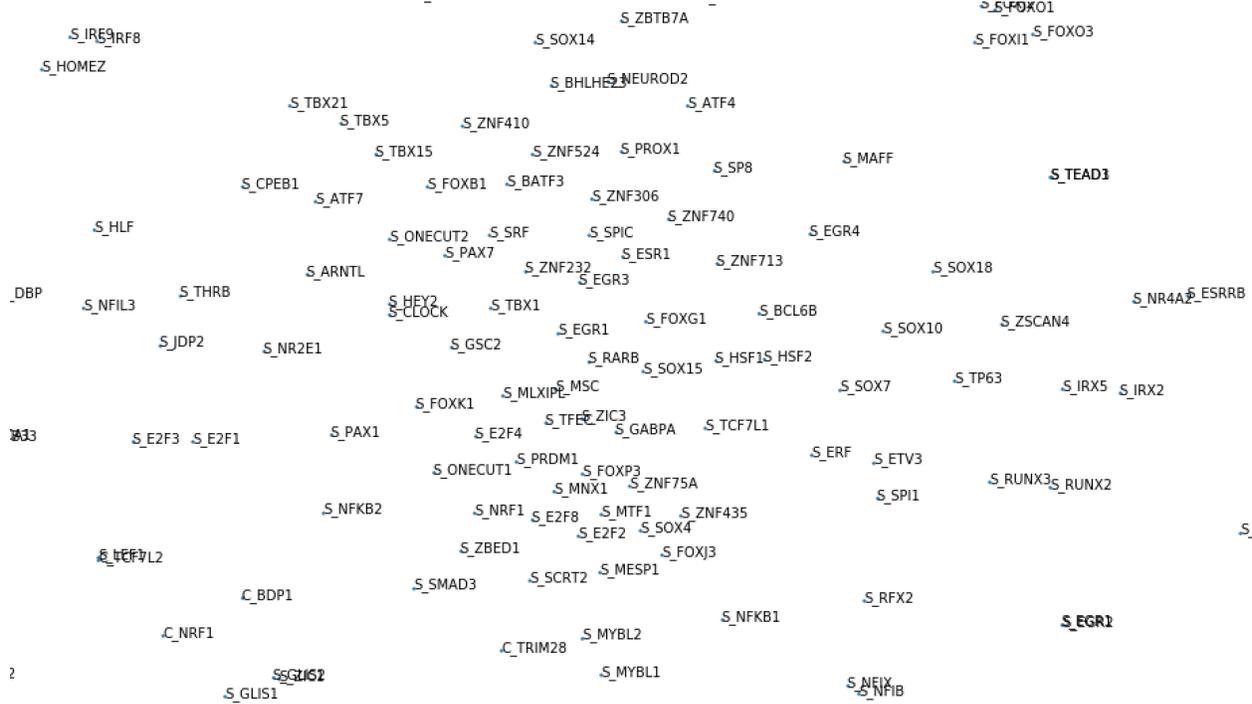


Figure 10 Zoomed in view of the middle of the t-SNE representation of the transcription factors. We notice certain small patches of similar transcription factors (FOX and SOX), but it's difficult to tell if there's an overall unifying theme behind this cluster.

This cluster features many members from the zinc-finger class of transcription factors. We also see a lot of the heat shock factor transcription factors, as well as transcription factors from the FOX and SOX class of transcription factors. It's tough to fully parse what links these transcription factors, but it is nice to see that similar proteins often appear very close to one another in this plot.

GWAS, Location, TF Tensor Analysis:

We performed two main analyses on the GWAS, Location, TF tensor. The first was an LDA analysis, while the second was flattening the tensor and running standard dimensionality

reduction techniques. In this section, we'll focus on the analyses performed on the LD-corrected dataset. Equivalent plots for the dot product strategy will be left in the appendix.

Latent Dirichlet Allocation on the Tensor

In general, the first step of our analysis was to determine what were the words and what were the documents. There's 6 different ways we could slice the data: we could let documents be any of TFs, genomic locations, or GWAS traits, or we could let words be any of TFs, genomic locations, or GWAS traits. Let's work through an example of this to understand what the results of this analysis are.

Say we let each of the genomic locations be a document. Then, the words are going to be the presence or absence of a particular GWAS/TF tuple. We say that a GWAS/TF tuple is a word in a particular document if the absolute value of the score in that genomic location, GWAS, TF tuple is above some threshold we specify. This threshold can be set based on any number of factors, but our general strategy was to use it as a hyper-parameter we tuned so that the downstream analysis would be better. We used absolute value of the score instead of regular value because a particular GWAS/TF pair is related to a genomic location if the score is very large (so the GWAS and TF activities were positively correlated in that location), or if the score is very small (GWAS and TF activities negatively correlated.)

Now, we might ask what are the results of this toy analysis. In general, an LDA analysis returns the posterior of two different hidden variables. The first is the topic distribution of each document, which is essentially a soft clustering of all the documents into various classes. In our toy example, this means that for each genomic location we get a vector of scores representing how much that genomic location is in each of k classes. The number of classes is also a hyper

parameter we can tune. Additionally, the algorithm also outputs a vector of word-probabilities, representing how likely a given word is to appear in each of the different classes.

We ran the model in all 6 directions, but only two of them yielded information that we could use for follow-up analysis. The first was performing clustering on every genomic location. In order to determine the number of clusters that we used, we started by computing the model perplexity of a held-out test set of data across a number of different cluster sizes. Perplexity is a measure of how good the model fit is, with a lower perplexity representing that the model generalizes well to a test dataset. Unfortunately, the perplexity strictly decreased as we increased the number of topics in our model. In the end, we settled on using 10 different latent classes in order to make the model more interpretable.

The other direction that yielded interesting results was making each of the words be diseases. Then, the posterior probability of each word given the class essentially acts as a clustering of diseases. This is pretty interesting to look at, as because LDA was originally a topic model, it's nice to see that our topics consist of human-readable words. In order to ensure that our topics were easy to parse for humans, we used a smaller number of topics this time.

Now, we can again plot the topic distributions across each genomic location by now marginalizing over our document distributions. In this analysis, our documents are tuples of genomic locations transcription factors. In order to determine the topic distribution of one of the genomic locations, we simply average all of the topic distributions for documents containing the genomic location.

LDA Results – Documents are genomic locations:

To begin with, we first ran an analysis where the documents were genomic locations, and the words were tuples (transcription factor, GWAS) tuples. There are two components to the output of LDA. The first is it outputs a document breakdown which specifies a soft clustering of each document into topics. This soft clustering revealed that at very low threshold levels (threshold < 20), our topic model was not really differentiating between the different topics for each disease. With a small number of classes, we found that the distribution of topics was quite uniform when our threshold was set too low.

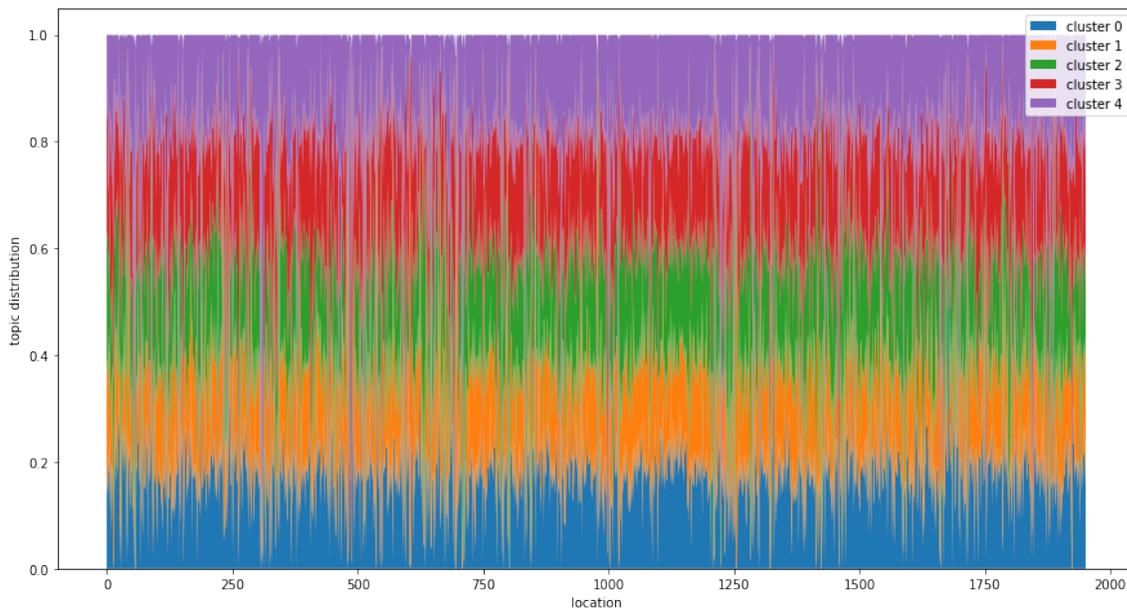


Figure 11: Clustering of each genomic location (x-axis) according to the GWAS variant perturbation of TF motif scores (words) and assignment of each location into a set of topics (colors) defined by their word frequencies. Topic breakdown using LDA for threshold=20 across all genomic locations in chromosome 1. The distribution on each genomic location is mostly uniform, showing that we need to increase our word threshold.

If we bump up the threshold higher, we get a much spikier plot, where now each genomic location is mostly dominated by a single cluster. We also overlaid a genomic annotation track

consisting of whether or not each 100kbp region had an enhancer present. These are the black lines in our figure. The large middle region in the figure only appear this way because of a lack of data in that region, so it's not worth reading into. We chose to use 10 clusters because it increases the interpretability of our results.

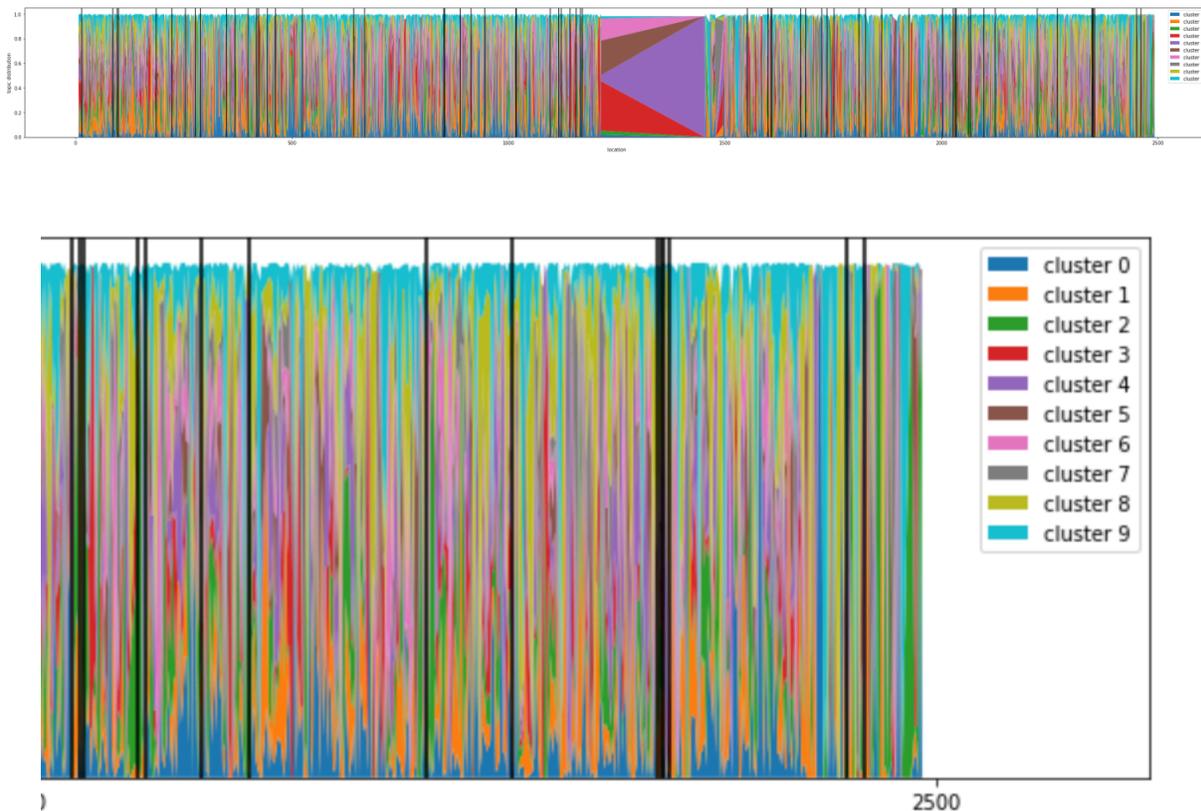


Figure 12: (Top) Re-clustering of chromosome 1 genomic locations by GWAS-TF effects (as in Figure 11) now using a word threshold value of 50 on the LD-corrected dataset. (bottom) zoomed in view of the last 500kbp region.

The next question we might ask is what do these clusters mean? Do they relate to any known diseases? The second result of an LDA analysis is that we are able to determine the word distribution for each topic. In this case, our words were transcription factor, GWAS pairs. So, one approach we can use to better understand our different topics is to marginalize over each of

our transcription factors / GWAS conditions tuples to determine the conditional probability of each GWAS given the topic. Another analysis we can run is marginalizing to get the probability of each transcription factor appearing given the topic.

When we do this on our GWAS traits, we notice a very unfortunate conclusion: all of the topics contain the exact same list of top diseases. At the top of every single topic appear the same traits every time: BMI, Weight, Platelet Count, Basal Metabolic Weight, and a few more. These are all very important GWAS traits, so perhaps it is not surprising that they appear at the top of every single cluster. Most clusters should be predicting these traits to be common words in our corpus, as these traits essentially function as stop-words. One approach we can use to get better results is to completely remove these stop-words.

Stop Word Removal

One potential explanation for this uniformity of topic distributions is that perhaps the most common words in each topic are actually our LDA use-case's version of stop words. In conventional topic modelling use-cases of LDA, we usually perform a pre-processing step to remove the most common words that appear in the English language. Words such as "and", "the", and "or" are so common that they're going to appear at the top of every topic simply because of how highly probable it is that a word is going to be one of these stop words. Generally, the way this is handled by statisticians is by filtering out a big list of stop words. Unfortunately that is not going to be possible for our particular use-case, so instead we're going to need to use a smarter method of removing these stop words.

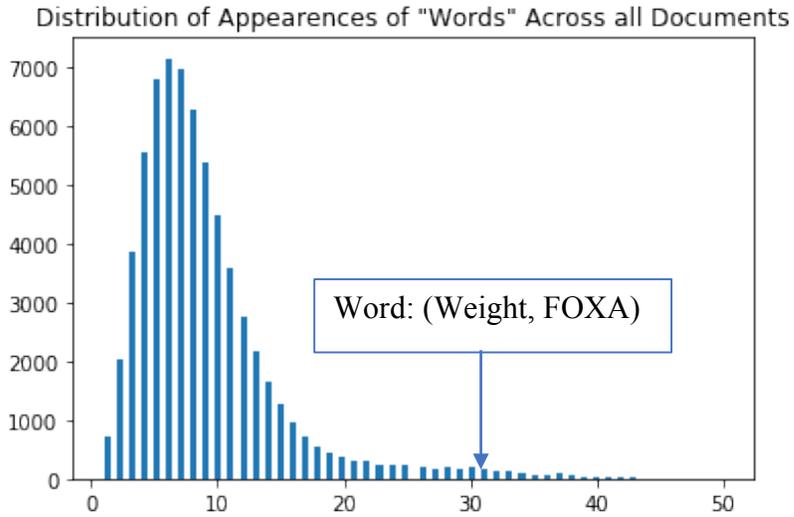


Figure 13: Distribution of number of genomic location (document) appearances of each TF-trait tuple (word). Because this distribution has a fat tail, a small number of words appear in far more documents than others. These are stop-words.

As we see in figure 13, most words don't appear in many documents; however, quite a few words seem to appear in a massive number of documents. Looking into these words reveals that we've found the source of our problem. Here are some of the top words and their counts:

Trait	TF	Evidence	Count
Basal metabolic rate	MEF2D	Selex	N=50
Basal metabolic rate	NFKB2	Selex	N=49
Basal metabolic rate	SOX21	Selex	N=49
Basal metabolic rate	HNF1A	Selex	N=48
Basal metabolic rate	DLX5	Selex	N=47
Basal metabolic rate	RCOR1	ChIP	N=47
Basal metabolic rate	MAX	Selex	N=47
Weight	PAX3	Selex	N=46
Weight	HEY1	Selex	N=46
Basal metabolic rate	TCF3	Selex	N=46
Basal metabolic rate	ZNF263	ChIP	N=46
Basal metabolic rate	TCF4	Selex	N=46
Basal metabolic rate	POU3F2	Selex	N=46
Platelet count	USF1	Selex	N=46
Weight	MEF2D	Selex	N=45
Weight	POU5F1P1	Selex	N=45
Weight	MLX	Selex	N=45
Basal metabolic rate	GBX2	Selex	N=45
Weight	PRDM1	Selex	N=44
Basal metabolic rate	PRDM1	Selex	N=44
Basal metabolic rate	EP300	ChIP	N=44
Basal metabolic rate	NKX3-1	Selex	N=44
Basal metabolic rate	PKN0X1	Selex	N=44
Basal metabolic rate	ESRRB	Selex	N=44
Basal metabolic rate	ZNF740	Selex	N=44
Platelet count	ZIC4	Selex	N=44
Platelet count	SREBF2	Selex	N=44
Weight	HNF1A	Selex	N=43
Weight	RFX3	Selex	N=43

All of the conditions in this list are the same as the ones that appear at the top of every topic distribution. If we consider the histogram in figure 13 as a probability distribution, the elbow of its CDF appears at around 15. Let's filter out all of the words that appear in more than 15 documents.

Unfortunately, while this does fix the problem of conditions like BMI, Weight, and Platelet Count being in every list of topics, it does not fix the overall issue. Although these topics get removed almost altogether from the list of marginalized top diseases for each topic, we now find that a new set of top diseases appear at the top of each topic. It's definitely less uniform, which is good, but similar diseases are still appearing in each list.

Topic 0	2.302E-03	Neuroticism score
	2.210E-03	Non-cancer illness code- self-reported: asthma
	2.077E-03	Fluid intelligence score
	2.019E-03	Morning/evening person (chronotype)
	1.881E-03	Sleep duration
	1.867E-03	Irritability
	1.862E-03	Lifetime number of sexual partners
	1.686E-03	Alcohol intake frequency.
	1.625E-03	Nap during day
	1.613E-03	Facial ageing
	Topic 1	2.498E-03
2.006E-03		Morning/evening person (chronotype)
1.961E-03		Non-cancer illness code- self-reported: asthma
1.792E-03		Mean time to correctly identify matches
1.714E-03		Non-cancer illness code- self-reported: hypertension
1.669E-03		Fluid intelligence score
1.542E-03		Alcohol intake frequency.
1.530E-03		Irritability
1.511E-03		Sleep duration
1.489E-03		Mood swings
Topic 2	2.143E-03	Morning/evening person (chronotype)
	2.078E-03	Alcohol intake frequency.
	1.744E-03	Mean time to correctly identify matches
	1.742E-03	Non-cancer illness code- self-reported: hypertension
	1.738E-03	Nap during day
	1.737E-03	Current tobacco smoking
	1.694E-03	Lifetime number of sexual partners
	1.680E-03	Time spend outdoors in summer
	1.671E-03	Miserableness
	1.668E-03	Non-cancer illness code- self-reported: asthma
Topic 3	2.227E-03	Facial ageing
	2.207E-03	Sleep duration
	1.913E-03	Fluid intelligence score
	1.870E-03	Nap during day
	1.863E-03	Morning/evening person (chronotype)
	1.813E-03	Neuroticism score
	1.751E-03	Current tobacco smoking
	1.661E-03	Skin colour
	1.582E-03	Non-cancer illness code- self-reported: hypertension
	1.577E-03	Miserableness
Topic 4	2.141E-03	Non-cancer illness code- self-reported: hypertension
	2.058E-03	Nap during day
	1.875E-03	Pulse rate- automated reading
	1.854E-03	Leisure/social activities: Religious group
	1.804E-03	Length of mobile phone use
	1.800E-03	Sensitivity / hurt feelings
	1.772E-03	Drive faster than motorway speed limit
	1.631E-03	Sleeplessness / insomnia
	1.480E-03	Current tobacco smoking
	1.477E-03	Facial ageing

One explanation for the reason this happens is that although individual topics will vary pretty substantially based on the individual words that appear, when we aggregate across each disease to get a marginal probability of the different GWAS trait, we're smoothing over the variations that exist within the topics, and instead getting only meta results. The topic distributions also look even spikier, with even less similarities between adjacent genomic regions. (Figure 15). This in general seems to show that stop word removal is pretty unhelpful.

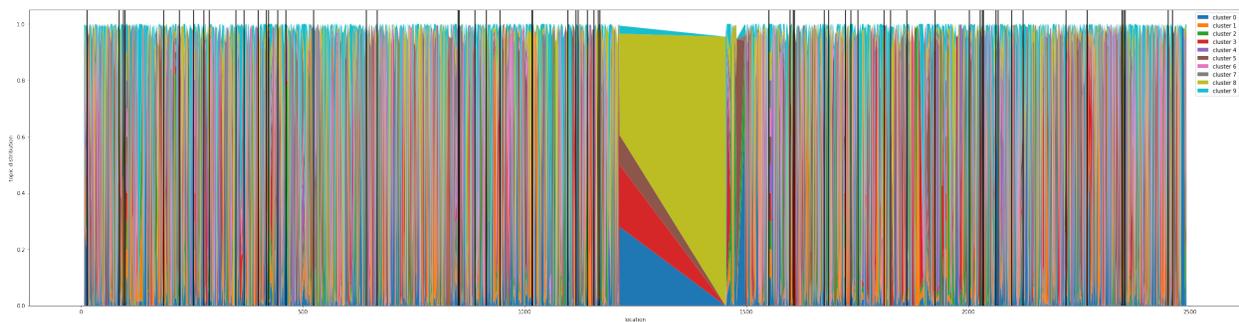


Figure 14: Re-clustering of chromosome 1 genomic locations by GWAS-TF effects (as in Figure 12) after stop-words are removed. We observe that the per-document distributions are even spikier now because each document has very few words. Coupled with the lack of GWAS cluster structure, this shows the results are mostly noise.

Increasing the Number of Clusters

In the previous section, we mentioned that the number of topics we chose for our LDA model was ten. This number was chosen almost completely arbitrarily. One alternative approach we used to get potentially better results was bumping up the number of clusters massively. One measurement of how good our LDA model is at representing our data is the model perplexity. By holding out a sample of the documents then training the LDA model on the rest of the data, we can compute the probability of our LDA model generating each of the documents. The larger this probability, the better our model is. This concept can be understood using a metric known as perplexity, with a lower perplexity representing a better model. Testing the model perplexity for

cluster values between 1 and 1000 revealed that the perplexity decreases as the number of clusters increases without an end in sight, so unfortunately, we were unable to use the minimum perplexity value as a means of choosing the number of clusters.

So again rather arbitrarily, we explored choosing an enormous number of clusters: 100 unique different topics for our LDA algorithm. Strangely enough, two of the topics strongly dominate the document distributions across all the different genomic regions. Diving into the actual topics, we find that one of these topics relates to GWAS traits about size (BMI, Weight, etc.) while the other topic relates to GWAS traits about happiness (miserableness, depression, etc.). The top diseases in all other topics were also about a person's size.

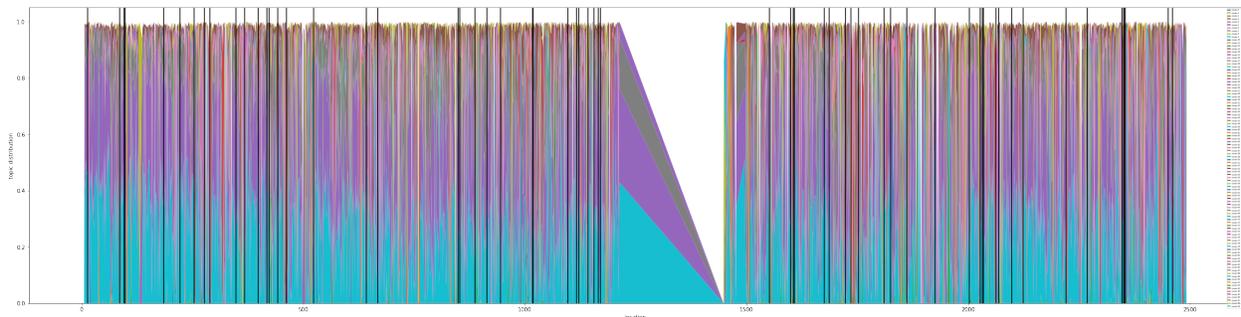


Figure 15: Re-clustering of chromosome 1 genomic locations by GWAS-TF effects (as in Figure 15) now with 100 unique clusters. Two clusters are extremely dominant in the topic breakdowns for all of the documents, while all other topics appear at very negligible frequencies. One cluster contains GWAS traits related to a person's physical size, while the other contains GWAS traits related to a person's general happiness.

LDA Results – Words are Diseases:

There's a second way we can perform this analysis, which is where instead of having each genomic location be a document and each tuple of GWAS, transcription factor be a word, we can flip the script. Now we let each word be a transcription factor, and each tuple of GWAS, genomic location be a document. Again, there are two different outputs of the LDA results, a

posterior topic distribution for each document as well as a posterior word distribution conditional on each topic. The reason why this breakdown is so compelling comes from looking at the word breakdown of each topic. This is the topic distribution for LDA run with a cutoff value of 50 on the LD-corrected dataset with 7 topics.

Topic 0	0.481951	Platelet count
	0.200881	Haemoglobin concentration
	0.025895	Non-cancer illness code- self-reported: high cholesterol
	0.019837	Mean time to correctly identify matches
	0.016221	Recent poor appetite or overeating
	0.015241	Trouble falling or staying asleep- or sleeping too much
	0.014754	Non-cancer illness code- self-reported: hypertension
	0.014031	Recent feelings of tiredness or low energy
	0.012645	General happiness
	0.011955	Ever believed in an un-real conspiracy against self
Topic 1	0.150023	Age first had sexual intercourse
	0.102997	Leisure/social activities: Religious group
	0.084135	Lifetime number of sexual partners
	0.078177	Knee pain for 3+ months
	0.074402	Wears glasses or contact lenses
	0.059822	Chest pain or discomfort when walking uphill or hurrying
	0.058271	Time employed in main current job
	0.051867	Fluid intelligence score
	0.030373	Ever addicted to any substance or behaviour
	0.028028	Ever had same-sex intercourse
Topic 2	0.118415	Frequency of consuming six or more units of alcohol
	0.117539	Amount of alcohol drunk on a typical drinking day
	0.110931	Frequency of memory loss due to drinking alcohol in last year
	0.053892	Substances taken for anxiety: Drugs or alcohol (more than once)
	0.041148	Prospective memory result
	0.038535	Non-cancer illness code- self-reported: migraine
	0.03538	Non-cancer illness code- self-reported: angina
	0.031885	Non-cancer illness code- self-reported: back problem
	0.031256	Non-cancer illness code- self-reported: muscle/soft tissue problem
	0.024762	Non-cancer illness code- self-reported: cardiomyopathy
Topic 3	0.311381	Basal metabolic rate
	0.218323	Weight
	0.17303	Body mass index (BMI)
	0.158437	Body fat percentage
	0.021485	Platelet count
	0.017368	Age first had sexual intercourse
	0.015192	Haemoglobin concentration
	0.006042	Leisure/social activities: Religious group
	0.003994	Non-cancer illness code- self-reported: asthma
	0.003827	Lifetime number of sexual partners
Topic 4	0.112081	Health satisfaction
	0.085921	Friendships satisfaction
	0.036619	Ever addicted to prescription or over-the-counter medication
	0.032858	Non-cancer illness code- self-reported: hayfever/allergic rhinitis
	0.031428	Non-cancer illness code- self-reported: allergy/hypersensitivity/anaphylaxis
	0.029227	Non-cancer illness code- self-reported: alcohol dependency
	0.026534	Mental health problems ever diagnosed by a professional: Obsessive compulsive disorder (OCD)
	0.025698	Ever self-harmed
	0.025371	Ever taken cannabis
	0.025235	Non-cancer illness code- self-reported: multiple sclerosis
Topic 5	0.143063	ECG- heart rate
	0.106225	Neuroticism score
	0.093032	Guilty feelings
	0.062111	Loneliness- isolation
	0.033069	Miserableness
	0.029942	Mood swings
	0.028556	Non-cancer illness code- self-reported: depression
	0.026663	Sensitivity / hurt feelings
	0.023821	Irritability
	0.017388	Haemoglobin concentration
Topic 6	0.125505	Current employment status: Unemployed
	0.093204	Current employment status: In paid employment or self-employed
	0.043178	Longest period spent worried or anxious
	0.042427	Ever physically dependent on alcohol
	0.04058	Ever addicted to illicit or recreational drugs
	0.039553	Risk taking
	0.038806	Ever heard an un-real voice
	0.036233	Bipolar disorder status
	0.031918	Non-cancer illness code- self-reported: stress
0.031789	Ever attempted suicide	

This distribution seems to actually correlate with our understanding of how GWAS traits relate to one another. Topic 3 is very self-contained, with a lot of the GWAS traits related to weight and metabolism appearing at the top of the topic. Topic 5 also has some nice features, especially when we notice that guilty feelings, loneliness, miserableness, and mood swings all appear next to one another. Topic 6 seems to relate to employment, while topic 1 has a lot to do with sexual intercourse.

In order to again create a similar plot of the topic distribution across every location, we now look towards the other part of the analysis, or the breakdown of each document by topic. Again, documents are tuples of locations and TF activity, so we can marginalize across the per document topic distribution to get a probability of topic for each location. This allows us to plot our topic distribution for each genomic location across the whole genome.

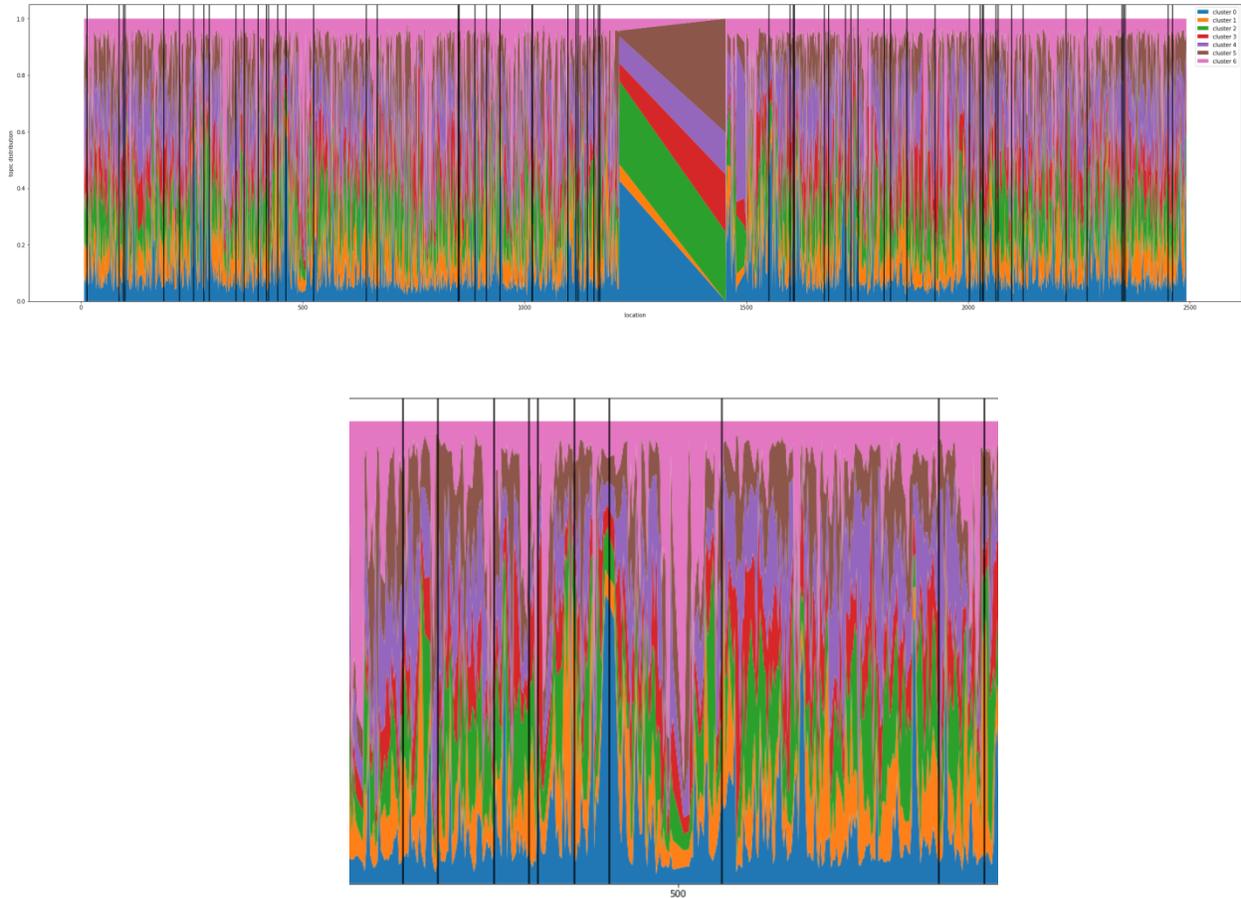


Figure 16: (top) Re-clustering of chromosome 1 genomic locations by GWAS-TF effects (as in Figure 11), under a new structure that documents are genomic location-TF effect tuples, and words are diseases. Each topic distribution for the different genomic locations is the average topic breakdown of all documents containing that genomic location.

(Bottom) Zoomed in view of the top plot. We see that nearby genomic locations tend to have similar topic distributions, a good sign for our overall model.

Story: Inspecting a genomic region of interest

When inspecting the plot of topic breakdowns across each of the genomic locations when we view GWAS traits as words, we noticed a very large region of red (figure 18). This is abnormal, as most of the plot shows that the document distribution of red regions is quite small across most of the genome. The red cluster is topic 3, which contains a lot of the GWAS traits related to a person's physical size.

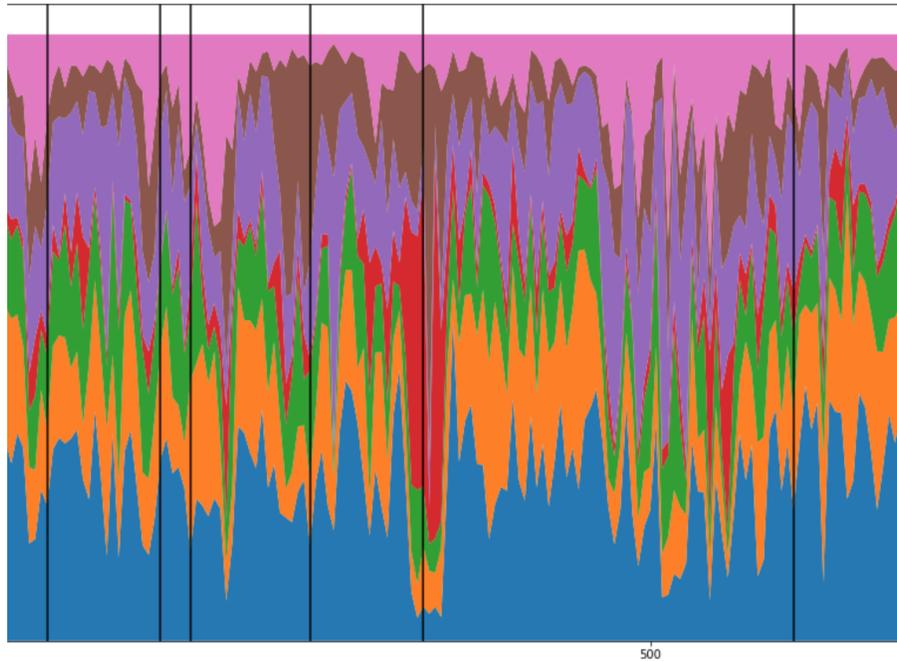


Figure 17: Looking through the plot in figure 16 (clustering of all genomic regions when GWAS traits are viewed as words), we find a large region of red indicating that an area of the genome strongly relates to topic 3. Topic 3 covers weight-related traits.

Another interesting thing is that this region of cluster 3 activity occurs right near an enhancer region. Using the Broad Institute’s Integrative Genomics Viewer, we find that there’s a gene located at that exact location called PIK3R3 (Figure 17).

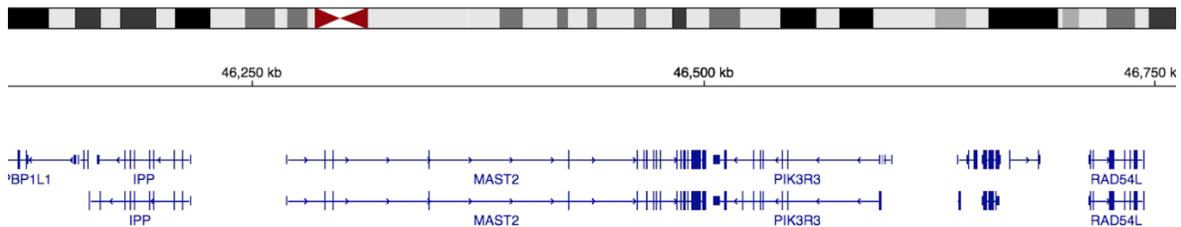


Figure 18: Screenshot of the same genomic region (large red block of cluster 3 activity from figure 18) inspected with the Broad Institute Integrative Genomics Viewer. A number of relevant genes are located in this region, including PIK3R3, a gene related to growth-signaling. This gene might be why genomic regions are tagged with cluster 3 (weight-related GWAS) activity in this region.

The PIK3R3 (Phosphoinositide-3-Kinase Regulatory Subunit 3) gene codes for a protein which serves as a secondary messenger for growth signaling pathways. It binds with protein tyrosine kinases to regulate their activity. In general, the fact that this gene was exactly at the location where we noticed a hotspot of cluster 3 activity is pretty amazing. This LDA technique and visualization has allowed us to detect an important location of DNA within chromosome 1.

Ontology Term Enrichment Using GREAT

Another way we can validate that our discovered clusters are biologically meaningful is by performing ontology term enrichment on the strongest genomic locations associated with each cluster. GREAT is a tool which will efficiently and automatically do this for us.^{xxiii} The way the tool works, is that given a list of genomic regions, the framework starts by first determining the set of all genes associated with the input list. Then it loads a diverse set of gene annotation datasets, each of which consists of a list of topics (or annotations) about the genes. GREAT then answers the following question, do any of these annotations appear in our set of regions statistically more often than they do relative to a background list of genomic areas? Once it calculates whether this is the case, the tool cleanly outputs the list of statistically enriched terms.

Running this analysis is very easy thanks to the web portal constructed by the authors of the GREAT tool. In order to run the analysis, we simply need to input a list of background genomic regions (in this case all 100 kbp regions in chromosome 1), as well as a list of target genomic regions. The next question we need to answer is how do we go from a soft clustering across all genomic regions to a hard list of locations for each cluster? We tried two different approaches. The first approach was to select the top n different locations for each topic based on the topic distributions. The second approach was to determine what the average distribution was for each topic, and then select all genomic locations where the topic distribution was

substantially greater than this background level. It turns out that both approaches yielded very similar results, implying that the ontology term enrichments were robust to adjustments in the hyperparameters / data-analysis strategy.

Of the seven different topics in our topic model, the ontology results were by far the best for topic 3. This makes a lot of sense considering the GWAS traits most related to topic three relate to a person’s overall size and body shape. A good sign that the genomic locations are actually relevant to real genes comes from the very first plot, which shows that our selection of genomic regions related to topic 3 are generally closer to transcription start sites than the background data.

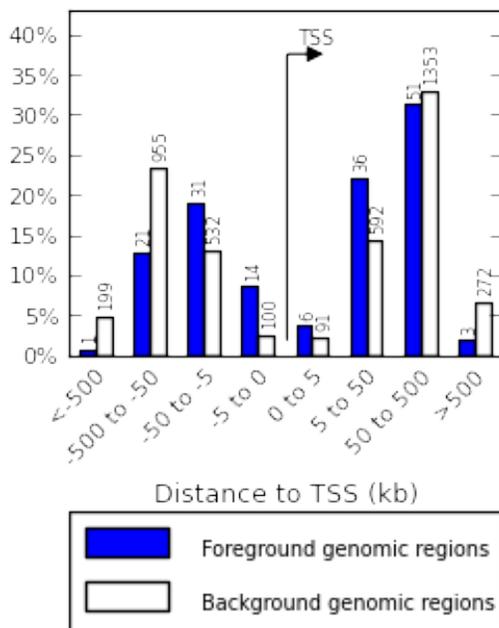


Figure 19: Distance to Transcription Start Site for genomic regions in topic 3 (blue) compared to all 100kbp regions in chromosome 1 (white). In general, the genomic regions in topic three are closer to transcription start sites, implying that our selected genomic regions are more likely to be biologically relevant than the background regions.

Digging in to the actual ontology terms that showed statistically significant enrichment for this cluster, we find that many of the terms are related to catabolic processes (Figure 21).

From the Biological Processes gene ontology dataset, we find that terms like ATP catabolic processes, organic substance catabolic processes, nucleotide catabolic processes, and many more are all significantly enriched. From the MSigDB gene ontology dataset, we find that the “ABC Transporters” term is significantly enriched. ABC transporters are a large class of proteins that help with ATP hydrolysis. The rest of the ontology terms are all very similar to the ones listed.

Term Name	Hyper Rank	Hyper Raw P-Value	Hyper FDR Q-Val
organophosphate catabolic process	1	1.4281e-9	1.4910e-5
cellular nitrogen compound catabolic process	2	1.2301e-8	6.4211e-5
aromatic compound catabolic process	3	1.4066e-8	4.8951e-5
heterocycle catabolic process	3	1.4066e-8	4.8951e-5
organic cyclic compound catabolic process	3	1.4066e-8	4.8951e-5
purine nucleotide catabolic process	6	2.3408e-8	4.0730e-5
nucleotide catabolic process	6	2.3408e-8	4.0730e-5
nucleoside phosphate catabolic process	6	2.3408e-8	4.0730e-5
ATP catabolic process	9	3.0983e-8	3.5940e-5
purine ribonucleoside monophosphate catabolic process	9	3.0983e-8	3.5940e-5
purine nucleotide metabolic process	11	4.5459e-8	4.3145e-5
nucleoside monophosphate catabolic process	12	4.8151e-8	4.1892e-5
purine nucleoside monophosphate catabolic process	12	4.8151e-8	4.1892e-5
nucleobase-containing compound catabolic process	14	5.6342e-8	4.2015e-5
organic substance catabolic process	15	8.8945e-8	6.1906e-5
purine ribonucleotide catabolic process	16	1.0323e-7	6.7358e-5
ribonucleotide catabolic process	16	1.0323e-7	6.7358e-5
organonitrogen compound catabolic process	18	1.1432e-7	6.6304e-5
purine ribonucleotide metabolic process	19	1.8343e-7	1.0079e-4
purine-containing compound catabolic process	20	3.2378e-7	1.6901e-4

Figure 20: Small subset of all the significantly enriched ontology terms for genomic locations in cluster 3. These are the gene ontology terms enriched from the biological processes ontology dataset. We notice that a lot of the terms relate to catabolic processes, a good sign considering cluster 3 relates to BMI / weight GWAS traits.

To reiterate, we used topic modeling on a dataset linking diseases and transcription factors to determine a set of latent clusters behind each region of DNA in chromosome 1. Because of the inherent setup of our topic-modelling algorithm, we knew certain clusters related to specific classes of diseases. One in particular stood out: topic 3 strongly related to weight, BMI, and other size-related traits. Digging into this cluster revealed that the top genomic locations for that cluster overwhelmingly related to catabolic processes and the production of ATP. This is a pretty amazing result.

Unfortunately, none of the other clusters had nearly as striking of results. To begin with topics 2,4, and 5 had very few terms that showed statistically significant enrichment at all. Topic 1 had slightly interesting results. Recall that topic 1 was loosely related to GWAS traits related to sex. One term that emerged as significantly enriched for this topic included prostaglandin F receptor activity. Prostaglandin naturally occurs in semen, and is a hormone related to childbirth. The other term that appeared was abnormal sex gland secretion. Both of these terms definitely relate to the GWAS traits most relevant to topic 1.

Topic 0 also featured some interesting ontology enrichments. The main two GWAS traits in topic 0 are hemoglobin concentration as well as platelet count. The enrichment terms for this trait included “response to human growth stimulus.” In mice (mouse phenotype ontology dataset), the enriched ontology terms included “decreased liver weight”, “decreased food intake”, “decreased lean body mass”, and “abnormal lean body mass”. One of the genes related to this cluster is NEGR1, a gene related to high blood pressure. High blood pressure is in turn related to high platelet counts and hemoglobin concentration. However, we will admit that connecting the significant ontology terms that GREAT detected back to topic 0 definitely requires some out of the box thinking.

The only other topic which had significant ontology enrichment terms was topic 5. Topic 5 involved the gwas traits which were loosely related to sadness and depression. The only enriched terms for the genomic locations related to topic 6 include “RNA processing”, “RNA splicing”, and “mRNA production.” It’s tough to see how this relates to the GWAS traits important to topic 5.

In summary, this table shows the top ontology terms and associated p-value for all the different clusters. Bolded terms are statistically significant. As we see, all terms that appear at the top of one cluster have very low significance for the other clusters. Even when not statistically significant, the top biological process term usually seems to be loosely related to the cluster description. If we no longer restrict our view to just the top ontology term, we end up finding that a lot of the strongest associated ontology terms are in fact related to the cluster descriptions. This is a very good sign that our cluster descriptions are representative of real biological phenomenon.

Cluster	Associated GWAS Trait	Enriched Biological Process							Min P-Value
		Response to Growth Hormone Stimulus	Multi-organism reproductive processes	Response to Mechanical Stimulus	Organophosphate Catabolic Processes	Regulation of Neurotransmitter Secretion	RNA Splicing	Membrane Hyperpolarization	
1	Blood	6.05	0.40	0.40	1.70	0.40	0.40	0.70	9.20E-07
2	Sex	0.40	4.70	1.40	0.40	0.40	0.70	0.40	1.65E-05
3	Alcohol / Risk Taking	0.52	0.40	3.15	0.40	0.40	0.40	0.40	7.32E-04
4	Weight	0.40	0.40	0.40	9.00	1.52	1.52	0.40	1.43E-09
5	Happiness / Depression	0.40	0.40	0.40	0.40	3.05	0.40	0.40	8.92E-04
6	Moodiness	1.00	1.30	0.40	0.40	0.40	5.22	0.40	5.65E-06
7	Employment	1.70	0.40	1.70	1.10	0.40	0.40	2.52	2.96E-03
Min P-value		9.20E-07	1.65E-05	7.32E-04	1.43E-09	8.92E-04	5.65E-06	2.96E-03	

Figure 21: For the genomic locations in each cluster, we perform ontology term enrichment. For each cluster, we find the \log_{10} p-value of association between the GWAS cluster and the top ontology term of every other cluster. In general, we can tell that each cluster enriches very different ontology terms, as the off-diagonal scores are all very close to zero.

Flattening the Tensor along both the TF and GWAS Axis

It is very easy to take a tensor and reduce its order by simply flattening it. In our use case, we're taking a 3rd order tensor and turning it into a 2-D matrix by either taking each row to be a transcription factor, or each row to be a GWAS disease. We then have a column for each pair of genomic location, GWAS condition or genomic location, transcription factor respectively.

Now, we can view each of these rows as an incredibly high-dimensional vector. In order to analyze this extremely high-dimensional space, we can again use our favorite dimensionality-

reduction techniques: PCA followed by t-SNE. Plotting the first couple dimensions of PCA will allow us to see what factors are most important for explaining the data, while the t-SNE plot will allow us to see what hidden clusters exist within the data. Ideally this should correspond with our human-understanding of disease and transcription factor clusters.

One way we can interpret this tensor is to flatten it across one of the orders to get a 2-D matrix. For example, we could flatten it along the TF axis to have one row for each transcription factor and one column for each GWAS, genomic location pair. This is an incredibly high-dimensional space, so the first step is to perform principal component analysis in order to reduce the dimension to something more manageable. Plotting the first two principal components for the transcription factor flattening of the LD-corrected tensor yields an interesting structure:

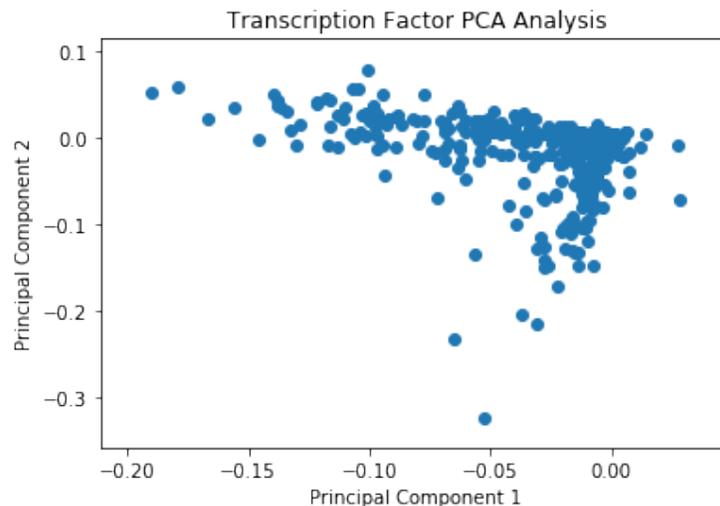


Figure 22: If we take the tensor of TF-GWAS effects and flatten it along the TF axis, we get a massive vector for each transcription factor. Performing principal component analysis on this tensor reveals a pretty interesting structure, In the upper left, ChIP-Seq experiments tend to congregate, while the bottom area of the figure contains a lot of homeobox SELEX experiments.

When we dive into the individual points on this plot, we notice that the transcription factors in the upper left corner are all from SELEX experiments, while the transcription factors

in the bottom right corner are from ChIP-seq experiments. If we study the deeper principal components, we notice that transcription factor 3 and 4 somewhat separates the transcription factors into known classes of transcription factors. For example, filtering all data points such that the 4th principal component is less than -0.08 yields the following list of transcription factors:

```
'D00356.005_SELEX_ELF1', 'D00371.003_SELEX_ERG',
'D00382.003_SELEX_ETV4', 'D00383.002_SELEX_ETV5',
'D00393.003_SELEX_FOXC2', 'D00394.003_SELEX_FOXD2',
'D00446.009_ChIP-seq_HNF4A', 'D00529.004_ChIP-seq_MYBL2',
'D00571.002_SELEX_PAX2', 'D00574.002_SELEX_PAX5',
'D00755.005_ChIP-seq_EP300', 'D00761.001_ChIP-seq_FOXA1',
'D00762.001_ChIP-seq_FOXA2', 'D00770.002_ChIP-seq_HDAC2',
'D00773.001_ChIP-seq_HNF4G', 'D00787.002_ChIP-seq_NFIC',
'D00818.003_ChIP-seq_TCF12'
```

If we take this PCA applied to transcription factors and then reduce the dimension even further using tSNE, we get some interesting structures:

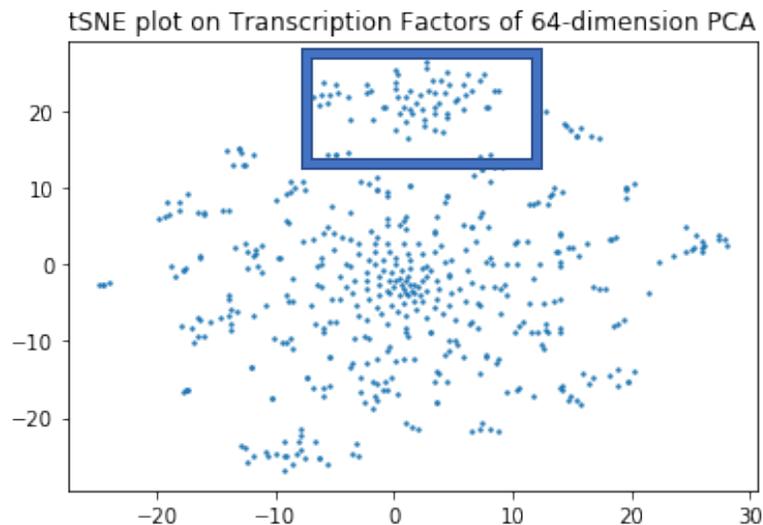


Figure 23: If we perform t-SNE visualization on our 64-dimension PCA, we find that the transcription factors form pretty clear clusters. In particular, the top cluster stands out as a region where a lot of the homeobox transcription factors cluster together.

If we look into the top large cluster and pull out the list of individual transcription factors, we see that many of the TFs in this list are related to one another. Many of these transcription factors fall into the LHX, HOX, and HMX class of transcription factors. This is the full list:

```

'D00290.003_SELEX_ALX3', 'D00292.001_SELEX_ALX4',
'D00295.002_SELEX_ARX', 'D00302.003_SELEX_BARHL2',
'D00303.002_SELEX_BARX1', 'D00313.002_SELEX_BSX',
'D00314.001_SELEX_CART1', 'D00335.003_SELEX_DLX2',
'D00336.001_SELEX_DLX3', 'D00337.001_SELEX_DLX4',
'D00338.002_SELEX_DLX5', 'D00339.002_SELEX_DLX6',
'D00342.003_SELEX_DRGX', 'D00364.002_SELEX_EMX1',
'D00365.003_SELEX_EMX2', 'D00366.003_SELEX_EN1',
'D00368.007_SELEX_EN2',
'D00377.001_SELEX_ESX1', 'D00385.002_SELEX_EVX1',
'D00386.001_SELEX_EVX2', 'D00413.003_SELEX_GBX1',
'D00416.005_SELEX_GBX2', 'D00427.003_SELEX_GSX1',
'D00428.002_SELEX_GSX2', 'D00431.002_SELEX_HESX1',
'D00440.003_SELEX_HMX1', 'D00441.003_SELEX_HMX2',
'D00442.003_SELEX_HMX3', 'D00451.002_SELEX_HOXA1',
'D00453.002_SELEX_HOXA2', 'D00455.001_SELEX_HOXB2',
'D00456.001_SELEX_HOXB3', 'D00457.003_SELEX_HOXB5',
'D00484.003_SELEX_ISL2', 'D00485.004_SELEX_ISX',
'D00491.003_SELEX_LBX2', 'D00493.003_SELEX_LHX2',
'D00495.003_SELEX_LHX6', 'D00497.002_SELEX_LHX9',
'D00498.002_SELEX_LMX1A', 'D00499.001_SELEX_LMX1B',
'D00513.002_SELEX_MEOX1', 'D00515.002_SELEX_MEOX2',
'D00518.002_SELEX_MIXL1', 'D00524.002_SELEX_MSX1',
'D00525.002_SELEX_MSX2', 'D00543.005_SELEX_NKX2-3',
'D00544.005_SELEX_NKX2-8', 'D00546.003_SELEX_NKX3-1',
'D00547.003_SELEX_NKX3-2', 'D00549.002_SELEX_NKX6-1',
'D00550.002_SELEX_NKX6-2', 'D00551.002_SELEX_NOTO',
'D00572.002_SELEX_PAX3', 'D00573.002_SELEX_PAX4',
'D00578.002_SELEX_PDX1', 'D00579.002_SELEX_PHOX2A',
'D00580.005_SELEX_PHOX2B', 'D00585.003_SELEX_PKNOX2',
'D00599.006_SELEX_POU6F2', 'D00602.008_SELEX_PROP1',
'D00604.004_SELEX_PRRX1', 'D00606.002_SELEX_PRRX2',
'D00612.001_SELEX_RAXL1', 'D00613.002_SELEX_RAX',
'D00632.001_SELEX_SHOX2', 'D00634.001_SELEX_SHOX',
'D00689.003_SELEX_TGIF1', 'D00698.001_SELEX_UNCX',
'D00701.001_SELEX_VAX1', 'D00702.001_SELEX_VAX2',
'D00705.003_SELEX_VENTX', 'D00707.004_SELEX_VSX1',
'D00708.001_SELEX_VSX2', 'D00713.003_SELEX_ZBTB49'

```

We can similarly repeat all of this analysis along the GWAS axis to see if the data clusters into recognizable structures. First, we notice that if we run a PCA analysis and plot the first two components, a small cluster emerges in the upper right corner of the plot (Figure 23). Digging into this cluster reveals that four traits in this cluster are BMI, basal metabolic rate, weight, and body fat percentage. All of this makes sense that these traits would cluster together. Another obvious group of points we notice is that below -0.2 in the PC2 axis. Digging in we notice a lot of very similar traits: miserableness, sensitivity / hurt feelings, irritability, mood swings, neuroticism score, loneliness / isolation, and guilty feelings. It's pretty amazing that such a clear clustering emerged from just a simple PCA analysis.

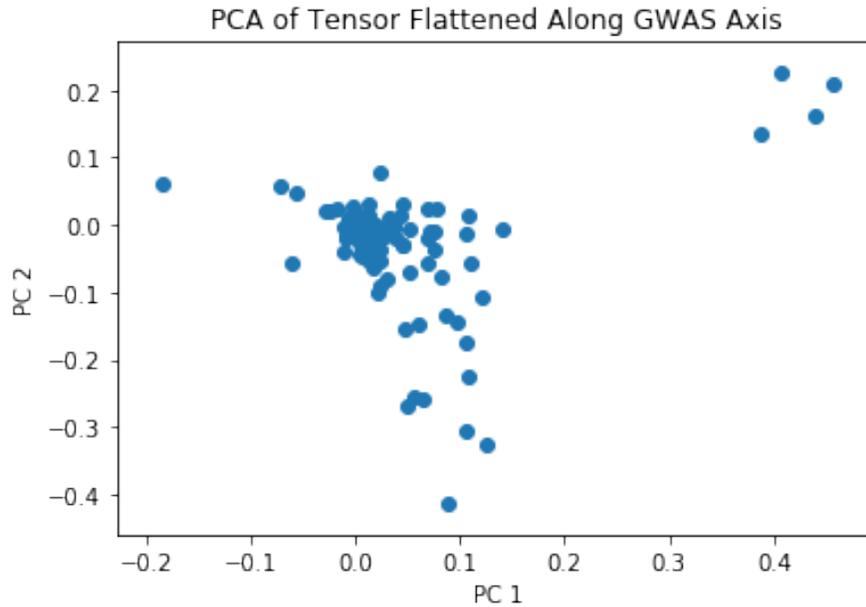


Figure 24: When we perform PCA on the GWAS traits, we see that traits related to size cluster substantially far away from the rest of the data. Many of the diseases which separate along the vertical axis (PC 2) are related to happiness.

After seeing these very promising results, we were excited to see what clusters a tSNE plot would reveal. Normally the tSNE algorithm runs into the issue where it makes fake clusters emerge where there is simply noise, but surprisingly the tSNE results were the exact opposite. Every point in the reduced-dimensional space was essentially uniformly spread out from every other point. This we thought was very bizarre, especially considering how well the PCA seemed to have been working at separating the GWAS traits into clusters which match our human intuition of the diseases.

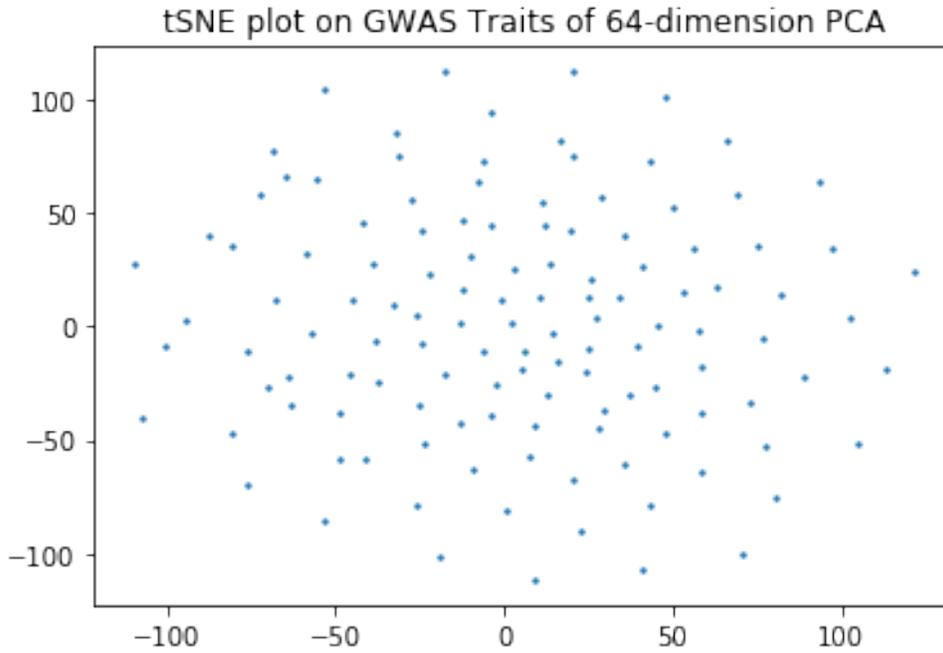


Figure 25: Interestingly, even though the PCA results looked good, a t-SNE dimensionality reduction did not reveal any noticeable structure for our GWAS traits.

A slight amount of structure still exists within the t-SNE plot. The central data points seem to be loosely related to weight/height/BMI, while a sampling of some of the points on the left edge of the graph reveal some related to depression.

Conclusion:

In today's era of big data and commoditized machine learning, we now have the capabilities to answer research questions that before were considered computationally impossible. Thanks to the Kipoi suite, it is now easier than ever to conduct experiments leveraging state of the art machine-learning models in new contexts. The work in this thesis started thanks to the knowledge that it is now incredibly easy to impute intermediate phenotypes across the entire genome. While much work has been done using intermediate phenotypes like histone modification and DNA accessibility, not nearly as much work has been done exploring transcription factor binding affinity as an intermediate phenotype.

This idea led to the construction of our first dataset: the variant effect dataset. For each of the transcription factors with pre-trained models in the Kipoi suite, let's compute the variant effect score of the transcription factor across a massive set of known genomic variants. After four days of intense computation on the Broad Cluster, we were able to analyze this dataset.

We ran two different post-analyses on the variant effect dataset. First, we first wanted to determine whether or not there were any latent structures behind the dataset. To do this, we used a popular technique known as an autoencoder to reduce the dimensionality substantially. Afterwards, we performed t-SNE to view and understand the data, as well as observe what clusters exist within the data. First, we ran this analysis pipeline on DNA locations, but unfortunately our final plots looked like noise (which in retrospect is expected). When we ran this same pipeline on the transcription factors, we found that many of the known classes of transcription factors clustered together in the t-SNE representation of the data. Additionally, the two different experiment types separated out from one another.

The second analysis we ran was correlating all of the variant effect scores for each transcription factor against all of the z-scores for each GWAS trait. This allows us to score how related each transcription factor is to each disease. Clustering both the diseases and transcription factors using UPGMA allows us to validate whether or not these correlations are at all meaningful. Overall, they were loosely meaningful, as the clusters did contain semantically similar diseases, but the hierarchical cluster was by no means ideal.

The next dataset we built was the GWAS, Transcription Factor, Genomic Location Tensor. For each tuple of GWAS Trait, Transcription Factor, and 100 kbp region in chromosome 1, we calculated an association score, which was essentially an LD-corrected dot product of the variant effect scores for the genomic region with the GWAS z-scores. This dataset was a computational beast to construct, taking over 3-weeks of continuous computation on 150 virtual-machines.

After this massive computational job was completed, we performed two analyses on the tensor. The first was a Latent Dirichlet Allocation analysis. There were two different ways we could perform this study, the first was viewing each genomic location as a document, while the second was viewing each disease as a word. There were quite a lot of hyperparameters to adjust in order to make the LDA analysis meaningful. Although we spent a lot of time tuning these different hyperparameters for the first strategy, and even performed a stop-word removal analysis, the overall results weren't very meaningful because we were unable to determine the actual human-understandable meaning of the different automatically discovered topics. However, the second approach was actually quite instructive. Although most of the genome looks like noise, when the topic distributions at a particular location deviate substantially from the background, we have shown that real biological phenomenon can occur at these locations,

and relate to the topics learned by the overall model. This could allow us to discover new regions of interest in the human genome.

The second analysis was a flatten analysis. We flattened the tensor in two different ways: the first was so that each row was a transcription factor, and the second was to make each row a GWAS trait. Performing the standard dimensionality reduction technique of Principal Component Analysis showed us that the dataset we constructed was actually meaningful. Simple classifiers based on these principal components showed us that the transcription factors would separate based first on the experiment type, then based on whether they were part of the large class of transcription factors known as homeobox TFs, while the GWAS traits separated based on biologically and semantically similar traits.

We hope that the lab will explore both of these datasets in more depth in the future, as there are a number of follow-up analyses that can be performed with this data. Essentially any study that in the past has been run with a different intermediate phenotype like histone modification can be run again with the intermediate phenotype of transcription factor variant effect sizes. There is a lot more that can be studied relating this intermediate phenotype with GWAS traits, and we hope that the tensor will hopefully reveal even more interesting insights.

Taking a step back, this year has been an incredible experience for me because I was able to immerse myself in a field that I've never touched before. I have a lot of experience with statistics, big data, and machine learning, but I've never gotten the opportunity to apply these areas of interest to such an interesting area of research. First semester was incredible in that I actually got to teach computational biology in a classroom setting, and nothing motivated me more than knowing that I was responsible for my students' educations. And nothing forced me to learn a piece of material more in-depthly than the fear of being asked tricky questions by a

classroom full of hungry graduate students. The overall thesis process was very difficult, but also intensely rewarding. At the beginning of the year, it felt like I was staring at a black hole, unable to know where to start or how to proceed, but as the pieces of the project slowly fell into place, it became easier and easier to work on the final result.

I'd like to thank Manolis Kellis and Yongjin Park for all of their help this semester.

Without them, none of this would've been possible.

Table of Figures

Figure 1: Plate Notation for Latent Dirichlet Allocation. The outer plates represent documents, while the inner plates represent the repeated choice of topics and words.	15
Figure 2: Pictorial Representation of a Sparse Autoencoder. An autoencoder will reduce the dimension of a dataset n to the number of nodes in the bottleneck layer.	18
Figure 3: t-SNE helps distinguish and separate data clusters that are not clearly separable with other visualization schemes (in this example hand-written digits). (top) t-SNE representation of dataset. (bottom) Sammon Mapping of the same data.....	20
Figure 4: Distribution of Correlation values between GWAS z-scores and transcription factor effect sizes across chromosome 1 shows lack of outlier correlations, indicating that the relationship between TF binding and GWAS is not immediately visible.	38
Figure 5: The pairwise correlation coefficients for each pair of disease and transcription factor does not immediately show any structure, indicating that the relationship between the two is not immediately visible. Hierarchical clustering was performed on both the transcription factors and the GWAS traits using UPGMA to get the ordering along both axes.....	39
Figure 6: Viewing each SNP as a vector of transcription factor variant effects, we perform autoencoding dimensionality-reduction followed by t-SNE visualization. Across all different settings of the hyper-parameters, we do not observe any observable structure in the plot based on genomic location.....	41
Figure 7: We now consider each transcription factor to be a vector of variant effect-sizes across all different SNPs, and perform an autoencoding dimensionality reduction followed by t-	

SNE visualization. Each labeled point is a transcription factor, with the first letter representing whether or not the model was built using a SELEX (S) experiment or a CHIP-seq (C) experiment. In the next couple figures we zoom into three noticeable large clusters.

..... 42

Figure 9: Zoomed in view of the upper left corner of figure 8 (t-SNE plot of transcription factors). This cluster is well structured as we see that many of the transcription factors have similar names, and also that most TFs are homeobox transcription factors, a class of transcription factor related to embryonic growth..... 43

Figure 9: Zoomed in view of the bottom right corner of figure 8. This cluster is also well-structured: we see that all of the transcription factor models built using CHIP-seq data cluster together. The upper left corner of this picture shows a dense cluster of transcription repressors. 44

Figure 10 Zoomed in view of the middle of the t-SNE representation of the transcription factors. We notice certain small patches of similar transcription factors (FOX and SOX), but it's difficult to tell if there's an overall unifying theme behind this cluster..... 45

Figure 11: Clustering of each genomic location (x-axis) according to the GWAS variant perturbation of TF motif scores (words) and assignment of each location into a set of topics (colors) defined by their word frequencies. Topic breakdown using LDA for threshold=20 across all genomic locations in chromosome 1. The distribution on each genomic location is mostly uniform, showing that we need to increase our word threshold. 48

Figure 12: (Top) Re-clustering of chromosome 1 genomic locations by GWAS-TF effects (as in Figure 11) now using a word threshold value of 50 on the LD-corrected dataset. (bottom) zoomed in view of the last 500kbp region. 49

Figure 13: Distribution of number of genomic location (document) appearances of each TF-trait tuple (word). Because this distribution has a fat tail, a small number of words appear in far more documents than others. These are stop-words. 51

Figure 14: Re-clustering of chromosome 1 genomic locations by GWAS-TF effects (as in Figure 12) after stop-words are removed. We observe that the per-document distributions are even spikier now because each document has very few words. Coupled with the lack of GWAS cluster structure, this shows the results are mostly noise..... 54

Figure 15: Re-clustering of chromosome 1 genomic locations by GWAS-TF effects (as in Figure 15) now with 100 unique clusters. Two clusters are extremely dominant in the topic breakdowns for all of the documents, while all other topics appear at very negligible frequencies. One cluster contains GWAS traits related to a person’s physical size, while the other contains GWAS traits related to a person’s general happiness. 55

Figure 16: (top) Re-clustering of chromosome 1 genomic locations by GWAS-TF effects (as in Figure 11), under a new structure that documents are genomic location-TF effect tuples, and words are diseases. Each topic distribution for the different genomic locations is the average topic breakdown of all documents containing that genomic location..... 59

Figure 17: Looking through the plot in figure 16 (clustering of all genomic regions when GWAS traits are viewed as words), we find a large region of red indicating that an area of the genome strongly relates to topic 3. Topic 3 covers weight-related traits. 60

Figure 18: Screenshot of the same genomic region (large red block of cluster 3 activity from figure 18) inspected with the Broad Institute Integrative Genomics Viewer. A number of relevant genes are located in this region, including PIK3R3, a gene related to growth-signaling. This gene might be why genomic regions are tagged with cluster 3 (weight-related GWAS) activity in this region. 60

Figure 19: Distance to Transcription Start Site for genomic regions in topic 3 (blue) compared to all 100kbp regions in chromosome 1 (white). In general, the genomic regions in topic three are closer to transcription start sites, implying that our selected genomic regions are more likely to be biologically relevant than the background regions. 62

Figure 20: Small subset of all the significantly enriched ontology terms for genomic locations in cluster 3. These are the gene ontology terms enriched from the biological processes ontology dataset. We notice that a lot of the terms relate to catabolic processes, a good sign considering cluster 3 relates to BMI / weight GWAS traits. 63

Figure 21: For the genomic locations in each cluster, we perform ontology term enrichment. For each cluster, we find the log10 p-value of association between the GWAS cluster and the top ontology term of every other cluster. In general, we can tell that each cluster enriches very different ontology terms, as the off-diagonal scores are all very close to zero. 66

Figure 22: If we take the tensor of TF-GWAS effects and flatten it along the TF axis, we get a massive vector for each transcription factor. Performing principal component analysis on this tensor reveals a pretty interesting structure, In the upper left, ChIP-Seq experiments tend to congregate, while the bottom area of the figure contains a lot of homeobox SELEX experiments. 67

Figure 23: If we perform t-SNE visualization on our 64-dimension PCA, we find that the transcription factors form pretty clear clusters. In particular, the top cluster stands out as a region where a lot of the homeobox transcription factors cluster together..... 68

Figure 24: When we perform PCA on the GWAS traits, we see that traits related to size cluster substantially far away from the rest of the data. Many of the diseases which separate along the vertical axis (PC 2) are related to happiness. 70

Figure 25: Interestingly, even though the PCA results looked good, a t-SNE dimensionality reduction did not reveal any noticeable structure for our GWAS traits. 71

Appendix

Gwas Names:

100007_raw	Polyunsaturated fat
100008_raw	Total sugars
1001	Duration of strenuous sports
100150	Drinking water intake
100240	Coffee consumed
100580	Alcohol consumed
1011	Frequency of light DIY in last 4 weeks
102_raw	Pulse rate- automated reading
102330	Sweets intake
1031	Frequency of friend/family visits
104400	Fruit consumers
104900	Time spent doing vigorous physical activity
1050	Time spend outdoors in summer
1090	Time spent driving
1100	Drive faster than motorway speed limit
1110	Length of mobile phone use
1160	Sleep duration
1180	Morning/evening person (chronotype)
1190	Nap during day
1200	Sleeplessness / insomnia
1239	Current tobacco smoking
1558	Alcohol intake frequency.
1707_1	Handedness (chirality/laterality): Right-handed
1727	Ease of skin tanning
1717	Skin colour
1757	Facial ageing
1920	Mood swings
1930	Miserableness
1940	Irritability
1950	Sensitivity / hurt feelings
20001_1001	Cancer code- self-reported: lung cancer
20001_1003	Cancer code- self-reported: skin cancer
20001_1006	Cancer code- self-reported: larynx/throat cancer
20001_1020	Cancer code- self-reported: large bowel cancer/colorectal cancer
20001_1002	Cancer code- self-reported: breast cancer

20001_1044 Cancer code- self-reported: prostate cancer
 20001_1048 Cancer code- self-reported: leukaemia
 20001_1059 Cancer code- self-reported: malignant melanoma
 20002_1065 Non-cancer illness code- self-reported: hypertension
 20002_1066 Non-cancer illness code- self-reported: heart/cardiac problem
 20002_1067 Non-cancer illness code- self-reported: peripheral vascular disease
 20002_1074 Non-cancer illness code- self-reported: angina
 20002_1075 Non-cancer illness code- self-reported: heart attack/myocardial infarction
 20002_1077 Non-cancer illness code- self-reported: heart arrhythmia
 20002_1079 Non-cancer illness code- self-reported: cardiomyopathy
 20002_1081 Non-cancer illness code- self-reported: stroke
 20002_1111 Non-cancer illness code- self-reported: asthma
 20002_1123 Non-cancer illness code- self-reported: sleep apnoea
 20002_1135 Non-cancer illness code- self-reported: stomach disorder
 20002_1136 Non-cancer illness code- self-reported: liver/biliary/pancreas problem
 20002_1154 Non-cancer illness code- self-reported: irritable bowel syndrome
 20002_1155 Non-cancer illness code- self-reported: hepatitis
 20002_1201 Non-cancer illness code- self-reported: bladder problem (not cancer)
 20002_1220 Non-cancer illness code- self-reported: diabetes
 20002_1224 Non-cancer illness code- self-reported: thyroid problem (not cancer)
 20002_1261 Non-cancer illness code- self-reported: multiple sclerosis
 20002_1262 Non-cancer illness code- self-reported: parkinsons disease
 20002_1264 Non-cancer illness code- self-reported: epilepsy
 20002_1265 Non-cancer illness code- self-reported: migraine
 20002_1286 Non-cancer illness code- self-reported: depression
 20002_1287 Non-cancer illness code- self-reported: anxiety/panic attacks
 20002_1288 Non-cancer illness code- self-reported: nervous breakdown
 20002_1289 Non-cancer illness code- self-reported: schizophrenia
 20002_1290 Non-cancer illness code- self-reported: deliberate self-harm/suicide attempt
 20002_1291 Non-cancer illness code- self-reported: mania/bipolar disorder/manic depression
 20002_1294 Non-cancer illness code- self-reported: back problem
 20002_1297 Non-cancer illness code- self-reported: muscle/soft tissue problem
 20002_1374 Non-cancer illness code- self-reported: allergy/hypersensitivity/anaphylaxis
 20002_1377 Non-cancer illness code- self-reported: polymyalgia rheumatica
 20002_1387 Non-cancer illness code- self-reported: hayfever/allergic rhinitis
 20002_1396 Non-cancer illness code- self-reported: enlarged prostate
 20002_1408 Non-cancer illness code- self-reported: alcohol dependency
 20002_1425 Non-cancer illness code- self-reported: cerebral aneurysm
 20002_1440 Non-cancer illness code- self-reported: tuberculosis (tb)
 20002_1446 Non-cancer illness code- self-reported: anaemia

20002_1453 Non-cancer illness code- self-reported: psoriasis
 20002_1462 Non-cancer illness code- self-reported: crohns disease
 20002_1464 Non-cancer illness code- self-reported: rheumatoid arthritis
 20002_1473 Non-cancer illness code- self-reported: high cholesterol
 20002_1500 Non-cancer illness code- self-reported: vertigo
 20002_1534 Non-cancer illness code- self-reported: back pain
 20002_1614 Non-cancer illness code- self-reported: stress
 20016_irnt Fluid intelligence score
 20018 Prospective memory result
 20023_irnt Mean time to correctly identify matches
 20122 Bipolar disorder status
 20127_irnt Neuroticism score
 2020 Loneliness- isolation
 2030 Guilty feelings
 2040 Risk taking
 20401 Ever addicted to any substance or behaviour
 20404 Ever physically dependent on alcohol
 20403 Amount of alcohol drunk on a typical drinking day
 20408 Frequency of memory loss due to drinking alcohol in last year
 20416 Frequency of consuming six or more units of alcohol
 20420_irnt Longest period spent worried or anxious
 20453 Ever taken cannabis
 20454 Maximum frequency of taking cannabis
 20456 Ever addicted to illicit or recreational drugs
 20458 General happiness
 20463 Ever heard an un-real voice
 20468 Ever believed in an un-real conspiracy against self
 20480 Ever self-harmed
 20483 Ever attempted suicide
 20503 Ever addicted to prescription or over-the-counter medication
 20511 Recent poor appetite or overeating
 20517 Trouble falling or staying asleep- or sleeping too much
 20519 Recent feelings of tiredness or low energy
 20544_17 Mental health problems ever diagnosed by a professional: Agoraphobia
 20544_12 Mental health problems ever diagnosed by a professional: Bulimia nervosa
 Mental health problems ever diagnosed by a professional: Obsessive compulsive disorder
 20544_7 (OCD)
 20549_4 Substances taken for anxiety: Drugs or alcohol (more than once)
 21001_irnt Body mass index (BMI)
 21002_raw Weight

2139_irnt Age first had sexual intercourse
 2149 Lifetime number of sexual partners
 2159 Ever had same-sex intercourse
 2207 Wears glasses or contact lenses
 23099_irnt Body fat percentage
 23105_irnt Basal metabolic rate
 30020_raw Haemoglobin concentration
 30080_raw Platelet count
 3751 Chest pain or discomfort when walking uphill or hurrying
 3773 Knee pain for 3+ months
 4570 Friendships satisfaction
 4548 Health satisfaction
 5983_raw ECG- heart rate
 6142_5 Current employment status: Unemployed
 6142_1 Current employment status: In paid employment or self-employed
 6160_3 Leisure/social activities: Religious group
 757_irnt Time employed in main current job
 93_raw Systolic blood pressure- manual reading

Transcription Factors

D00290.003_SELEX_ALX3	D00316.003_SELEX_CDX2	D00340.003_SELEX_DMBX1
D00292.001_SELEX_ALX4	D00317.003_SELEX_CEBPB	D00341.003_SELEX_DPRX
D00293.003_SELEX_ARNTL	D00317.009_ChIP-seq_CEBPB	D00342.003_SELEX_DRGX
D00295.002_SELEX_ARX	D00318.001_SELEX_CEBPD	D00343.003_SELEX_DUXA
D00296.006_SELEX_AR	D00318.004_ChIP-seq_CEBPD	D00344.002_SELEX_E2F1
D00298.003_SELEX_ATF4	D00319.003_SELEX_CEBPE	D00344.005_ChIP-seq_E2F1
D00299.003_SELEX_ATF7	D00320.001_SELEX_CEBPG	D00345.002_SELEX_E2F2
D00302.003_SELEX_BARHL2	D00321.003_SELEX_CENPB	D00346.003_SELEX_E2F3
D00303.002_SELEX_BARX1	D00322.002_SELEX_CLOCK	D00347.003_SELEX_E2F4
D00304.003_SELEX_BATF3	D00323.002_SELEX_CPEB1	D00347.005_ChIP-seq_E2F4
D00305.003_SELEX_BCL6B	D00324.003_SELEX_CREB3L	D00348.003_SELEX_E2F7
D00306.003_SELEX_BHLHA1	1	D00349.002_SELEX_E2F8
5	D00326.002_SELEX_CREB3	D00350.002_SELEX_EBF1
D00308.003_SELEX_BHLHB2	D00328.003_SELEX_CTCF	D00350.005_ChIP-seq_EBF1
D00309.002_SELEX_BHLHB3	D00328.018_ChIP-seq_CTCF	D00351.001_SELEX_EGR1
D00310.002_SELEX_BHLHE2	D00329.002_SELEX_CUX1	D00351.006_SELEX_EGR1
2	D00330.002_SELEX_CUX2	D00351.009_ChIP-seq_EGR1
D00311.003_SELEX_BHLHE2	D00331.006_SELEX_DBP	D00353.005_SELEX_EGR2
3	D00332.003_SELEX_DLX1	D00354.002_SELEX_EGR3
D00312.002_SELEX_BHLHE4	D00335.003_SELEX_DLX2	D00355.002_SELEX_EGR4
1	D00336.001_SELEX_DLX3	D00356.005_SELEX_ELF1
D00313.002_SELEX_BSX	D00337.001_SELEX_DLX4	D00356.010_ChIP-seq_ELF1
D00314.001_SELEX_CART1	D00338.002_SELEX_DLX5	D00357.002_SELEX_ELF3
D00315.003_SELEX_CDX1	D00339.002_SELEX_DLX6	D00359.004_SELEX_ELF5

D00360.004_SELEX_ELK1	D00417.005_SELEX_GCM1	D00480.002_SELEX_IRF9
D00360.007_ChIP-seq_ELK1	D00418.003_SELEX_GCM2	D00481.003_SELEX_IRX2
D00361.001_SELEX_ELK3	D00419.003_SELEX_GLI2	D00483.001_SELEX_IRX5
D00363.003_SELEX_ELK4	D00420.004_SELEX_GLI1	D00484.003_SELEX_ISL2
D00363.004_ChIP-seq_ELK4	D00421.003_SELEX_GLI2	D00485.004_SELEX_ISX
D00364.002_SELEX_EMX1	D00422.003_SELEX_GLI3	D00487.003_SELEX_JDP2
D00365.003_SELEX_EMX2	D00423.005_SELEX_GMEB2	D00488.003_SELEX_KLF13
D00366.003_SELEX_EN1	D00424.002_SELEX_GRHL1	D00489.002_SELEX_KLF14
D00368.007_SELEX_EN2	D00425.002_SELEX_GSC2	D00490.003_SELEX_KLF16
D00369.003_SELEX_EOMES	D00426.003_SELEX_GSC	D00491.003_SELEX_LBX2
D00370.003_SELEX_ERF	D00427.003_SELEX_GSX1	D00492.003_SELEX_LEF1
D00371.003_SELEX_ERG	D00428.002_SELEX_GSX2	D00493.003_SELEX_LHX2
D00372.003_SELEX_ESR1	D00429.003_SELEX_HES5	D00495.003_SELEX_LHX6
D00374.003_SELEX_ESRRA	D00430.003_SELEX_HES7	D00497.002_SELEX_LHX9
D00375.003_SELEX_ESRRB	D00431.002_SELEX_HESX1	D00498.002_SELEX_LMX1A
D00376.003_SELEX_ESRRG	D00432.003_SELEX_HEY1	D00499.001_SELEX_LMX1B
D00377.001_SELEX_ESX1	D00433.003_SELEX_HEY2	D00501.003_SELEX_MAFF
D00378.007_SELEX_ETS1	D00435.003_SELEX_HIC2	D00501.004_ChIP-seq_MAFF
D00378.009_ChIP-seq_ETS1	D00436.002_SELEX_HINFP1	D00502.002_SELEX_MAFG
D00379.001_SELEX_ETV1	D00438.003_SELEX_HLF	D00503.004_SELEX_MAFK
D00380.002_SELEX_ETV2	D00440.003_SELEX_HMX1	D00503.014_ChIP-seq_MAFK
D00381.003_SELEX_ETV3	D00441.003_SELEX_HMX2	D00504.002_SELEX_MAX
D00382.003_SELEX_ETV4	D00442.003_SELEX_HMX3	D00504.005_ChIP-seq_MAX
D00383.002_SELEX_ETV5	D00443.001_SELEX_HNF1A	D00505.003_SELEX_MEF2A
D00384.002_SELEX_ETV6	D00444.001_SELEX_HNF1B	D00505.006_ChIP-
D00385.002_SELEX_EVX1	D00446.005_SELEX_HNF4A	seq_MEF2A
D00386.001_SELEX_EVX2	D00446.009_ChIP-seq_HNF4A	D00506.003_SELEX_MEF2B
D00387.001_SELEX_FEV	D00447.002_SELEX_HOMEZ	D00507.002_SELEX_MEF2D
D00388.003_SELEX_FIGLA	D00448.002_SELEX_HOXA10	D00508.004_SELEX_MEIS1
D00389.005_SELEX_FLI1	D00450.003_SELEX_HOXA13	D00510.003_SELEX_MEIS2
D00390.002_SELEX_FOXB1	D00451.002_SELEX_HOXA1	D00512.001_SELEX_MEIS3
D00392.002_SELEX_FOXC1	D00453.002_SELEX_HOXA2	D00513.002_SELEX_MEOX1
D00393.003_SELEX_FOXC2	D00454.003_SELEX_HOXB13	D00515.002_SELEX_MEOX2
D00394.003_SELEX_FOXD2	D00455.001_SELEX_HOXB2	D00516.002_SELEX_MESP1
D00395.002_SELEX_FOXD3	D00456.001_SELEX_HOXB3	D00517.002_SELEX_MGA
D00396.003_SELEX_FOXC1	D00457.003_SELEX_HOXB5	D00518.002_SELEX_MIXL1
D00397.003_SELEX_FOXI1	D00459.003_SELEX_HOXC10	D00519.002_SELEX_MLXIPL
D00398.001_SELEX_FOXJ2	D00460.003_SELEX_HOXC11	D00520.002_SELEX_MLX
D00400.003_SELEX_FOXJ3	D00461.003_SELEX_HOXC12	D00521.003_SELEX_MNT
D00401.003_SELEX_FOXK1	D00462.002_SELEX_HOXC13	D00522.002_SELEX_MNX1
D00403.001_SELEX_FOXL1	D00463.003_SELEX_HOXD11	D00523.003_SELEX_MSC
D00404.002_SELEX_FOXO1	D00464.001_SELEX_HOXD12	D00524.002_SELEX_MSX1
D00405.003_SELEX_FOXO3	D00465.003_SELEX_HOXD13	D00525.002_SELEX_MSX2
D00406.002_SELEX_FOXO4	D00468.003_SELEX_HOXD8	D00527.003_SELEX_MTF1
D00407.003_SELEX_FOXO6	D00470.005_SELEX_HSF1	D00528.001_SELEX_MYBL1
D00408.003_SELEX_FOXP3	D00471.002_SELEX_HSF2	D00529.001_SELEX_MYBL2
D00409.003_SELEX_GABPA	D00472.003_SELEX_HSF4	D00529.004_ChIP-
D00409.004_ChIP-	D00473.003_SELEX_HSFY2	seq_MYBL2
seq_GABPA	D00474.003_SELEX_ID4	D00530.003_SELEX_MYF6
D00410.003_SELEX_GATA3	D00475.003_SELEX_IRF3	D00531.003_SELEX_NEURO
D00410.009_ChIP-	D00475.006_ChIP-seq_IRF3	D2
seq_GATA3	D00476.002_SELEX_IRF4	D00532.003_SELEX_NEURO
D00411.003_SELEX_GATA4	D00476.004_ChIP-seq_IRF4	G2
D00412.002_SELEX_GATA5	D00477.003_SELEX_IRF5	D00533.003_SELEX_NFAT5
D00413.003_SELEX_GBX1	D00478.003_SELEX_IRF7	D00534.003_SELEX_NFATC1
D00416.005_SELEX_GBX2	D00479.006_SELEX_IRF8	

D00534.004_ChIP-seq_NFATC1
D00535.003_SELEX_NFE2
D00535.004_ChIP-seq_NFE2
D00536.003_SELEX_NFIA
D00537.001_SELEX_NFIB
D00538.001_SELEX_NFIL3
D00539.001_SELEX_NFIX
D00540.002_SELEX_NFKB1
D00541.001_SELEX_NFKB2
D00542.005_SELEX_NHLH1
D00543.005_SELEX_NKX2-3
D00544.005_SELEX_NKX2-8
D00546.003_SELEX_NKX3-1
D00547.003_SELEX_NKX3-2
D00549.002_SELEX_NKX6-1
D00550.002_SELEX_NKX6-2
D00551.002_SELEX_NOTO
D00552.002_SELEX_NR2C2
D00552.006_ChIP-seq_NR2C2
D00553.003_SELEX_NR2E1
D00554.003_SELEX_NR2F1
D00555.002_SELEX_NR2F6
D00556.003_SELEX_NR3C1
D00557.002_SELEX_NR3C2
D00558.002_SELEX_NR4A2
D00559.001_SELEX_NRF1
D00559.006_ChIP-seq_NRF1
D00560.003_SELEX_NRL
D00561.001_SELEX_OLIG1
D00562.006_SELEX_OLIG2
D00563.002_SELEX_OLIG3
D00564.004_SELEX_ONECUT1
D00565.002_SELEX_ONECUT2
D00566.003_SELEX_ONECUT3
D00568.001_SELEX_OTX1
D00569.002_SELEX_OTX2
D00570.002_SELEX_PAX1
D00571.002_SELEX_PAX2
D00572.002_SELEX_PAX3
D00573.002_SELEX_PAX4
D00574.002_SELEX_PAX5
D00574.006_ChIP-seq_PAX5
D00575.003_SELEX_PAX6
D00576.005_SELEX_PAX7
D00577.002_SELEX_PAX9
D00578.002_SELEX_PDX1
D00579.002_SELEX_PHOX2A
D00580.005_SELEX_PHOX2B
D00581.002_SELEX_PITX1
D00582.001_SELEX_PITX3
D00583.002_SELEX_PKNOX1
D00585.003_SELEX_PKNOX2
D00586.004_SELEX_POU1F1
D00587.002_SELEX_POU2F1
D00589.001_SELEX_POU2F2
D00589.005_ChIP-seq_POU2F2
D00590.001_SELEX_POU2F3
D00591.002_SELEX_POU3F1
D00592.003_SELEX_POU3F2
D00593.003_SELEX_POU3F3
D00594.002_SELEX_POU3F4
D00595.001_SELEX_POU4F1
D00596.005_SELEX_POU4F2
D00597.001_SELEX_POU4F3
D00598.001_SELEX_POU5F1P1
D00599.006_SELEX_POU6F2
D00600.001_SELEX_PRDM1
D00600.004_ChIP-seq_PRDM1
D00601.003_SELEX_PRDM4
D00602.008_SELEX_PROP1
D00603.003_SELEX_PROX1
D00604.004_SELEX_PRRX1
D00606.002_SELEX_PRRX2
D00607.006_SELEX_RARA
D00610.003_SELEX_RARB
D00611.001_SELEX_RARG
D00612.001_SELEX_RAXL1
D00613.002_SELEX_RAX
D00614.001_SELEX_RFX2
D00616.002_SELEX_RFX3
D00618.003_SELEX_RFX4
D00619.003_SELEX_RFX5
D00619.007_ChIP-seq_RFX5
D00621.003_SELEX_RHOXF1
D00622.003_SELEX_RORA
D00623.003_SELEX_RUNX2
D00624.007_SELEX_RUNX3
D00624.009_ChIP-seq_RUNX3
D00626.005_SELEX_RXRA
D00626.009_ChIP-seq_RXRA
D00628.002_SELEX_RXRB
D00629.006_SELEX_RXRG
D00630.003_SELEX_SCRT1
D00631.002_SELEX_SCRT2
D00632.001_SELEX_SHOX2
D00634.001_SELEX_SHOX
D00635.001_SELEX_SMAD3
D00636.003_SELEX_SNAI2
D00637.003_SELEX_SOX10
D00639.003_SELEX_SOX14
D00640.003_SELEX_SOX15
D00642.003_SELEX_SOX18
D00644.003_SELEX_SOX21
D00645.006_SELEX_SOX2
D00646.002_SELEX_SOX4
D00647.002_SELEX_SOX7
D00648.004_SELEX_SOX8
D00649.002_SELEX_SOX9
D00650.005_SELEX_SP1
D00650.007_ChIP-seq_SP1
D00651.003_SELEX_SP3
D00652.003_SELEX_SP4
D00652.004_ChIP-seq_SP4
D00653.003_SELEX_SP8
D00654.003_SELEX_SPDEF
D00655.002_SELEX_SPI1
D00655.006_ChIP-seq_SPI1
D00656.001_SELEX_SPIB
D00658.003_SELEX_SPIC
D00659.002_SELEX_SREBF2
D00660.005_SELEX_SRF
D00660.007_ChIP-seq_SRF
D00661.003_SELEX_SRY
D00662.003_SELEX_TBR1
D00663.002_SELEX_TBX15
D00664.002_SELEX_TBX19
D00665.002_SELEX_TBX1
D00666.002_SELEX_TBX20
D00667.001_SELEX_TBX21
D00668.003_SELEX_TBX2
D00669.002_SELEX_TBX4
D00670.002_SELEX_TBX5
D00672.001_SELEX_TCF3
D00673.001_SELEX_TCF4
D00674.003_SELEX_TCF7L1
D00677.003_SELEX_TEAD1
D00678.001_SELEX_TEAD3
D00679.004_ChIP-seq_TEAD4
D00680.001_SELEX_TEF
D00681.002_SELEX_TFAP2A
D00681.004_ChIP-seq_TFAP2A
D00682.002_SELEX_TFAP2B
D00683.002_SELEX_TFAP2C
D00683.007_ChIP-seq_TFAP2C
D00684.005_SELEX_TFAP4
D00685.003_SELEX_TFCP2
D00686.002_SELEX_TFE3
D00687.001_SELEX_TFEB
D00688.002_SELEX_TFEC
D00689.003_SELEX_TGIF1
D00690.003_SELEX_TGIF2LX
D00691.001_SELEX_TGIF2
D00692.002_SELEX_THRA
D00693.002_SELEX_THRB
D00695.002_SELEX_TP63
D00697.002_SELEX_T
D00698.001_SELEX_UNCX

D00700.001_SELEX_USF1
D00700.006_ChIP-seq_USF1
D00701.001_SELEX_VAX1
D00702.001_SELEX_VAX2
D00704.001_SELEX_VDR
D00705.003_SELEX_VENTX
D00707.004_SELEX_VSX1
D00708.001_SELEX_VSX2
D00709.002_SELEX_XBP1
D00710.002_SELEX_YY1
D00710.007_ChIP-seq_YY1
D00711.003_SELEX_YY2
D00712.002_SELEX_ZBED1
D00713.003_SELEX_ZBTB49
D00714.003_SELEX_ZBTB7A
D00714.004_ChIP-seq_ZBTB7A
D00715.002_SELEX_ZBTB7B
D00716.004_SELEX_ZBTB7C
D00719.002_SELEX_ZIC1
D00721.002_SELEX_ZIC3
D00722.002_SELEX_ZIC4
D00723.003_SELEX_ZNF143
D00723.006_ChIP-seq_ZNF143
D00724.003_SELEX_ZNF232
D00725.004_SELEX_ZNF238
D00726.003_SELEX_ZNF282
D00727.003_SELEX_ZNF306
D00728.002_SELEX_ZNF410
D00729.003_SELEX_ZNF435
D00730.002_SELEX_ZNF524
D00731.003_SELEX_ZNF713
D00732.005_SELEX_ZNF740
D00733.003_SELEX_ZNF75A
D00734.003_SELEX_ZNF784
D00735.002_SELEX_ZSCAN4
D00736.002_ChIP-seq_ARID3A
D00737.001_ChIP-seq_ATF1
D00738.002_ChIP-seq_ATF2
D00739.001_ChIP-seq_ATF3
D00740.002_ChIP-seq_BACH1
D00741.001_ChIP-seq_BATF
D00742.002_ChIP-seq_BCL11A
D00743.001_ChIP-seq_BCL3
D00744.001_ChIP-seq_BCLAF1
D00745.001_ChIP-seq_BDP1
D00746.004_ChIP-seq_BHLHE40
D00747.001_ChIP-seq_BRCA1
D00748.001_ChIP-seq_CBX3
D00749.001_ChIP-seq_CCNT2
D00750.001_ChIP-seq_CHD1
D00751.005_ChIP-seq_CHD2
D00752.001_ChIP-seq_CTBP2
D00753.001_ChIP-seq_CTCFL
D00754.003_ChIP-seq_E2F6
D00755.003_ChIP-seq_EP300
D00755.005_ChIP-seq_EP300
D00756.007_ChIP-seq_EZH2
D00757.001_ChIP-seq_FAM48A
D00758.001_ChIP-seq_FOSL1
D00759.001_ChIP-seq_FOSL2
D00760.003_ChIP-seq_FOS
D00761.001_ChIP-seq_FOXA1
D00762.001_ChIP-seq_FOXA2
D00763.001_ChIP-seq_FOXM1
D00764.001_ChIP-seq_FOXP2
D00765.001_ChIP-seq_GATA1
D00766.002_ChIP-seq_GATA2
D00767.001_ChIP-seq_GTF2B
D00768.001_ChIP-seq_GTF2F1
D00769.001_ChIP-seq_HDAC1
D00770.002_ChIP-seq_HDAC2
D00771.001_ChIP-seq_HDAC6
D00772.001_ChIP-seq_HMGN3
D00773.001_ChIP-seq_HNF4G
D00774.001_ChIP-seq_IKZF1
D00775.001_ChIP-seq_JUNB
D00776.005_ChIP-seq_JUND
D00777.002_ChIP-seq_JUN
D00778.002_ChIP-seq_KAP1
D00779.001_ChIP-seq_KDM5A
D00780.001_ChIP-seq_KDM5B
D00781.003_ChIP-seq_MAZ
D00782.001_ChIP-seq_MEF2C
D00783.001_ChIP-seq_MTA3
D00784.004_ChIP-seq_MXI1
D00785.001_ChIP-seq_MYC
D00786.001_ChIP-seq_NANOG
D00787.002_ChIP-seq_NFIC
D00788.001_ChIP-seq_NFYA
D00789.003_ChIP-seq_NFYB
D00790.001_ChIP-seq_NR2F2
D00791.001_ChIP-seq_PBX3
D00792.001_ChIP-seq_PHF8
D00793.002_ChIP-seq_PML
D00794.047_ChIP-seq_POLR2A
D00795.001_ChIP-seq_POU5F1
D00796.001_ChIP-seq_RAD21
D00797.001_ChIP-seq_RBBP5
D00798.001_ChIP-seq_RCOR1
D00799.001_ChIP-seq_REST
D00800.001_ChIP-seq_RPC155
D00801.001_ChIP-seq_SAP30
D00802.001_ChIP-seq_SETDB1
D00803.002_ChIP-seq_SIN3AK20
D00804.002_ChIP-seq_SIN3A
D00805.001_ChIP-seq_SIRT6
D00806.003_ChIP-seq_SIX5
D00807.001_ChIP-seq_SMARCB1
D00808.004_ChIP-seq_SMC3
D00809.002_ChIP-seq_SP2
D00810.001_ChIP-seq_STAT3
D00811.001_ChIP-seq_STAT5A
D00812.001_ChIP-seq_SUZ12
D00813.009_ChIP-seq_TAF1
D00814.001_ChIP-seq_TAF7
D00815.001_ChIP-seq_TAL1
D00816.001_ChIP-seq_TBL1XR1
D00817.001_ChIP-seq_TBP
D00818.003_ChIP-seq_TCF12
D00819.002_ChIP-seq_TCF7L2
D00820.001_ChIP-seq_THAP1
D00821.001_ChIP-seq_TRIM28
D00822.001_ChIP-seq_UBTF
D00823.002_ChIP-seq_USF2
D00824.001_ChIP-seq_WRNIP1
D00825.001_ChIP-seq_ZBTB33
D00826.001_ChIP-seq_ZEB1
D00827.001_ChIP-seq_ZNF217
D00828.002_ChIP-seq_ZNF263
D00829.002_ChIP-seq_ZNF274

References

- ⁱ Park, Yongjin, and Manolis Kellis. "Deep learning for regulatory genomics." *Nature biotechnology* 33.8 (2015): 825.
- ⁱⁱ Avsec, Ziga, et al. "Kipoi: accelerating the community exchange and reuse of predictive models for genomics." *bioRxiv* (2018): 375345.
- ⁱⁱⁱ Alipanahi, Babak, et al. "Predicting the sequence specificities of DNA-and RNA-binding proteins by deep learning." *Nature biotechnology* 33.8 (2015): 831.
- ^{iv} Zeng, Haoyang, and David K. Gifford. "Predicting the impact of non-coding variants on DNA methylation." *Nucleic acids research* 45.11 (2017): e99-e99.
- ^v Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." *Journal of machine Learning research* 3.Jan (2003): 993-1022.
- ^{vi} Rogers, Simon, et al. "The latent process decomposition of cDNA microarray data sets." *IEEE/ACM transactions on computational biology and bioinformatics* 2.2 (2005): 143-156.
- ^{vii} Bicego, Manuele, et al. "Expression microarray classification using topic models." *Proceedings of the 2010 ACM Symposium on Applied Computing*. ACM, 2010.
- ^{viii} Pinoli, Pietro, Davide Chicco, and Marco Masseroli. "Latent Dirichlet allocation based on Gibbs sampling for gene function prediction." *Computational Intelligence in Bioinformatics and Computational Biology, 2014 IEEE Conference on*. IEEE, 2014.
- ^{ix} Liu, Bing, et al. "Identifying functional miRNA–mRNA regulatory modules with correspondence latent dirichlet allocation." *Bioinformatics* 26.24 (2010): 3105-3111.

-
- ^x Blei, David M., and Michael I. Jordan. "Modeling annotated data." *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*. ACM, 2003.
- ^{xi} Pu, Yunchen, et al. "Variational autoencoder for deep learning of images, labels and captions." *Advances in neural information processing systems*. 2016.
- ^{xii} Semeniuta, Stanislau, Aliaksei Severyn, and Erhardt Barth. "A hybrid convolutional variational autoencoder for text generation." *arXiv preprint arXiv:1702.02390* (2017).
- ^{xiii} Kingma, Diederik P., and Max Welling. "Auto-encoding variational bayes." *arXiv preprint arXiv:1312.6114* (2013).
- ^{xiv} Maaten, Laurens van der, and Geoffrey Hinton. "Visualizing data using t-SNE." *Journal of machine learning research* 9.Nov (2008): 2579-2605.
- ^{xv} Ge, Tian, et al. "Phenome-wide heritability analysis of the UK Biobank." *PLoS genetics* 13.4 (2017): e1006711.
- ^{xvi} Bulik-Sullivan, Brendan K., et al. "LD 1Score regression distinguishes confounding from polygenicity in genome-wide association studies." *Nature genetics* 47.3 (2015): 291.
- ^{xvii} Xia, Hanzhang, et al. "Building a schizophrenia genetic network: Evidence that Transcription Factor 4 regulates Schizophrenia Risk Genes." *Human Molecular Genetics*
- ^{xviii} Georgieva, Lyudmila, et al. "Convergent evidence that oligodendrocyte lineage transcription factor 2 (OLIG2) and interacting genes influence susceptibility to schizophrenia." *Proceedings of the National Academy of Sciences* 103.33 (2006): 12469-12474.
- ^{xix} Sandelin, Albin, et al. "JASPAR: an open-access database for eukaryotic transcription factor binding profiles." *Nucleic acids research* 32.suppl_1 (2004): D91-D94

^{xx} Rao, Anjana, Chun Luo, and Patrick G. Hogan. "Transcription factors of the NFAT family: regulation and function." *Annual review of immunology* 15.1 (1997): 707-747.

^{xxi} Purcell, Shaun, et al. "PLINK: a tool set for whole-genome association and population-based linkage analyses." *The American journal of human genetics* 81.3 (2007): 559-575.

^{xxii} McLean, Cory Y., et al. "GREAT improves functional interpretation of cis-regulatory regions." *Nature biotechnology* 28.5 (2010): 495.