



Technische Universität München

Department of Bioinformatics and Computational Biology

Technische Universität München

Master's Thesis in Bioinformatics

Variation of HERV elements in the KORA cohort

Julian Schmidt

Department of Bioinformatics and Computational Biology

Technische Universität München

Master's Thesis in Bioinformatics

Variation of HERV elements in the KORA cohort

Variation von HERV elementen in der KORA Kohorte

Author: Julian Schmidt
Supervisor: Dr. Matthias Heinig
Advisor: Johann Hawe
Submitted: 15.04.2018

Contents

1	Introduction	1
1.1	Epigenetics	1
1.2	Human endogenous retroviruses	1
1.3	Effect network analysis	2
2	Data	3
2.1	HERV annotation	3
2.2	KORA	4
2.2.1	Expression	5
2.2.2	Methylation	6
2.2.3	Genotypes	6
2.2.4	Covariates	7
2.2.5	Methylation quantitative trait loci	8
2.2.6	Data:meQTL	8
2.3	Transcription factor binding	8
2.4	Chromatin states	9
2.5	Data:Chromatin states	9
3	Methods	11
3.1	Overlaps	11
3.2	Data normalization	11
3.3	eQTL/eQTM calculation	12
3.4	Functional Analysis of Gene Sets	12
3.5	Gaussian Graphical Models	12
4	Results	13
4.1	Normalized Data	13
4.2	HERV region features	13
4.2.1	Expression	14
4.2.2	Methylation	15
4.2.3	Genotypes	16
4.2.4	Chromatin states	16
4.3	eQTLs	16
4.4	eQTMs	18
4.5	meQTLs	19
4.6	HERV related regulatory networks	19
4.6.1	Data collection	19
5	Discussion	21

5.1	HERV region features	21
5.1.1	Expression	21
5.1.2	Methylation	21
5.1.3	Genotypes	21
6	Conclusion and Outlook	23

Chapter 1

Introduction

1.1 Epigenetics

- central dogma of biology
- rising importance of other factors atop DNA → epigenetics
- quick overview epigenetic marks
- chromatin states
- DNA methylation
 - most understood and researched epigenetic mark [1]
 - mechanism: methylation fifth position of cytosine, almost exclusively in CpG context in mammals, installed and maintained by 3 enzymes: DNA methyltransferase 1/3A/3B (DNMT1/...)
 - generally considered as silencing mark
 - CpG islands = groups of cg pairs, often in promoter regions, usually hypomethylated
 - promoter regions kept free of methylation
- snp – > cpg – > expression pattern
- TF < – > cpg interaction

1.2 Human endogenous retroviruses

When the draft for the first human genome was published in 2001[2] there was the expectation, that it would allow to fully understand the human genome. But soon after genes were annotated less than 5% of the genome was found to be protein coding[3] and later on the amount of exonic protein coding DNA was estimated to be as low as 1.2%[4]. The remaining majority of genome was often labeled as non-functional "junk DNA" [5].

Since then this notion has generally been revoked. Especially within the two phases of the Encyclopedia of DNA Elements (ENCODE) project[4], participation in some functional role has been assigned to up 80.4% of the human genome.

Still up to this day there are areas in the genome that are hard to analyze. Especially repetitive sequences cause huge problems for any method that relies on sequencing[6]. In human repetitive DNA makes up about 50% of the whole Genome[6]. There are five major classes of repetitive elements: Satellites, short interspersed nuclear elements (SINE), long interspersed nuclear elements (LINE), DNA transposons and long terminal repeat (LTR) retrotransposons.

In this work we focused on a subclass of the latter repeat class. Human endogenous retroviruses (HERVs) belong to the class of LTR retrotransposons and make up about 8% of the human genome[7]. HERVs originate from retroviruses that infected the germ line of humans or human ancestors up to 30 million years ago[8]. Therefore, they became dormant in the genome and are inherited in a Mendelian fashion. Newly integrated endogenous retroviruses (ERVs) share the typical structure of the viruses they originate from. This means they contain at least the four main genes in the order $5' - gag - pro - pol - env - 3'$. Gag encodes the matrix and capsid proteins. Pro codes for a protease. Pol is the reverse transcriptase and integrase and env codes for the surface proteins. Furthermore, they have two flanking LTRs that originate from a duplication of the sequence, where the retrovirus was integrated[8].

However, endogenous retroviruses in human are generally not able to replicate. This is due to a high level of degeneration of the HERV gene sequences due to mutations and deletion introducing frame shifts and premature stop codons. Homologous recombination of the flanking LTRs can also lead to elements containing only a single LTR and losing the entire coding body[8]. Additionally, the promoter activity of LTRs is commonly silenced by hypermethylation[1].

Regardless of their loss of capability to replicate there still are multiple ways in which HERVs can influence human cells. LTRs of HERVs can act as promoters or enhancers to regulate gene expression. This has been associated to various cancers[9]. While most HERV genes are degenerated, some are still expressed. A change in the transcription of several HERV subclasses was found in schizophrenia subjects, but the mechanisms of HERV involvement are not yet well understood[10]. Furthermore, non-coding HERV RNAs have been found to play a role in control of stem cell properties[7].

1.3 Effect network analysis

- many bioinformatics methods find correlations, but not direct cause
- attempt to discern direct connections from bigger data webs
- hope to find possible biological mechanisms of gene regulation = path in model
- used approach: Gaussian Graphical models

Chapter 2

Data

2.1 HERV annotation

HERV annotation was pertained from RepeatMasker[11] repeat library. RepeatMasker is a tool that screens DNA sequences against a library interspersed repeats and low complexity DNA sequences. It generates an annotation of identified repeats and masks them in the query sequence.

The track representing all identified elements from the RepeatMasker library for human genome hg19 was downloaded from the UCSC genome browser download section (<http://hgdownload.soe.ucsc.edu/goldenPath/hg19/database/rmsk.txt.gz>). It contains a total of 5298130 occurrences of repeats. Each entry consists of 17 values. These are repeat name, the repeat class and family, as well as the chromosome, strand, and the genomic start and end position of the repeat occurrence. Furthermore it contains the area of the known repeat sequence, that is covered by the occurrence. It also contains the quality of the alignment of the repeat sequence to the annotated position using the Smith Waterman alignment score[12] and the number of base mismatches, deletions and insertions per thousand base pairs. Finally, there is an indexing field used to speed up chromosome range queries and the first digit of the id field in the RepeatMasker output file. The first 10 lines of the track are shown in table 2.1.

To extract HERV elements the annotation was filtered on different columns generating three sets of variable size. Multiple HERV sets were created as an attempt to cover different possible definitions of HERV elements.

This work only discerned between different HERV element types by defining different HERV sets. Therefore, within each set annotations, whose genomic positions overlap or are directly adjacent, were merged into one element.

A first set, HERV set 1 (HERV S1) was constructed by extracting all elements that contained "ERV" in the repeat name column. This set The resulting 42508 annotations

bin	swScore	milliDiv	milliDel	milliIns	genoName	genoStart	genoEnd	genoLeft	strand	repName	repClass	repFamily	repStart	repEnd	repLeft	id
585	1,504	13	4	13	chr1	10,000	10,468	$-2.49 \cdot 10^8$	+	(CCCTAA)n	Simple_repeat	Simple_repeat	1	463	0	1
585	3,612	114	270	13	chr1	10,468	11,447	$-2.49 \cdot 10^8$	-	TAR1	Satellite	telo	-399	1,712	483	2
585	437	235	186	35	chr1	11,503	11,675	$-2.49 \cdot 10^8$	-	L1MC	LINE	L1	-2,236	5,646	5,449	3
585	239	294	19	10	chr1	11,677	11,780	$-2.49 \cdot 10^8$	-	MER5B	DNA	hAT-Charlie	-74	104	1	4
585	318	230	38	0	chr1	15,264	15,355	$-2.49 \cdot 10^8$	-	MIR3	SINE	MIR	-119	143	49	5
585	203	162	0	0	chr1	16,712	16,749	$-2.49 \cdot 10^8$	+	(TGG)n	Simple_repeat	Simple_repeat	1	37	0	6
585	239	338	148	0	chr1	18,906	19,048	$-2.49 \cdot 10^8$	+	L2a	LINE	L2	2,942	3,104	-322	7
585	652	346	85	42	chr1	19,947	20,405	$-2.49 \cdot 10^8$	+	L3	LINE	CR1	3,042	3,519	-970	8
585	270	331	7	27	chr1	20,530	20,679	$-2.49 \cdot 10^8$	+	Plat_L3	LINE	CR1	2,802	2,947	-639	9
585	254	279	47	39	chr1	21,948	22,075	$-2.49 \cdot 10^8$	+	MLT1K	LTR	ERV1-MaLR	15	142	-453	1

Table 2.1: First ten rows of the RepeatMasker annotation on hg19

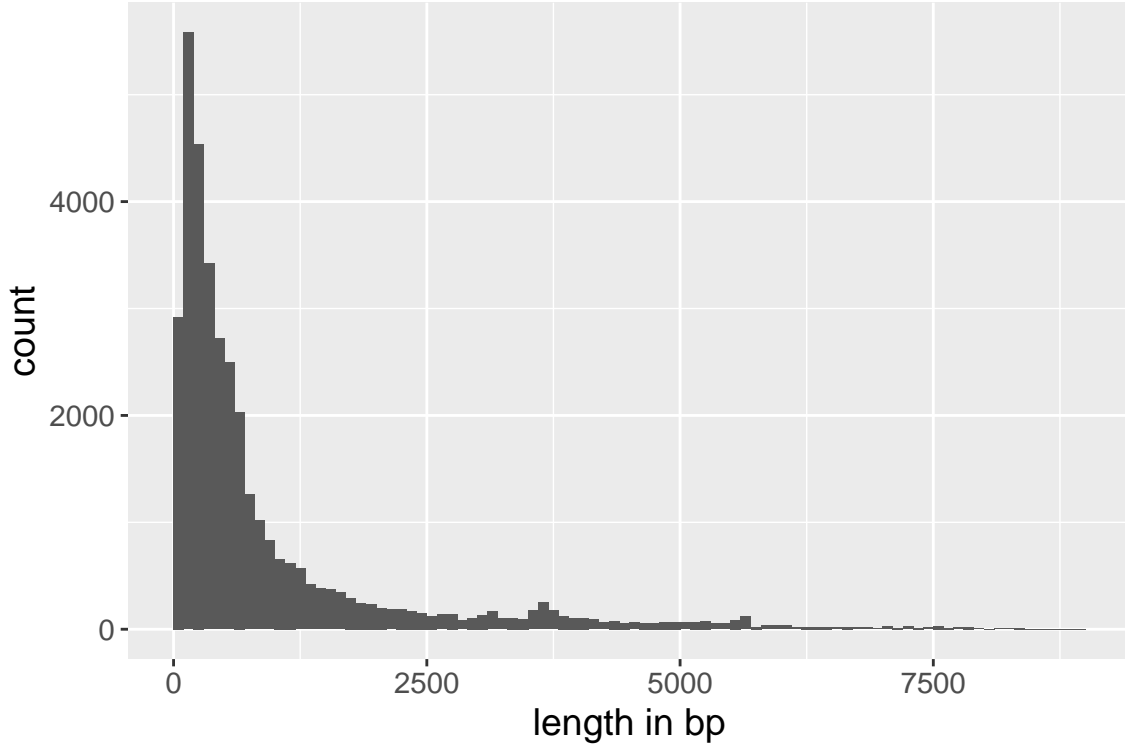


Figure 2.1: Length distribution of elements in HERV S2”

condensed to 35358 elements after merging. The elements have a mean width of 956 bp and cover a total of 33.8 Mbp, which is about 1.04% of the human genome. The distribution of element lengths in HERV S1 is shown in Figure 2.1

Alternatively, filtering the annotation for ”ERV” in the super family column leads to 696689 annotations. HERV set 2 (HERV S2) contains all endogenous retroviral sequences found in the human genome. Therefore, it is the most comprehensive set and it is the focus of detailed results that are shown and discussed. After merging overlapping and adjacent annotations this led to 612594 elements. Their mean width is 430 bp and they make up to 263.4 Mbp or ca 8.13% of the human genome.

A third set ”HERV set 3” (HERV S3) was constructed by filtering the repeat name column for ”HERV”, which resulted in 21361 annotations. As this set contains only elements that are explicitly named ”HERV”, we expect it to contain elements, that were inserted directly into the human genome. Merging overlapping and adjacent annotations resulted in 15284 elements with an average width of 1480 bp and a combined length of 22.6 Mbp.

2.2 KORA

The functional data for our analysis of gene regulatory mechanisms are made up up expression, methylation and genotype data. These were collected in the light of the platform for Cooperative Health Research in the Region of Augsburg, short KORA. It contains health surveys as well as examinations of individuals of German nationality living in the area of Augsburg, Bavaria. The objective of KORA is to track changes in health

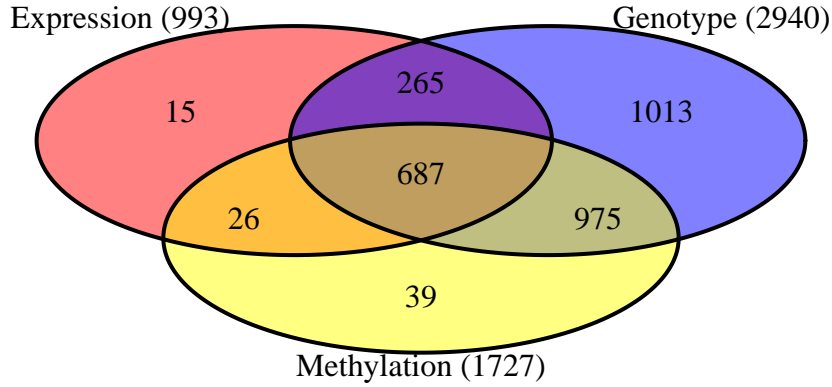


Figure 2.2: Number of samples with genotype, expression and methylation measurements in KORA F4 Survey

conditions over a long period in order to identify and examine the causes, effects and development of chronic diseases.

The data used in this work originates from the KORA F4 Survey, which was conducted from 2006 to 2008 and comprised samples of 3080 individuals. Of these data was available for 3020 individuals. F4 is a follow up study to the KORA S4 Survey performed from 1999 to 2001. It contains 4261 individuals.

All measurements were performed on whole blood samples. Houseman blood counts[13] describing the composition of different cell types for each individual had previously been calculated from the methylation data. They were used to weight certain cell type specific data.

Not all assays are available for all samples. Therefore different analyses were performed on varying sets of individuals according to availability of the required data types. A diagram of the samples available for each assay can be seen in figure 2.2.

2.2.1 Expression

The expression data was generated using the HumanHT-12 v3.0 Gene Expression Bead-Chip. The chip can measure expression values for 49576 probes. However, only 47864 probes represent an actual genomic location. The remaining probes are control probes used to assess the quality of the measurements and infer the background for measurements.

Measurements for 993 individuals are available from the KORA F4 survey. The comprise values for a total of 48803 probes per sample. Probes that do not map to a genomic location were excluded in all analyses, leaving 47864 probes. Of these 29521 are annotated to a total of 19288 genes.

The available data had already been quantile normalized, corrected for the background and log-transformed.

Figure 2.3 shows the distribution and coefficient of variation over all probes and samples. Apart from a tail of some probes with very high expression values figure 2.3A roughly resembles a standard distribution.

Probes that could not be mapped to genes were not discarded to avoid losing data for HERV regions, which are only sparsely annotated with genes and were the focus of this

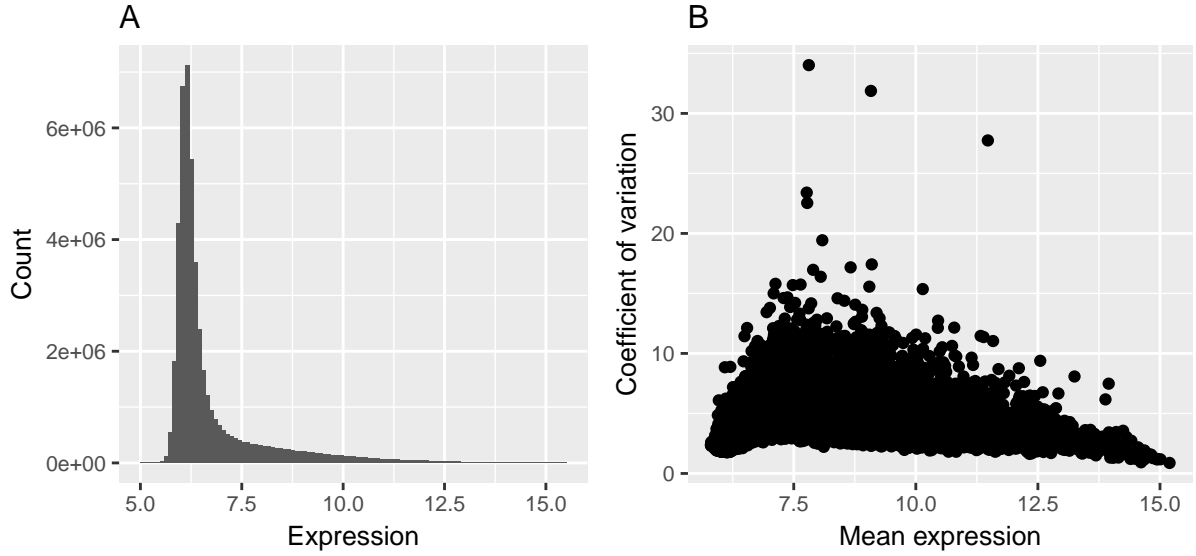


Figure 2.3: Distribution of expression values (A) and coefficient of expression variation (B) of 47864 probes in 993 individuals.

work and. Therefore, most analyses were performed on probe level or a mix of probes and genes.

2.2.2 Methylation

DNA methylation was measured using the Infinium HumanMethylation450K BeadChip, which interrogates methylation levels at 485577 genomic locations. Methylation intensities had already been corrected and transformed to beta values. Beta values have a range from 0 to 1 and represent the fraction of copies of a CpG site that are methylated in a sample.

No imputation of missing measurements was performed for the methylation data. Therefore, only 44335 probes have measurements in all samples available. Overall there are about 11.2 million missing values, which makes up 1.3% of all measurements and on average 23 samples missing per probe.

Methylation data was available for 1727 individuals and 485512 sites, which make up all 'cg' and 'ch' probe type probes. The distribution of beta values and the variances over all samples and probes is shown in figure 2.4. As seen in figure 2.4A, the interrogated CpG sites tend to be either entirely methylated or unmethylated within single samples. However, as B shows there is some variance between individuals in most CpG-sites. The low variance for very low and very high betas is explained by the fact that for a CpG site to have a mean close to 0 all samples have to have a value close to 0. Therefore, the variance also is very low.

2.2.3 Genotypes

Genotyping was performed with the Affymetrix Axiom array. On the data set used in this work genotypes had already been called using the Illuminus calling algorithm and missing values had been imputed using the IMPUTE2 software[14]. Furthermore, imputation results had already been filtered at IMPUTE value of 0.4.

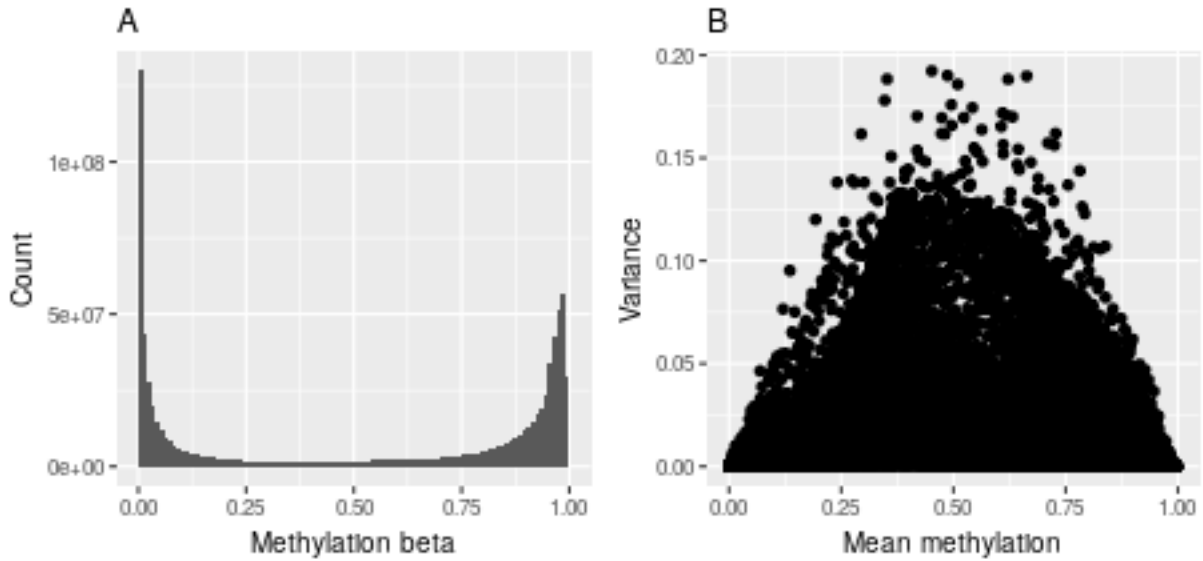


Figure 2.4: Distribution of all methylation beta values (A) and variance (B) of 485577 CpGs between 1727 individuals.

Additionally SNPs with a minor allele frequency of less than one percent had been removed, which allows used association analyses to be have more powerful test statistics.

In total measurements of 9533127 SNPs for 2940 individuals were available in the form of continuous values from 0 to 2 with 0 representing no occurrence of the SNP and 2 meaning the SNP is present in both chromosomes. Non integer values mean that different genotypes were measured in different cells in the sample and the resulting value describes the fraction of the occurrence of the variant.

2.2.4 Covariates

Several covariates were known for each sample. They were used to correct expression and methylation values for common effects.

The sex of all participants was known. Of the 3020 participants with any available data 1458 were SEX1 and the remaining 1562 were SEX2. The 697 individuals that had all three essays available contained 348 men/women and 339 women/men. Further covariates specific to participating individuals are sex, age, body mass index (BMI) and white blood cell count. An overview of these covariate values for all 3020 individuals and the individuals with all three essays available is shown in table 2.2.

Covariate	Minimum	Maximum	Mean	Covariate	Minimum	Maximum	Mean
Age	31	82	56.37	Age	61	81	69.04
BMI	16.05	55.99	27.65	BMI	18.99	47.58	28.87
WBC	2.5	40.9	5.98	WBC	2.7	12.6	5.9

Table 2.2: Covariate overview

Finally, two experimental factors for the expression measurements, storage time and RNA integrity number (RIN) and plate, were known.

2.2.5 Methylation quantitative trait loci

2.2.6 Data:meQTL

Genotypic variants are known to effect methylation patterns in humans [1]. Therefore, we used previously processed methylation quantitative trait loci (meQTL) data to explore the mechanism of trans acting SNP-CpG associations. Associated pairs were used as a seed for trans-acting interaction networks.

For meQTLs a stringent definition of trans effects was used: A SNP-CpG pair was called a trans-meQTL only if they are located on different chromosomes. This definition differs from the one used for eQTM and eQTL calculation as described later.

The data set contained a total of 11165559 significantly associated SNP-CpG pairs. 70709 distinct CpG-sites and 2709428 SNPs were part of at least one meQTL. Of these 467915 pairs consisting of 3592 CpG-sites and 200761 SNPs were between different chromosomes.

2.3 Transcription factor binding

DNA methylation plays a major role in the specifics of transcription factor binding sites[15]. We used transcription factor binding sites obtained from two publicly available sources to enrich information about CpG-sites of interest with transcription factors that bind nearby:

The first source was the third version of the track "Transcription Factor ChIP-seq (161 factors) from ENCODE with Factorbook Motifs"[16] downloaded from the UCSC genome browser download section (<http://hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeRegTfbsClustered/wgEncodeRegTfbsClusteredWithCellsV3.bed.gz>).

It combines 690 high quality ENCODE ChIP-seq data sets, which were processed with the Factorbook motif discovery and annotation pipeline[16]. The pipeline uses the tools MEME-ChIP[17] and FIMO[18] from the MEME software suite and merges discovered motifs with known motifs from Jaspar[19] and TransFac[20] using machine learning methods and manual curation.

The track contains a total of 438044 distinct peaks for 161 transcription factors in 91 cell types. For our analyses we filtered and combined the peaks for 23 blood related cell types. This leaves a total of 2173371 peaks for 125 transcription factors.

The second source was the ReMap project[21]. It combines 395 publicly available ChIP-seq data sets covering 132 different transcription factors across 83 cell lines. ReMap uses Bowtie2[2] to map reads to the human genome and the tool MACS[3] for peak calling. The final data set was downloaded from the ReMap website(http://tagc.univ-mrs.fr/remap/download/All/filPeaks_public.bed.gz).

There were a total of xxxxxx peaks for xxx transcription factors in the data set. After filtering for 19 blood related cell types xxxxx peaks and xxx different transcription factors remained.

Combining both filtered data sets lead to a total of 3847561 peaks of 145 different transcription factors.

2.4 Chromatin states

2.5 Data:Chromatin states

Chromatin state annotations were downloaded from Roadmap Epigenomics Core 15-state model (<http://egg2.wustl.edu/roadmap/data/byFileType/chromhmmSegmentations/ChmmModels/coreMarks/jointModel/final>). The Model provides a whole genome chromatin state annotation of 200 bp wide windows to the following 15 states: Active Transcription Start Site (TSS), Flanking Active TSS, Transcription at gene 5' and 4', Strong transcription, Weak transcription, Genic enhancers, Enhancers, ZNF genes & repeats, Heterochromatin, Bivalent/Poised TSS, Flanking Bivalent TSS/Enhancer, Bivalent Enhancer, Repressed PolyComb, Weak Repressed PolyComb and Quiescent/Low. The model is available for 127 diverse cell lines.

It was generated using ChromHMM v1.10[22] on the chromatin marks H3K4me1, H3K4me3, H3K27me3, H3K9me3, and H3K36me3. ChromHMM is based on a multivariate Hidden Markov Model.

In this work the annotations for 27 blood related cell lines were used. When combining the annotation for a genomic location of interest like a HERV element or a SNP, houseman blood counts were used to account of the cell composition of the whole blood samples.

Chapter 3

Methods

Most computations were done using the statistical computing environment and programming language R, version 3.4.1[23]. Pre filtering steps on large files, that could not be loaded into random access memory, were performed using the bash command "awk", which runs scripts written in the "AWK" programming language.

3.1 Overlaps

The main objective of this thesis was to analyze the effects of HERV elements on the regulation of cell activity.

To identify functional features like expression probes, genes, SNPs and CpG-sites that were associated with HERV elements, overlaps between these features and elements were calculated. Features were considered to be of interest for the analysis of an HERV element, if they overlapped by at least one base pair.

The calculation was performed using the function "findOverlaps" from the Bioconductor package "GenomicRanges"[24]. The function allows to identify overlaps between "GRanges" objects. These were constructed from BED files for all HERV sets as well as SNPs and retrieved from the Bioconductor packages available for analyses on the expression and methylation essays, "illuminaHumanv3.db"[25] and "FDb.InfiniumMethylation.hg19"[26] respectively.

The search for overlapping features was also performed on "GRanges" objects including all HERV elements as well as the 1kb or 2kb flanking regions of all elements. This was done to include functional features, that while not directly within HERV regions, are still likely to be influenced by them. Furthermore, this analysis was performed on all three HERV sets.

3.2 Data normalization

Before expression and methylation values were used for analyses, they were corrected for batch effects. For this we used the available covariates to calculate a linear models and regress out the respective residuals. The model used for expression values is shown in equation 3.1 and considers the covariates age, sex, RNA integrity number (RIN), plate and storage time. The model for methylation values, shown in equation 3.2, includes the cell composition for the five blood cell types known from the houseman blood counts, as

well as the first 20 principal components of the Illumina 450k array control-probes.

$$Expr \sim 1 + age + sex + RIN + plate + storage.time \quad (3.1)$$

$$Meth \sim 1 + CD4T + CD8T + NK + Bcell + Mono + \sum_{i=1}^{20} PC_i \quad (3.2)$$

3.3 eQTL/eQTM calculation

As a preliminary analysis to interrogate the effect of genotype and methylation variants on expression, expression quantitative trait loci (eQTL) and expression quantitative trait methylation (eQTM) were calculated.

Calculations were performed using the Bioconductor package MatrixEQTL[27]. MatrixEQTL tests for association of SNP-transcript pairs.

We set the parameter `useModel = modelLINEAR`, which leads to the use of an additive linear model. The association is modeled as simple linear regression and the absolute value of the sample correlation is used as test statistic.

After calculating the test statistics the p-values for all pairs that pass a defined significance threshold are calculated. These are corrected multiple testing using a Benjamini-Hochberg procedure[28], adapted for not recording all p-values, as shown in equation ??

MatrixEQTL is rather efficient because it manages to reduce the calculation of the sample correlation over all SNPs and transcripts to one single large matrix multiplication. This is achieved by transforming the genotype and transcription values.

MatrixEQTL also allows to include covariates in the QTL calculation. As the expression and methylation values used are residuals and therefore already corrected this option is not used. Furthermore, MatrixEQTL can differentiate between cis- and trans-interactions based on distance. The maximal distance to consider a pair on the same chromosome as cis was set to 500 kb.

The threshold for significant cis-QTLs during calculation was set to 10e-6 and 10e-8 for trans.

3.4 Functional Analysis of Gene Sets

In multiple analyses functional Gene Ontology enrichments were performed.

First a set of all GO annotations with any evidence code for gene symbols was retrieved from the Bioconductor package AnnotationDbi[29].

Then a hypergeometric test[30] for overrepresentation is performed on a set of genes of interest. For most enrichments a custom background set of genes specific to the analysis is given. Finally, the p-values for overrepresented GO terms were adjusted for multiple testing using the Holm method[31]. Terms with an adjusted p-value of less than 0.05 were considered significantly overrepresented.

3.5 Gaussian Graphical Models

Chapter 4

Results

4.1 Normalized Data

The distributions of the expression and methylation residuals, calculated as described in chapter 3.2, over all probes and samples can be seen in figure ?? . As expected the residuals follow a normal distribution. This is important for the calculation of the Gaussian graphical models, as the normal distribution of the data is one of base assumptions for Gaussian[].

4.2 HERV region features

In this section I will describe the expression probes that overlap with any HERV element and the cpgs and SNPs that lie within any HERV element and/or their flanking regions. The results for the set of all endogenous retroviral elements, HERV set 2, without flanking regions are described in detail. The results for the other sets defined in chapter 2.1 and including flanking regions will be shown in tables or in the supplementary data.

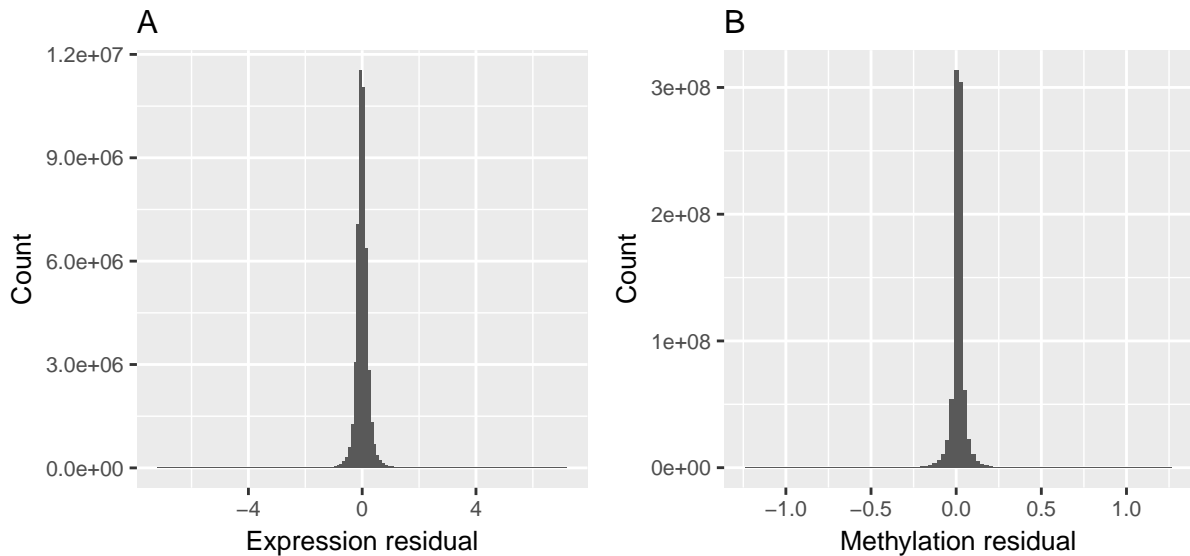


Figure 4.1: Distribution of expression (A) and methylation (B) residuals over all samples and probes

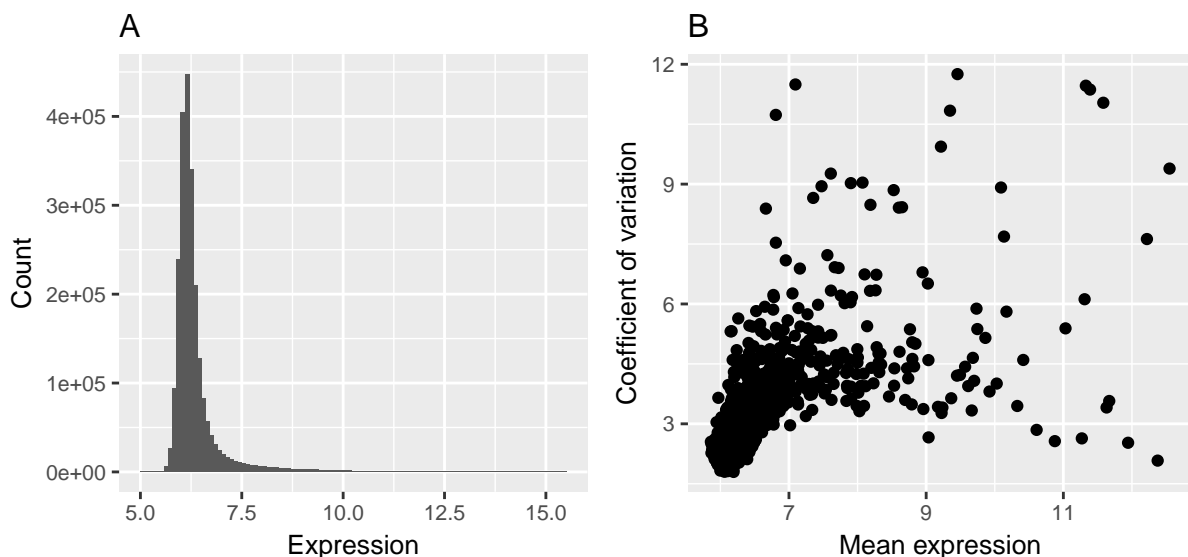


Figure 4.2: Distribution of expression values (A) and coefficient of expression variation (B) of 2343 expression probes overlapping with HERV S2 in 993 individuals.

4.2.1 Expression

A total of 2343 expression probes overlap directly with at least one of 2271 HERV elements from HERV S2. This means ca 4.50% of all expression probes overlap with HERV elements and expression measurements are known for parts of ca 0.37% of all HERV elements.

The distribution of expression values of the probes overlapping with HERV elements, as shown in figure 4.2A is almost identical to the distribution of all probes (figure 2.3A). The coefficient of variation (figure 4.2B), however, is lower on average for the considered subset compared to the whole set of probes (figure 2.3B).

Of these 2343 probes 510 were annotated to one of 449 different genes. I performed a GO enrichment for the biological process ontology on these genes with the set of all genes with available expression data as background. After correcting for multiple testing only the term defense response (GO:0006952, p-value=1.37e-5, fdr=0.047) was significantly enriched.

Term ID	Term	p	fdr
GO:0006952	defense response	$4.82 \cdot 10^{-11}$	$5.2 \cdot 10^{-7}$
GO:0045087	innate immune response	$5.5 \cdot 10^{-11}$	$5.94 \cdot 10^{-7}$
GO:0006955	immune response	$2.8 \cdot 10^{-9}$	$3.02 \cdot 10^{-5}$
GO:0098542	defense response to other organism	$1 \cdot 10^{-7}$	$1.08 \cdot 10^{-3}$
GO:0009615	response to virus	$3.2 \cdot 10^{-7}$	$3.45 \cdot 10^{-3}$
GO:0051607	defense response to virus	$3.45 \cdot 10^{-6}$	$3.73 \cdot 10^{-2}$

Table 4.1: Significantly enriched GO biological process terms among genes overlapping with HERV S2.

However, when including the 2kb flanking regions of HERV S2 and performing the GO enrichment on the 5518 identified genes, the six GO terms shown in table 4.1 are

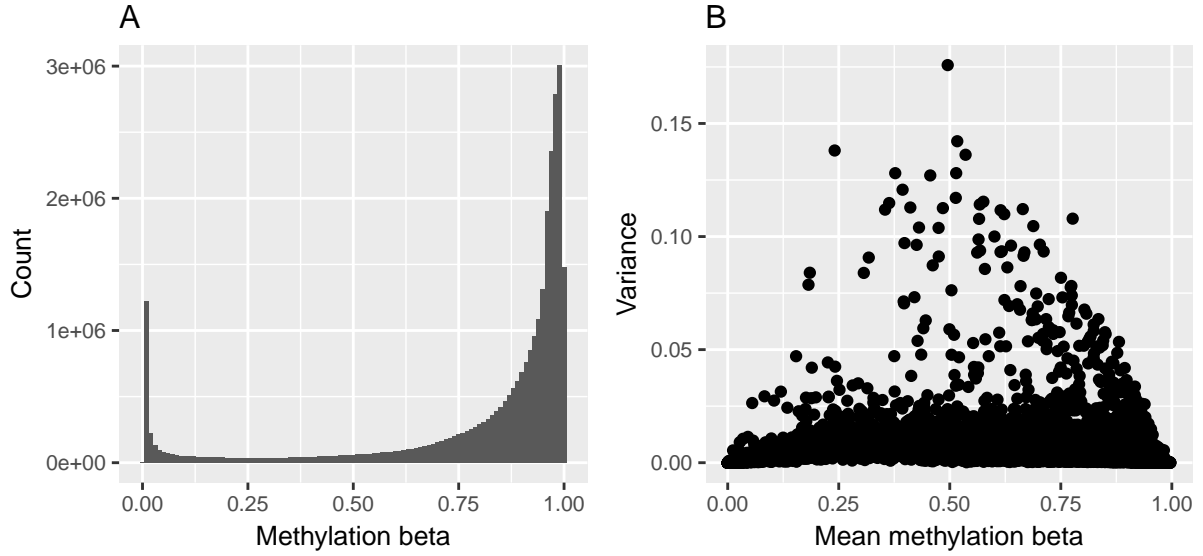


Figure 4.3: Distribution of methylation beta values (A) and coefficient of methylation variation (B) of 17077 GpG-sites located within HERV S2 in 1727 individuals.

significantly overrepresented. They are all connected to defense or immune response.

The results for all HERV sets and including flanking regions is shown in table 4.2.

Set	S1	S1.1kb	S1.2kb	S2	S2.1kb	S2.2kb	S3	S3.1kb	S3.2kb
Pairs	239	766	1,496	2,368	13,601	27,476	165	390	703
HERVs	219	706	1,336	2,271	12,551	24,275	146	350	615
Probes	239	563	970	2,343	9,364	15,466	165	302	487
Genes	21	131	293	449	2,968	5,517	15	69	138

Table 4.2: Number of expression probes overlapping with different HERV sets and flanking regions. "Pairs" is the total number of overlaps occurring, "HERVs" is the number of distinct HERV elements that have an overlap with any of the expression probes, Probes describes the number of distinct expression probes that overlap with the HERV elements or their flanking regions, "Genes" is the number of distinct Genes that are annotated to these probes.

4.2.2 Methylation

17077 CpG sites, equaling 3.52% of all measured CpGs, were found within HERV S2 and 12871 distinct HERV elements (2.10%) contain at least one interrogated CpG site.

CpGs associated with HERV S2, as shown in figure 4.3A, tend to be highly methylated. The variance also is (figure 4.3B) is on average lower than the one for all CpG-sites (figure 2.4B).

The number of CpG sites within all HERV sets and including flanking regions is shown in table 4.3

Set	S1	S1.1kb	S1.2kb	S2	S2.1kb	S2.2kb	S3	S3.1kb	S3.2kb
Pairs	1,587	6,105	12,005	17,077	$1.11 \cdot 10^5$	$2.39 \cdot 10^5$	1,152	3,465	6,408
HERVs	973	3,091	4,885	12,871	60,727	$1.01 \cdot 10^5$	614	1,561	2,425
CpGs	1,587	4,497	7,790	17,077	78,466	$1.39 \cdot 10^5$	1,152	2,671	4,404

Table 4.3: Number of CpGs overlapping with different HERV sets and flanking regions. "Pairs" is the total number of overlaps occurring, "HERVs" is the number of distinct HERV elements that have an overlap with any of the expression probes, "CpGs" is the number of distinct CpG sites that lie within the HERV elements or their flanking regions.

4.2.3 Genotypes

A total of 890780 the considered SNPs are located within elements of HERV S2. This constitutes 9.34% of all SNPs. These SNPs are found in 330744 distinct HERV elements. Therefore, 53.99% contain at least one SNP.

The results for all sets and flanking regions are shown in table 4.4

Set	S1	S1.1kb	S1.2kb	S2	S2.1kb	S2.2kb	S3	S3.1kb	S3.2kb
Pairs	$1.25 \cdot 10^5$	$3.84 \cdot 10^5$	$6.37 \cdot 10^5$	$8.91 \cdot 10^5$	$5.15 \cdot 10^6$	$9.35 \cdot 10^6$	89,066	$2.08 \cdot 10^5$	$3.21 \cdot 10^5$
HERVs	21,805	31,412	31,601	$3.31 \cdot 10^5$	$5.55 \cdot 10^5$	$5.59 \cdot 10^5$	10,139	13,125	13,189
SNPs	$1.25 \cdot 10^5$	$2.64 \cdot 10^5$	$3.72 \cdot 10^5$	$8.91 \cdot 10^5$	$3.31 \cdot 10^6$	$4.79 \cdot 10^6$	89,066	$1.54 \cdot 10^5$	$2.05 \cdot 10^5$

Table 4.4: Number of SNPs overlapping with different HERV sets and flanking regions. "Pairs" is the total number of overlaps occurring, "HERVs" is the number of distinct HERV elements that have an overlap with any of the expression probes, "SNPs" is the number of distinct considered SNPs that lie within the HERV elements or their flanking regions.

4.2.4 Chromatin states

Data in /storage/groups/groups_epigenreg/users/julian.schmidt ...

4.3 eQTLs

Associations were calculated for 156.4 millions cis acting SNP-expression probe pairs and 456.1 billion possible pairs in trans.

812147 of cis-pairs were significantly associated with p-values of $1e - 6$ or less. These significant eQTLs have a false discovery rate of less than $1.93e - 4$. A total of 551728 distinct SNPs and 4903 expression probes were part of at least one cis-eQTL. 4145 of these probes are annotated to 3552 different genes. The remaining 758 can not be assigned to a specific gene.

There were a total of 1511235 significant trans-eQTL with an p-value of less than $1e - 8$ and a FDR of less than $3.02e - 4$. These are made up by 229332 distinct SNPs and 21338 expression probes. 11505 of the probes found in at least one trans-eQTL are annotated to 9389 different genes, while 9767 probes are not assigned to a gene.

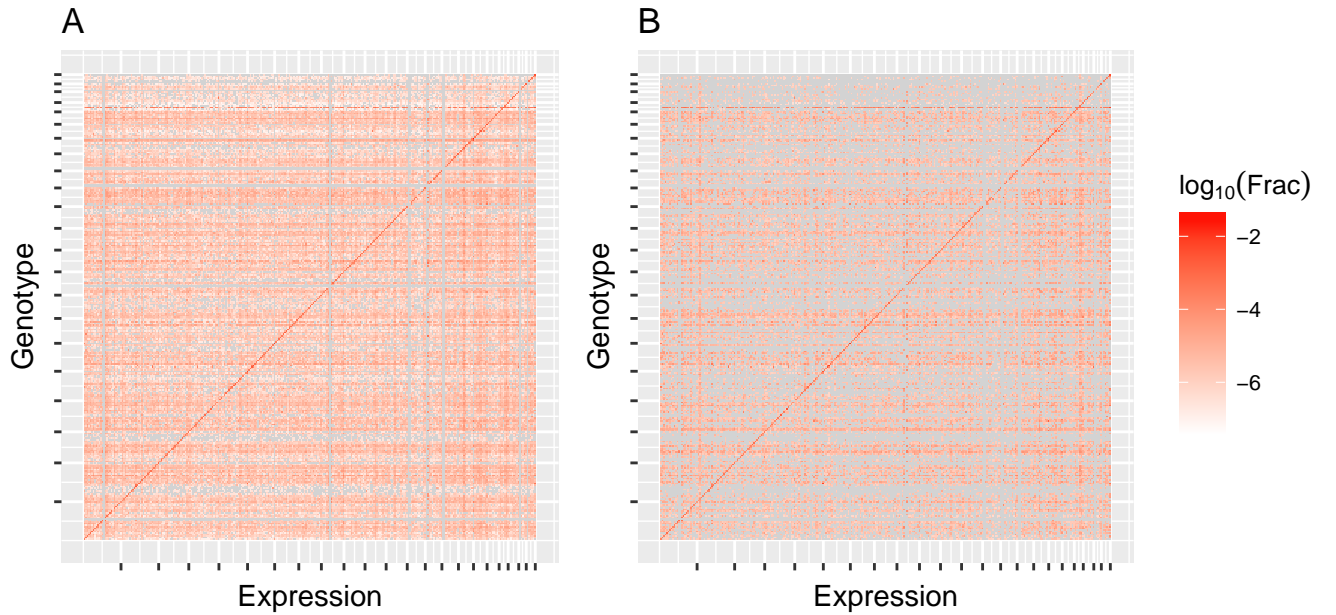


Figure 4.4: Fraction of SNP-expression probe pairs that were significantly associated sorted into bins according to their genomic locations. Ticks on the axis denote chromosome boundaries. Grey means there were no eQTLs in the given pair of bins. A considers all eQTLs an pairs, while B is limited on the ones whose SNP and/or expression probes are connected to HERV set 2.

In a previous work Schramm et al. [32] calculated eQTLs on expression and genotype data from 890 KORA F4 samples. They found significant cis-eQTLs for 4116 probes and considered only the pair of the most strongly associated SNPs with the probe. Of these pairs 2680 were also found in my analysis.

Figure 4.4A shows the fraction of significant to all possible SNP-expression probe pairs, mapped to 10 Mpb windows. As expected the values at the diagonal are comparatively higher, containing all cis-eQTLs.

The signatures of wildtype, heterogenous and homogenous mutated SNP site for the two eQTLs with the best and worst p-values in cis and trans are shown in figure ???. All eQTLs show clear differences in expression between different genotypic variants.

45862 of the SNPs and 166 of the expression probes found in at least one cis-eQTL are located within HERV S2. When considering all significant associations, where either the SNP or the expression probe lie within a HERV element, a total of 112604 pairs remain. There are 5855 cis-eQTLs constituting one of 4748 SNPs within a HERV element controlling one of 128 expression probes, that overlap with a HERV element.

When limiting trans-eQTLs to the ones that contain a HERV S2 SNP (23396) and/or expression probe (1227), 240301 remain. In 9177 of these pairs both are HERV related.

The fraction of HERV S2 related SNP-expression probe pairs that proved to be significantly associated is shown in figure 4.4B. The distribution of this subset of eQTLs is very similar to the whole set. At least at this resolution there does not seem to be a difference between HERV regions and the whole genome when it comes to the probability of a SNP or expression probe to be part of an eQTL.

A GO enrichment was performed on all 8211 genes associated to HERV related eQTLs using the set of all genes found in eQTLs as background. However, there were no signifi-

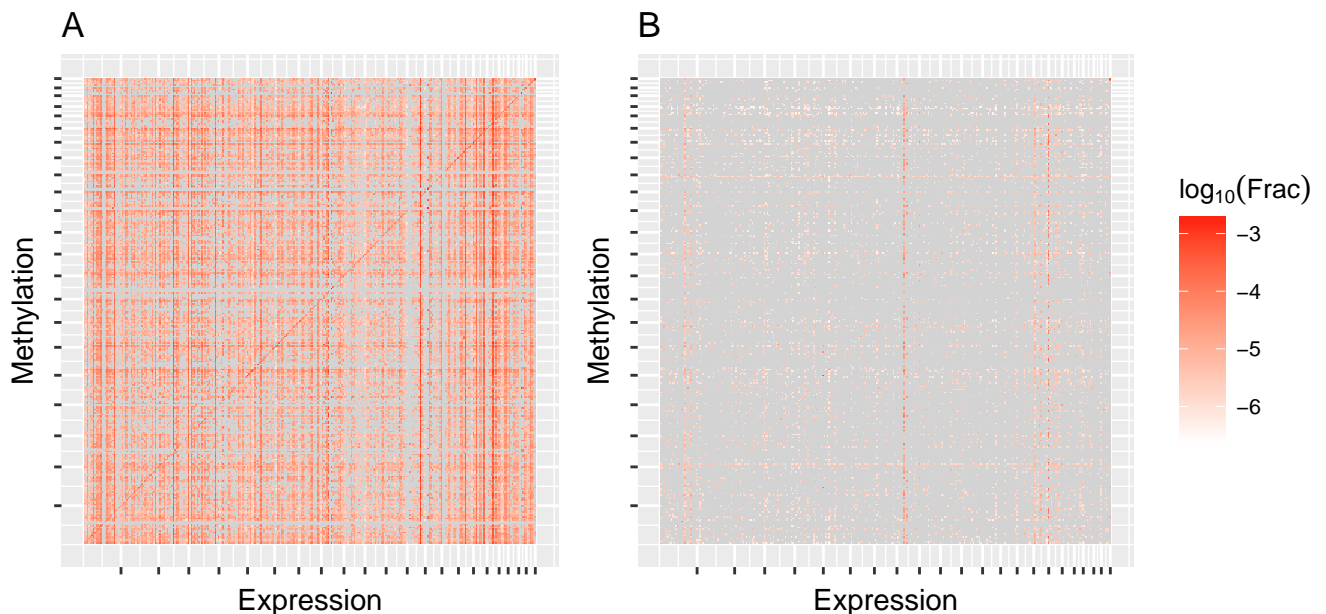


Figure 4.5: Fraction of CpG-expression probe pairs that were significantly associated sorted into bins according to their genomic locations. Ticks on the axis denote chromosome boundaries. Grey means there were no eQTM in the given pair of bins. A considers all eQTLs an pairs, while B is limited on the ones whose SNP and/or expression probes are connected to HERV set 2.

cantly enriched genes.

4.4 eQTM

When calculating expression quantitative trait methylation there were a total of 13.88 millions CpG-expression probe pairs within a distance of 50 Kpb or less of each other. The number of potential trans-acting pairs equaled around 23.22 billions.

Calculating eQTM with a significance threshold of $10e-6$ for cis resulted in 8187 significant associations ($fdr < 1.7e - 3$) consisting of 5957 distinct CpG sites and 1959 different expression probes. 1658 of these probes are annotated to 1461 genes, while there are no gene annotations for the remainder.

361485 of the potential trans-action CpG-expression probe pairs proved to be significant at a p-value threshold of $10e-8$ ($fdr < 6.5e - 4$). These trans-eQTLs are made up by 48206 CpGs and 11673 expression probes, for which 5738 gene annotations are available.

The fractions of potential CpG-expression probe pairs that were significantly associated between their genomic locations are shown in figure 4.5A. In contrast to the eQTL results there is only a weaker preference for cis interactions.

HERV S2 contains 311 CpG sites and 420 expression probes, that were present in cis-eQTM. A total of 738 cis-eQTM are related to the set. 33 of these are associations between one of 33 CpG sites within a HERV element and one of 14 expression probes.

Considering trans-eQTM, 1109 significantly associated CpGs lie within HERV S2 and 5422 expression probes overlapping a HERV element are part of a trans-eQTM.

The results of GO biological process enrichment performed on the xxxx genes that the

expression probes found in HERV S2 related eQTM are shown in table ??.

4.5 meQTLs

4.6 HERV related regulatory networks

4.6.1 Data collection

Chapter 5

Discussion

5.1 HERV region features

5.1.1 Expression

- less probes than expected – > chip design

5.1.2 Methylation

- less data available than expected? chip design or less CpG-sites?
- mostly methylated – > inactivated LTRs

5.1.3 Genotypes

Chapter 6

Conclusion and Outlook

Bibliography

- [1] Zachary D. Smith and Alexander Meissner. Dna methylation: roles in mammalian development. *Nature Reviews Genetics*, 14:204 EP –, Feb 2013. Review Article.
- [2] J. Craig Venter, Mark D. Adams, Eugene W. Myers, Peter W. Li, Richard J. Mural, Granger G. Sutton, Hamilton O. Smith, Mark Yandell, Cheryl A. Evans, Robert A. Holt, Jeannine D. Gocayne, Peter Amanatides, Richard M. Ballew, Daniel H. Huson, Jennifer Russo Wortman, Qing Zhang, Chinnappa D. Kodira, Xiangqun H. Zheng, Lin Chen, Marian Skupski, Gangadharan Subramanian, Paul D. Thomas, Jinghui Zhang, George L. Gabor Miklos, Catherine Nelson, Samuel Broder, Andrew G. Clark, Joe Nadeau, Victor A. McKusick, Norton Zinder, Arnold J. Levine, Richard J. Roberts, Mel Simon, Carolyn Slayman, Michael Hunkapiller, Randall Bolanos, Arthur Delcher, Ian Dew, Daniel Fasulo, Michael Flanigan, Liliana Florea, Aaron Halpern, Sridhar Hannenhalli, Saul Kravitz, Samuel Levy, Clark Mobarry, Knut Reinert, Karin Remington, Jane Abu-Threideh, Ellen Beasley, Kendra Biddick, Vivien Bonazzi, Rhonda Brandon, Michele Cargill, Ishwar Chandramouliswaran, Rosane Charlab, Kabir Chaturvedi, Zuoming Deng, Valentina Di Francesco, Patrick Dunn, Karen Eilbeck, Carlos Evangelista, Andrei E. Gabrielian, Weiniu Gan, Wang-mao Ge, Fangcheng Gong, Zhiping Gu, Ping Guan, Thomas J. Heiman, Maureen E. Higgins, Rui-Ru Ji, Zhaoxi Ke, Karen A. Ketchum, Zhongwu Lai, Yiding Lei, Zhenya Li, Jiayin Li, Yong Liang, Xiaoying Lin, Fu Lu, Gennady V. Merkulov, Natalia Milshina, Helen M. Moore, Ashwinikumar K Naik, Vaibhav A. Narayan, Beena Nee-lam, Deborah Nusskern, Douglas B. Rusch, Steven Salzberg, Wei Shao, Bixiong Shue, Jingtao Sun, Zhen Yuan Wang, Aihui Wang, Xin Wang, Jian Wang, Ming-Hui Wei, Ron Wides, Chunlin Xiao, Chunhua Yan, Alison Yao, Jane Ye, Ming Zhan, Weiqing Zhang, Hongyu Zhang, Qi Zhao, Liansheng Zheng, Fei Zhong, Wenyan Zhong, Shiaoping C. Zhu, Shaying Zhao, Dennis Gilbert, Suzanna Baumhueter, Gene Spier, Christine Carter, Anibal Cravchik, Trevor Woodage, Feroze Ali, Huijin An, Aderonke Awe, Danita Baldwin, Holly Baden, Mary Barnstead, Ian Barrow, Karen Beeson, Dana Busam, Amy Carver, Angela Center, Ming Lai Cheng, Liz Curry, Steve Danaher, Lionel Davenport, Raymond Desilets, Susanne Dietz, Kristina Dodson, Lisa Doup, Steven Ferriera, Neha Garg, Andres Gluecksmann, Brit Hart, Jason Haynes, Charles Haynes, Cheryl Heiner, Suzanne Hladun, Damon Hostin, Jarrett Houck, Timothy Howland, Chinyere Ibegwam, Jeffery Johnson, Francis Kalush, Lesley Kline, Shashi Koduru, Amy Love, Felecia Mann, David May, Steven McCawley, Tina McIntosh, Ivy McMullen, Mee Moy, Linda Moy, Brian Murphy, Keith Nelson, Cynthia Pfannkoch, Eric Pratts, Vinita Puri, Hina Qureshi, Matthew Reardon, Robert Rodriguez, Yu-Hui Rogers, Deanna Romblad, Bob Ruhfel, Richard Scott, Cynthia Sitter, Michelle Smallwood, Erin Stewart, Renee Strong, Ellen Suh, Reginald Thomas, Ni Ni Tint, Sukyee Tse, Claire Vech, Gary Wang, Jeremy Wetter,

- Sherita Williams, Monica Williams, Sandra Windsor, Emily Winn-Deen, Keriellen Wolfe, Jayshree Zaveri, Karena Zaveri, Josep F. Abril, Roderic Guigó, Michael J. Campbell, Kimmen V. Sjolander, Brian Karlak, Anish Kejariwal, Huaiyu Mi, Betty Lazareva, Thomas Hatton, Apurva Narechania, Karen Diemer, Anushya Muruganujan, Nan Guo, Shinji Sato, Vineet Bafna, Sorin Istrail, Ross Lippert, Russell Schwartz, Brian Walenz, Shibu Yooseph, David Allen, Anand Basu, James Baxendale, Louis Blick, Marcelo Caminha, John Carnes-Stine, Parris Caulk, Yen-Hui Chiang, My Coyne, Carl Dahlke, Anne Deslattes Mays, Maria Dombroski, Michael Donnelly, Dale Ely, Shiva Esparham, Carl Fosler, Harold Gire, Stephen Glanowski, Kenneth Glasser, Anna Glodek, Mark Gorokhov, Ken Graham, Barry Gropman, Michael Harris, Jeremy Heil, Scott Henderson, Jeffrey Hoover, Donald Jennings, Catherine Jordan, James Jordan, John Kasha, Leonid Kagan, Cheryl Kraft, Alexander Levitsky, Mark Lewis, Xiangjun Liu, John Lopez, Daniel Ma, William Majoros, Joe McDaniel, Sean Murphy, Matthew Newman, Trung Nguyen, Ngoc Nguyen, Marc Nodell, Sue Pan, Jim Peck, Marshall Peterson, William Rowe, Robert Sanders, John Scott, Michael Simpson, Thomas Smith, Arlan Sprague, Timothy Stockwell, Russell Turner, Eli Venter, Mei Wang, Meiyuan Wen, David Wu, Mitchell Wu, Ashley Xia, Ali Zandieh, and Xiaohong Zhu. The sequence of the human genome. *Science*, 291(5507):1304–1351, 2001.
- [3] International Human Genome Sequencing Consortium. Initial sequencing and analysis of the human genome. *Nature*, 409:860 EP –, Feb 2001.
 - [4] The ENCODE Project Consortium. An integrated encyclopedia of dna elements in the human genome. *Nature*, 489:57 EP –, Sep 2012. Article.
 - [5] Elizabeth Pennisi. Encode project writes eulogy for junk dna. *Science*, 337(6099):1159–1161, 2012.
 - [6] Todd J. Treangen and Steven L. Salzberg. Repetitive dna and next-generation sequencing: computational challenges and solutions. *Nature Reviews Genetics*, 13:36 EP –, Nov 2011. Review Article.
 - [7] Robin A. Weiss. Human endogenous retroviruses: friend or foe? *APMIS*, 124(1-2):4–10, 2016.
 - [8] Norbert Bannert and Reinhard Kurth. The evolutionary dynamics of human endogenous retroviral families. *Annual Review of Genomics and Human Genetics*, 7(1):149–173, 2006. PMID: 16722807.
 - [9] H. Yu, Z. Zhao, and F. Zhu. The role of human endogenous retroviral long terminal repeat sequences in human cancer (review). *International Journal of Molecular Medicine*, 32(4):755–762, 2013.
 - [10] Gorjan Slokar and Gregor Hasler. Human endogenous retroviruses as pathogenic factors in the development of schizophrenia. 6, 01 2016.
 - [11] R Smit AFA, Hubley and Green P. Repeatmasker open-4.0. 2013-2015.
 - [12] T.F. Smith and M.S. Waterman. Identification of common molecular subsequences. *Journal of Molecular Biology*, 147(1):195 – 197, 1981.

- [13] Lindsay L. Waite, Benjamin Weaver, Kenneth Day, Xinrui Li, Kevin Roberts, Andrew W. Gibson, Jeffrey C. Edberg, Robert P. Kimberly, Devin M. Absher, and Hemant K. Tiwari. Estimation of cell-type composition including t and b cell subtypes for whole blood methylation microarray data. *Frontiers in Genetics*, 7:23, 2016.
- [14] Bryan N. Howie, Peter Donnelly, and Jonathan Marchini. A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLOS Genetics*, 5(6):1–15, 06 2009.
- [15] Yimeng Yin, Ekaterina Morgunova, Arttu Jolma, Eevi Kaasinen, Biswajyoti Sahu, Syed Khund-Sayeed, Pratyush K. Das, Teemu Kivioja, Kashyap Dave, Fan Zhong, Kazuhiro R. Nitta, Minna Taipale, Alexander Popov, Paul A. Ginno, Silvia Domcke, Jian Yan, Dirk Schübeler, Charles Vinson, and Jussi Taipale. Impact of cytosine methylation on dna binding specificities of human transcription factors. *Science*, 356(6337), 2017.
- [16] Jie Wang, Jiali Zhuang, Sowmya Iyer, XinYing Lin, Troy W. Whitfield, Melissa C. Greven, Brian G. Pierce, Xianjun Dong, Anshul Kundaje, Yong Cheng, Oliver J. Rando, Ewan Birney, Richard M. Myers, William S. Noble, Michael Snyder, and Zhiping Weng. Sequence features and chromatin structure around the genomic regions bound by 119 human transcription factors. *Genome Res*, 22(9):1798–1812, Sep 2012.
- [17] Philip Machanick and Timothy L. Bailey. Meme-chip: motif analysis of large dna datasets. *Bioinformatics*, 27(12):1696–1697, 2011.
- [18] Charles E. Grant, Timothy L. Bailey, and William Stafford Noble. Fimo: scanning for occurrences of a given motif. *Bioinformatics*, 27(7):1017–1018, 2011.
- [19] Aziz Khan, Oriol Fornes, Arnaud Stigliani, Marius Gheorghe, Jaime A Castro-Mondragon, Robin van der Lee, Adrien Bessy, Jeanne ChÃˆneby, Shubhada R Kulkarni, Ge Tan, Damir Baranasic, David J Arenillas, Albin Sandelin, Klaas Vandepoele, Boris Lenhard, BenoÃ®t Ballester, Wyeth W Wasserman, FranÃ§ois Parcy, and Anthony Mathelier. Jaspar 2018: update of the open-access database of transcription factor binding profiles and its web framework. *Nucleic Acids Research*, 46(D1):D260–D266, 2018.
- [20] V. Matys, O. V. Kel-Margoulis, E. Fricke, I. Liebich, S. Land, A. Barre-Dirrie, I. Reuter, D. Chekmenev, M. Krull, K. Hornischer, N. Voss, P. Stegmaier, B. Lewicki-Potapov, H. Saxel, A. E. Kel, and E. Wingender. Transfac(?) and its module transcompel(?): transcriptional gene regulation in eukaryotes. *Nucleic Acids Res*, 34(Database issue):D108–D110, Jan 2006. 16381825[pmid].
- [21] Aurelien Griffon, Quentin Barbier, Jordi Dalino, Jacques van Helden, Salvatore Spicuglia, and Benoit Ballester. Integrative analysis of public chip-seq experiments reveals a complex multi-cell regulatory landscape. *Nucleic Acids Research*, 43(4):e27, 2015.
- [22] Jason Ernst and Manolis Kellis. Chromhmm: automating chromatin-state discovery and characterization. *Nature Methods*, 9:215 EP –, Feb 2012. Correspondence.

- [23] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2017.
- [24] Michael Lawrence, Wolfgang Huber, Hervé Pagès, Patrick Aboyoun, Marc Carlson, Robert Gentleman, Martin Morgan, and Vincent Carey. Software for computing and annotating genomic ranges. *PLoS Computational Biology*, 9, 2013.
- [25] Mark Dunning, Andy Lynch, and Matthew Eldridge. *illuminaHumanv3.db: Illumina HumanHT12v3 annotation data (chip illuminaHumanv3)*, 2015. R package version 1.26.0.
- [26] Tim Triche, Jr. *FDb.InfiniumMethylation.hg19: Annotation package for Illumina Infinium DNA methylation probes*, 2014. R package version 2.2.0.
- [27] Andrey A. Shabalin. Matrix eqtl: ultra fast eqtl analysis via large matrix operations. *Bioinformatics*, 28(10):1353–1358, 2012.
- [28] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, 57(1):289–300, 1995.
- [29] Hervé Pagès, Marc Carlson, Seth Falcon, and Nianhua Li. *AnnotationDbi: Annotation Database Interface*, 2017. R package version 1.38.1.
- [30] S Falcon and R Gentleman. Using GOstats to test gene lists for GO term association. *Bioinformatics*, 23(2):257–8, 2007.
- [31] Sture Holm. A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*, 6(2):65–70, 1979.
- [32] Katharina Schramm, Carola Marzi, Claudia Schurmann, Maren Carstensen, Eva Reinmaa, Reiner Biffar, Gertrud Eckstein, Christian Gieger, Hans-Jürgen Grabe, Georg Homuth, Gabriele Kastenmüller, Reedik Mägi, Andres Metspalu, Evelin Mihailov, Annette Peters, Astrid Petersmann, Michael Roden, Konstantin Strauch, Karsten Suhre, Alexander Teumer, Uwe Völker, Henry Völzke, Rui Wang-Sattler, Melanie Waldenberger, Thomas Meitinger, Thomas Illig, Christian Herder, Harald Grallert, and Holger Prokisch. Mapping the genetic architecture of gene regulation in whole blood. *PLOS ONE*, 9(4):1–13, 04 2014.