

Online methods

1

2

3 Contents

4	1 Discovery data	2
5	1.1 Data description	2
6	1.2 Normalisation	2
7	2 Exhaustive 2D-eQTL analysis	3
8	2.1 Two stage search	3
9	2.1.1 Stage 1	3
10	2.1.2 Stage 2	4
11	2.2 Type 1 error rate	4
12	2.3 Population stratification	5
13	2.4 Probe mapping	5
14	3 Replication	5
15	3.1 Data description	5
16	3.2 Meta Analysis	6
17	3.3 Concordance of direction of effects	6
18	4 Effects of LD on detection and replication	7
19	4.1 Simulation 1	7
20	4.2 Simulation 2	8
21	4.3 Interpretation	9
22	5 Additive and non-additive variance estimation	9
23	5.1 Fixed effects	9
24	5.2 Pedigree estimates	9
25	6 Functional enrichment analysis	10
26	6.1 Tissue specific transcriptionally active regions	10
27	6.2 Chromosome interactions	10
28	6.3 SNP colocalisation with genomic features	11
29	6.4 Transcription factor enrichment	11
30	6.5 Defining previously identified SNP associations	12
31	7 References	12

1 Discovery data

1.1 Data description

The Brisbane Systems Genetics Study (BSGS) comprises 846 individuals of European descent from 274 independent families.¹ DNA samples from each individual were genotyped on the Illumina 610-Quad Beadchip by the Scientific Services Division at deCODE Genetics Iceland. Full details of genotyping procedures are given in Medland et al.² Standard quality control (QC) filters were applied and the remaining 528,509 autosomal SNPs were carried forward for further analysis.

Gene expression profiles were generated from peripheral blood collected with PAXgene TM tubes (QIAGEN, Valencia, CA) using Illumina HT12-v4.0 bead arrays. The Illumina HT-12 v4.0 chip contains 47,323 probes, although some probes are not assigned to RefSeq genes. We removed any probes that did not match the following criteria: contained a SNP within the probe sequence with MAF > 0.05 within 1000 genomes data; did not map to a listed RefSeq gene; were not significantly expressed (based on a detection p -value < 0.05) in at least 90% of samples. After this stringent QC 7339 probes remained for 2D-eQTL mapping.

1.2 Normalisation

Gene expression profiles were normalised and adjusted for batch and polygenic effects. Profiles were first adjusted for raw background expression in each sample. Expression levels were then adjusted using quantile and \log_2 transformation to standardise distributions between samples. Batch and polygenic effects were adjusted using the linear model

$$y = \mu + \beta_1 c + \beta_2 p + \beta_3 s + \beta_4 a + g + e \quad (1)$$

where μ is the population mean expression levels, c , p , s and a are vectors of chip, chip position, sex and generation respectively, fitted as fixed effects; and g is a random additive polygenic effect with a variance covariance matrix

$$G_{jk} = \begin{cases} \sigma_a^2 & j = k \\ 2\phi_{jk}\sigma_a^2 & j \neq k \end{cases} \quad (2)$$

The parameter σ_a^2 is the variance component for additive background genetic. Here, we are using family based pedigree information rather than SNP based IBD to account for relationships between individuals and so ϕ_{jk} is the kinship coefficient between individuals j and k . The residual, e , from equation 1 is assumed to follow a multivariate normal distribution with a mean of zero. Residuals were normalised by rank transformation and used as the adjusted phenotype for the pairwise epistasis scan to remove any skewness and avoid results being driven by outliers. The GenABEL package for R was used to perform the normalisation.³

2 Exhaustive 2D-eQTL analysis

2.1 Two stage search

We used epiGPU⁴ software to perform an exhaustive scan for pairwise interactions, such that each SNP is tested against all other SNPs for statistical association with the expression values for each of the 7339 probes. This uses the massively parallel computational architecture of graphical processing units (GPUs) to speed up the exhaustive search. For each SNP pair there are 9 possible genotype classes. We treat each genotype class as a fixed effect and fit an 8 *d.f.* *F*-test to test the following hypotheses:

$$H_0 : \sum_{i=1}^3 \sum_{j=1}^3 (\bar{x}_{ij} - \mu)^2 = 0; \quad (3)$$

$$H_1 : \sum_{i=1}^3 \sum_{j=1}^3 (\bar{x}_{ij} - \mu)^2 > 0; \quad (4)$$

where μ is the mean expression level and x_{ij} is the pairwise genotype class mean for genotype i at SNP 1 and genotype j at SNP 2. This type of test does not parameterize for specific types of epistasis, rather it tests for the joint genetic effects at two loci. This has been demonstrated to be statistically more efficient when searching for a wide range of epistatic patterns, although will also include any marginal effects of SNPs which must be dealt with post-hoc.⁵

2.1.1 Stage 1

The complete exhaustive scan for 7339 probes comprises 1.03×10^{15} *F*-tests. We used permutation analysis to estimate an appropriate significance threshold for the study. To do this we performed a further 1600 exhaustive 2D scans on permuted phenotypes to generate a null distribution of the extreme *p*-values expected to be obtained from this number of multiple tests given the correlation structure between the SNPs. We took the most extreme *p*-value from each of the 1600 scans and set the 5% FWER to be the 95% most extreme of these *p*-values, $T_* = 2.13 \times 10^{-12}$. The effective number of tests in one 2D scan being performed is therefore $N_* = 0.05/T_* \approx 2.33 \times 10^{10}$. To correct for the testing of multiple traits we established an experiment wide threshold of $T_e = 0.05/(N_* \times 7339) = 2.91 \times 10^{-16}$. This is likely to be conservative as it assumes independence between probes.

Filtering We used two approaches to filter SNPs from stage 1 to be tested for significant interaction effects in stage 2.

Filter 1 After keeping SNP pairs that surpassed the 2.91×10^{-16} threshold in stage 1 only SNP pairs with at least 5 data points in all 9 genotype classes were kept. We then calculated the LD between interacting SNPs (amongst

unrelated individuals within the discovery sample and also from 1000 genomes data) and removed any pairs with $r^2 > 0.1$ or $D'^2 > 0.1$ to avoid the inclusion of haplotype effects and to increase the accuracy of genetic variance decomposition. If multiple SNP pairs were present on the same chromosomes for a particular expression trait then only the sentinel SNP pair was retained, *i.e.* if a probe had multiple SNP pairs that were on chromosomes one and two then only the SNP pair with the most significant p -value was retained. At this stage 6404 filtered SNP pairs remained.

Filter 2 We also performed a second filtering screen applied to the list of SNP pairs from stage 1 that was identical to filter 1 but an additional step was included where any SNPs that had previously been shown to have a significant additive or dominant effect ($p < 1.29 \times 10^{-11}$) were removed,⁶ creating a second set of 4751 unique filtered SNP pairs.

2.1.2 Stage 2

To ensure that interacting SNPs were driven by epistasis and not marginal effects we performed a nested ANOVA on each pair in the filtered set to test if the interaction terms were significant. We did this by contrasting the full genetic model (8 *d.f.*) against the reduced marginal effects model which included the additive and dominance terms at both SNPs (4 *d.f.*). Thus, a 4 *d.f.* F -test was performed on the residual genetic variation, representing the contribution of epistatic variance. Significance of epistasis was determined using a Bonferroni threshold of $0.05/(6404 + 4751) = 4.48 \times 10^{-6}$. This resulted in 406 and 95 SNP pairs with significant interaction terms from filters 1 and 2, respectively.

2.2 Type 1 error rate

Using a Bonferroni correction of 0.05 in the second stage of the two stage discovery scan implies a type 1 error rate of $\alpha = 0.05$. However, this could be underestimated because the number tests performed in the second stage depends on the number of tests in the first stage, and this depends on statistical power and model choice. We performed simulations to estimate the type 1 error rate of this study design.

We assumed a null model where there was one true additive effect and 7 other terms with no effect. To simulate a test statistic we simulated 8 z -scores, $z_1 \sim N(\sqrt{NCP}, 1)$ and $z_{2..8} \sim N(0, 1)$. Thus $z_{full} = \sum_{i=1}^8 z_i \sim \chi_8^2$ (representing the 8 d.f. test) and $z_{int} = \sum_{i=5}^8 z_i \sim \chi_4^2$ (representing the 4 d.f. test where the null hypothesis of no epistasis is true). For a particular value of NCP we simulated 100,000 z values, and calculated the p_{full} -value for the z_{full} test statistic. The n_{int} test statistics with $p_{full} < 2.31 \times 10^{-16}$ were kept for the second stage, where the type 1 error rate of stage 2 was calculated as the proportion of $p_{int} < 0.05/n_{int}$. The power at stage 1 was calculated as $n_{int}/100,000$. This procedure was performed for a range of NCP parameters that represented power ranging from ~ 0 to ~ 1 .

142 2.3 Population stratification

143 We ruled out population stratification as a possible cause of inflated test statis-
144 tics. To test for cryptic relatedness driving the interaction terms we tested for
145 increased LD among the SNPs.⁷ We calculated the mean of the off-diagonal
146 elements of the correlation matrix of all unique SNPs from the 501 interactions
147 (731 SNPs) using only unrelated individuals, $r^2 = 0.0039$. This is not signifi-
148 cantly different from the null hypothesis of zero (sampling error = $1/n_{\text{unrelated}} =$
149 0.0039).

150 2.4 Probe mapping

151 To avoid possibility that epistatic signals might arise due to expression probes
152 hybridising in multiple locations we verified that probe sequences for genes with
153 significant interactions mapped to only a single location. As an initial verifi-
154 cation we performed a BLAST search of the full probe sequence against 1000
155 genomes phase 1 version 3 human genome reference and ensured that only one
156 genomic location aligned significantly ($p < 0.05$). As a second step, to mitigate
157 the possibility of weak hybridisation elsewhere in the genome we divided the
158 probe sequence into three sections (1-25bp, 13-37bp, 26-50bp) and performed
159 a BLAST search of these probe sequence fragments. No probe sequences or
160 probe sequence fragments mapped to positions other than the single expected
161 genomic target ($p < 0.05$).

162 3 Replication

163 3.1 Data description

164 We attempted replication of the 501 significant interactions from the discovery
165 set using three independent cohorts; Fehrmann, EGCUT, and CHDWB. It was
166 required that LD $r^2 < 0.1$ and $D'^2 < 0.1$ between interacting SNPs (as measured
167 in the replication sample directly), and all nine genotype classes had at least 5
168 individuals present in order to proceed with statistical testing for replication in
169 both datasets. We also excluded any putative SNPs that had discordant allele
170 frequencies in any of the datasets. Details of the cohorts are as follows.

171 **Fehrmann:** $n = 1240$ The Fehrmann dataset⁸ consists of peripheral blood
172 samples of 1240 unrelated individuals from the United Kingdom and the Nether-
173 lands. Some of these individuals are patients, while others are healthy controls.
174 Individuals were genotyped using the Illumina HumanHap300, Illumina Human-
175 Hap370CNV, and Illumina 610 Quad platforms. RNA levels were quantified
176 using the Illumina HT-12 V3.0 platform.

177 **EGCUT:** $n = 891$ The Estonian Genome Center of the University of Tartu
178 (EGCUT) study⁹ consists of peripheral blood samples of 891 unrelated individ-

179 uals from Estonia. They were genotyped using the Illumina HumanHap370CNV
180 platform. RNA levels were quantified using the Illumina HT-12 V3.0 platform.

181 **CDHWB:** $n = 139$ The Center for Health Discovery and Well Being (CD-
182 HWB) Study¹⁰ is a population based cohort consisting of 139 individuals of
183 European descent collected in Atlanta USA. Gene expression profiles were gen-
184 erated with Illumina HT-12 V3.0 arrays from peripheral blood collected from
185 Tempus tubes that preserve RNA. Whole genome genotypes were measured us-
186 ing Illumina OmniQuad arrays. Due to the small sample size, most SNP pairs
187 did not pass filtering in this dataset (20 SNP pairs remained) and so we have
188 excluded it from the rest of the analysis.

189 3.2 Meta Analysis

190 The 4 *d.f.* interaction p -values for each independent replication dataset were
191 calculated using the same statistical test as was performed in the discovery
192 dataset. We then took the interaction p -values from EGCUT and Fehrman
193 and calculated a joint p -value using Fisher’s method of combining p -values for
194 a meta analysis as $-2 \ln p_1 - 2 \ln p_2 \sim \chi^2_{4d.f.}$. As in the discovery analysis, all
195 gene expression levels were normalised using rank transformation to avoid skew
196 or outliers in the distribution.¹¹

197 3.3 Concordance of direction of effects

198 We used four methods to calculate the concordance of the direction of effects
199 between the discovery and replication datasets.

200 **Test 1** Is the most significant epistatic effect in the discovery set in the same
201 direction as the same epistatic effect in the replication sets? We decomposed
202 the genetic variance into 8 orthogonal effects, four of which are epistatic ($A \times A$,
203 $A \times D$, $D \times A$, $D \times D$). The sign of the epistatic effect that had the largest
204 variance in the discovery was recorded, and then was compared to the same
205 epistatic effect in the two replication datasets (regardless of whether or not the
206 same epistatic effect was the largest in the replication datasets). The probability
207 of the sign being the same in one dataset is 1/2. The probability of the sign
208 being the same in two is 1/4.

209 **Test 2** Is the most significant epistatic effect in the discovery the same as
210 the largest epistatic effect in the replication set with the sign being concordant.
211 As in Test 1, but this time we required that the largest effect was the same in
212 the discovery and the replication, and that they had the same sign (*e.g.* if the
213 largest effect in the discovery is $A \times A$, with a positive effect, then concordance is
214 achieved if the same is true in the replication). The probability of one replication
215 dataset being concordant by chance is 1/8, and concordance in both is 1/64.

216 **Test 3** Do the epistatic effects that are significant at nominal $p < 0.05$ in
 217 the discovery have the same direction of effect as in the replication? Here we
 218 count all the epistatic variance components in the discovery that have $p < 0.05$
 219 (1133 amongst the 434 discovery SNP pairs, *i.e.* each SNP pair has at least 1
 220 and at most 4 significant epistatic variance components). Then we compare
 221 the direction of the effect in the replication dataset. The probability of the
 222 sign being the same in one dataset for any one significant effect is $1/2$. The
 223 probability of the sign being the same in two is $1/4$.

224 **Test 4** If we count how many of the 4 epistatic effects are concordant between
 225 the discovery and replication data for each interaction then is this significant
 226 from what we expect by chance? There can be either 0, 1, 2, 3 or 4 concordant
 227 signs at each interaction, each with expectation of $p = 1/16, 4/16, 6/16, 4/16, 1/16$
 228 under the null, respectively. Observed counts are multinomially distributed,
 229 and we tested if the observed proportions were statistically different from the
 230 expected proportions using an approximation of the multinomial test.¹²

231 The probability of observing the number of concordant signs in tests 1-3 is
 232 calculated using a binomial test. All variance decompositions were calculated
 233 using the NOIA method.¹³

234 4 Effects of LD on detection and replication

235 The power to detect genetic effects, when the observed markers are in LD with
 236 the causal variants, is proportional to r^x . For additive effects $x = 2$, but for
 237 non-additive effects x is larger, *i.e.* $x = 4$ for dominance or $A \times A$, $x = 6$ for
 238 $A \times D$ or $D \times A$, and $x = 8$ for $D \times D$. Many biologically realistic GP maps
 239 may be comprised of all 8 variance components.⁵

240 This is important for both detection and for replication of epistasis. For
 241 detection, if the epistatic effect includes the $D \times D$ term then if the two causal
 242 variants are tagged by observed markers that are each in LD $r = 0.9$, then if
 243 the true variance is V_t then the observed variance V_o at the markers will be
 244 $0.9^8 V_t = 0.43 V_t$. Therefore, it is important to consider the sampling variation
 245 of \hat{r}^x in a sample given some true population value of r .

246 4.1 Simulation 1

247 For some values of fixed population parameters, p_1 (minor allele frequency at
 248 observed marker), q_1 (minor allele frequency at causal variant), and r (LD
 249 between marker and causal variant), the expected haplotype frequencies are

$$h_{11} = r\sqrt{p_1 q_1 p_2 q_2} + p_1 q_1 \quad (5)$$

$$h_{12} = p_1 q_2 - r\sqrt{p_1 q_1 p_2 q_2} \quad (6)$$

$$h_{21} = p_2 q_1 - r\sqrt{p_1 q_1 p_2 q_2} \quad (7)$$

$$h_{22} = r\sqrt{p_1 q_1 p_2 q_2} + p_2 q_2 \quad (8)$$

where $p_2 = 1 - p_1$ and $q_2 = 1 - q_1$. For a range of population parameters we randomly sampled $2n$ haplotypes where the expected haplotype frequencies were $h_{11}, h_{12}, h_{21}, h_{22}$. From the sample haplotype frequencies we then calculated sample estimates of \hat{r} where

$$\hat{r} = \frac{\hat{h}_{11} - \hat{p}_1 \hat{q}_1}{\sqrt{\hat{p}_1 \hat{q}_1 \hat{p}_2 \hat{q}_2}} \quad (9)$$

For each value of combination of the parameters p_1, q_1, r, n 1000 simulations were performed and the sampling mean and sampling standard deviation of $\hat{r}, \hat{r}^2, \hat{r}^4, \hat{r}^6, \hat{r}^8$ were recorded. It was observed that sampling variance increases for increasing x in \hat{r}^x .

4.2 Simulation 2

We assume that the discovery SNP pairs are ascertained (from a very large number of tests) have high \hat{r} between observed SNPs and causal variants because otherwise power of detection would be low. We can hypothesize that the distribution of \hat{r} in this ascertained sample will be a mixture of r that is high and r that is lower but with ascertained higher values from sampling. Therefore, we would expect those with truly high r to have a higher replication rate in independent datasets, and those with ascertained high \hat{r} to have lower replication because resampling is unlikely to result in the same extreme ascertainment. To obtain empirical estimates of \hat{r} in discovery and replication datasets we conducted the following simulation.

1. Using 1000 genomes data (phase 1, version 3, 379 European samples) we selected the 528,509 “markers” used in the original discovery analysis, plus 100,000 randomly chosen “causal variants” (CVs) with minor allele frequency > 0.05 .
2. The 379 individuals were split into discovery (190) and replication (189) sets.
3. For each CV the marker with the maximum \hat{r}_D^2 from the marker panel was recorded in the discovery set. This marker was known as the “discovery marker” (DM).
4. The \hat{r}_R^2 for each CV/DM pair was then calculated in the replication set where the discovery LD was ascertained to be high, such that $\hat{r}_D^2 > 0.9$.

We observed that there was an average decrease in \hat{r}_R^x relative to \hat{r}_D^x , and that this decrease was larger with increasing x . We observed that $(\hat{r}_R^2 - \hat{r}_D^2)/\hat{r}_D^2 = 0.029$ whereas $(\hat{r}_R^8 - \hat{r}_D^8)/\hat{r}_D^8 = 0.092$. The average drop in in replication \hat{r}^8 was 3 times higher than the drop in \hat{r}^2 .

284 4.3 Interpretation

285 Simulation 1 shows that sampling variance of r^x increases as x increases. Detec-
 286 tion of epistasis is highly dependent upon high \hat{r} . Amongst the discovery SNPs
 287 there will be a mixture of interactions where observed SNPs are either in true
 288 high LD with causal variants, or will have highly inflated sample \hat{r}^x compared
 289 to the population r^x . Simulation 2 shows that as x gets larger, the average
 290 decrease in \hat{r}^x between discovery and replication becomes larger, likely to be
 291 a result of ascertained high \hat{r} in the discovery and increased sampling variance
 292 with increasing x in the replication. These results demonstrate that if all else
 293 is equal, the impact of sampling variance of r alone will reduce the replication
 294 rate of epistatic effects compared to additive effects.

295 5 Additive and non-additive variance estimation

296 5.1 Fixed effects

297 To compare the relative contribution to the phenotypic variance of gene ex-
 298 pression levels between additive and epistatic effects we are constrained by the
 299 problem that non-additive variance components for a phenotype cannot be cal-
 300 culated directly. Here, we only have SNP pairs that exceed a threshold of
 301 $p < 2.91 \times 10^{-16} = T_e$. A strong conclusion cannot be made about the genome-
 302 wide variance contribution, but we can compare the variance explained by SNP
 303 effects at this threshold for additive scans and epistatic scans.

304 In Powell *et al* 2012¹ an expression quantitative trait locus (eQTL) study
 305 was performed searching for additive effects in the same BSGS dataset as was
 306 used for the discovery here. Using the threshold T_e for the additive eQTL study,
 307 453 of the 7339 probes analysed here had at least one significant additive effect.
 308 Assuming that the phenotypic variance for each of the probes is normalised to
 309 1, the total phenotypic variance of all 7339 explained by the significant additive
 310 effects was 1.73%.

311 Following the same procedure, at the threshold T_e there were 238 gene ex-
 312 pression probes with at least one significant pairwise epistatic interaction out
 313 of the 7339 tested. In total the proportion of the phenotypic variance explained
 314 by the epistatic effects at these SNP pairs was 0.25%.

315 5.2 Pedigree estimates

316 The gene expression levels for MBNL1, TMEM149, NAPRT1, TRAPPC5 and
 317 CAST are influenced by large *cis-trans* epistatic networks (eight interactions
 318 or more). Though it is not possible to orthogonally estimate the non-additive
 319 genetic variance for non-clonal populations, an approximation of a component of
 320 non-additive variance can be estimated using pedigree information. The BSGS
 321 data is comprised of some related individuals and standard quantitative genetic
 322 analysis was used to calculate the additive and dominance variance components
 323 for each gene expression phenotype in Powell *et al* 2013.⁶ The dominance effect

324 is likely to capture additive \times additive genetic variance plus some fraction of
 325 other epistatic variance components. We found that the aforementioned genes
 326 had dominance variance component estimates within the top 5% of all 17,994
 327 gene expression probes that were analysed in Powell *et al* 2013.

328 6 Functional enrichment analysis

329 6.1 Tissue specific transcriptionally active regions

330 We employed a recently published method ([http://www.broadinstitute.org/
 331 mpg/epigwas/](http://www.broadinstitute.org/mpg/epigwas/))¹⁴ that tests for cell-type-specific enrichment of active chromatin,
 332 measured through H3K4me3 chromatin marks¹⁵ in regions surrounding the 731
 333 SNPs that comprise the 501 discovery interactions. The exact method used to
 334 perform this analysis has been described previously.¹⁶ Briefly, we tested the
 335 hypothesis that the 731 SNPs were more likely to be in transcriptionally active
 336 regions (as measured by chromatin marks) than a random set of SNPs selected
 337 from the same SNP chip. This hypothesis was tested for 34 cell types across
 338 four broad tissue types (haematopoietic, gastrointestinal, musculoskeletal and
 339 endocrine, and brain).

340 6.2 Chromosome interactions

341 It has been shown¹⁷ that different regions on different chromosomes or within
 342 chromosomes spatially colocalise within the cell. We shall refer to the colocalisa-
 343 tion of two chromosome regions as a chromosome interaction. A map of pairwise
 344 chromosome interactions for K562 blood cell lines was recently produced,¹⁸ and
 345 we hypothesised that part of the underlying biological mechanism behind some
 346 of the 501 epistatic interactions may arise from chromosome interactions. We
 347 found that 44 of the putative epistatic interactions were amongst SNPs that
 348 were within 5Mb of known chromosome interactions. This means that SNP A
 349 was no more than 2.5Mb from the focal point of the chromosome interaction on
 350 chromosome A, and SNP B was no more than 2.5Mb from the focal point on
 351 chromosome B.

352 We performed simulations to test how extreme the observation of 44 epistatic
 353 interactions overlapping with chromosome interactions is compared to chance.
 354 Chromosome interactions fall within functional genomic regions,^{17,18} and the
 355 SNPs in our epistatic interactions are enriched for functional genomic regions.
 356 Therefore, we designed the simulations to ensure that the null distribution was
 357 of chromosome interactions between SNPs enriched for functional genomic re-
 358 gions but with no known epistatic interactions. To do this we used the 731 SNPs
 359 that form the 501 putative epistatic interactions and randomly shuffled them to
 360 create new sets of 501 pairs, disallowing any SNP combinations that were in the
 361 original set. Therefore, each new random set was enriched for functional regions
 362 but had no genetic interactions. We scanned the map of chromosome interac-
 363 tions for overlaps with the new sets and then repeated the random shuffling

process. We performed 1,000 such permutations to generate a null distribution of chromosome interaction overlaps.

We repeated this process, searching for overlaps within 1Mb, 250kb, and 10kb.

6.3 SNP colocalisation with genomic features

We tested for enrichment of genomic features for the 687 IndexSNPs that comprise the 434 epistatic interactions with data present in discovery and replication datasets. For each of the 687 IndexSNPs we calculated LD with all regional SNPs within a radius of 0.5Mb and kept all regional SNPs with LD $r^2 > 0.8$. We then cross-referenced the remaining regional SNPs with the annotated chromatin structure reference¹⁹) querying whether the regional SNPs fell in Predicted promoter region including TSS (TSS), Predicted promoter flanking region (PF), Predicted enhancer (E), Predicted weak enhancer or open chromatin cis regulatory element (WE), CTCF enriched element (CTCF), Predicted transcribed region (T), or Predicted Repressed or Low Activity region (R) positions. Therefore a particular IndexSNP might cover multiple genomic features through LD.

We then performed the whole querying process for each of the 528,509 SNPs present in the SNP chip used in the scan, and used the results from this second analysis to establish a null distribution for the expected proportion of SNPs for each genomic feature. We calculated p -values for enrichment of each of the seven genomic features independently, and for *cis*- and *trans*-SNPs separately, using a binomial test. For each genomic feature we used the expected proportion of SNPs as the expected probability of “success” (p). Here, a success is defined as an IndexSNP residing in a region that includes the genomic feature. The observed number of successes for each IndexSNP (k) out of the total count of IndexSNPs (n) was then modelled as $\Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$.

6.4 Transcription factor enrichment

To test for enrichment of transcription factor binding sites (TFBS) we followed a procedure similar to that described in Section 6.3. For each of the 687 IndexSNPs we extracted regional SNPs as previously described. We then used the PWMEnrich package in Bioconductor (<http://www.bioconductor.org/packages/2.12/bioc/html/PWMEnrich.html>) to identify which TFBSs each of the regional SNPs for one IndexSNP falls in (within a radius of 250bp). Thus, the number of occurrences of a particular TFBS was counted for each IndexSNP. We used the “Threshold-free affinity” method for identifying TFBSs.²⁰

We constructed a null distribution of expected TFBS occurrences based on the same null hypothesis as described in Section 6.3 - the probability of an IndexSNP covering a particular TFBS is identical to any of the 528,509 SNPs in the discovery SNP chip. To do this, we performed the same procedure for each SNP in the discovery SNP chip as was performed for each IndexSNP to obtain an expected probability of covering a particular TFBS. We then tested the

IndexSNPs for enrichment of each TFBS independently, and for *cis*- and *trans*-SNPs separately. *p*-values were obtained using Z-scores, calculated by using a normal approximation to the sum of binomial random variables representing motif hits along the sequence.²¹

6.5 Defining previously identified SNP associations

The discovery dataset (BSGS) had previously been analysed for additive and dominant marginal effects for all gene expression levels.^{1,6} To define SNPs that had been previously detected to have effects for a particular gene expression level we used a significance threshold accounting for multiple testing across SNPs and expression probes, $T_m = 0.05/(528509 \times 7339) = 1.29 \times 10^{-11}$. From this, we found that only nine of the 501 discovery interactions had known main effects, 64 were between SNPs that had no known marginal effects, and 439 were between a SNP with a known marginal effect and a SNP with no known marginal effect.

7 References

References

- ¹ Powell, J. E. *et al.* The Brisbane Systems Genetics Study: genetical genomics meets complex trait genetics. *PLoS one* **7**, e35430 (2012).
- ² Medland, S. E. *et al.* Common variants in the trichohyalin gene are associated with straight hair in Europeans. *American journal of human genetics* **85**, 750–5 (2009).
- ³ Aulchenko, Y. S., Ripke, S., Isaacs, A. & van Duijn, C. M. GenABEL: an R library for genome-wide association analysis. *Bioinformatics (Oxford, England)* **23**, 1294–6 (2007).
- ⁴ Hemani, G., Theocharidis, A., Wei, W. & Haley, C. EpiGPU: exhaustive pairwise epistasis scans parallelized on consumer level graphics cards. *Bioinformatics (Oxford, England)* **27**, 1462–5 (2011).
- ⁵ Hemani, G., Knott, S. & Haley, C. An Evolutionary Perspective on Epistasis and the Missing Heritability. *PLoS Genetics* **9**, e1003295 (2013).
- ⁶ Powell, J. E. *et al.* Congruence of Additive and Non-Additive Effects on Gene Expression Estimated from Pedigree and SNP Data. *PLoS Genetics* **9**, e1003502 (2013).
- ⁷ Yang, J. *et al.* Genome partitioning of genetic variation for complex traits using common SNPs. *Nature Genetics* **43**, 519–525 (2011).

- 440 ⁸ Fehrmann, R. S. N. *et al.* Trans-eQTLs reveal that independent genetic vari-
 441 ants associated with a complex phenotype converge on intermediate genes,
 442 with a major role for the HLA. *PLoS genetics* **7**, e1002197 (2011).
- 443 ⁹ Metspalu, A. The Estonian Genome Project. *Drug Development Research*
 444 **62**, 97–101 (2004).
- 445 ¹⁰ Preininger, M. *et al.* Blood-informative transcripts define nine common axes
 446 of peripheral blood gene expression. *PLoS genetics* **9**, e1003362 (2013).
- 447 ¹¹ Westra, H.-J. *et al.* MixupMapper: correcting sample mix-ups in genome-
 448 wide datasets increases power to detect small genetic effects. *Bioinformatics*
 449 (*Oxford, England*) **27**, 2104–11 (2011).
- 450 ¹² Williams, D. A. Improved likelihood ratio tests for complete contingency
 451 tables. *Biometrika* **63**, 33–37 (1976).
- 452 ¹³ Alvarez-Castro, J., Le Rouzic, A., Carlborg, O., Álvarez Castro, J. M. &
 453 Carlborg, O. How to perform meaningful estimates of genetic effects. *PLoS*
 454 *Genetics* **4**, e1000062 (2008).
- 455 ¹⁴ Trynka, G. *et al.* Chromatin marks identify critical cell types for fine mapping
 456 complex trait variants. *Nature genetics* **45**, 124–30 (2013).
- 457 ¹⁵ Koch, C. M. *et al.* The landscape of histone modifications across 1% of the
 458 human genome in five human cell lines. *Genome research* **17**, 691–707 (2007).
- 459 ¹⁶ Rietveld, C. A. *et al.* GWAS of 126,559 Individuals Identifies Genetic Variants
 460 Associated with Educational Attainment. *Science* (2013).
- 461 ¹⁷ Lieberman-Aiden, E. *et al.* Comprehensive mapping of long-range interactions
 462 reveals folding principles of the human genome. *Science (New York, N.Y.)*
 463 **326**, 289–93 (2009).
- 464 ¹⁸ Lan, X. *et al.* Integration of Hi-C and ChIP-seq data reveals distinct types
 465 of chromatin linkages. *Nucleic acids research* **40**, 7690–704 (2012).
- 466 ¹⁹ Hoffman, M., Buske, O., Wang, J. & Weng, Z. Unsupervised pattern dis-
 467 covery in human chromatin structure through genomic segmentation. *Nature*
 468 *Methods* **9**, 473–476 (2012).
- 469 ²⁰ Stormo, G. DNA binding sites: representation and discovery. *Bioinformatics*
 470 **16**, 16–23 (2000).
- 471 ²¹ Ho Sui, S. J. *et al.* oPOSSUM: identification of over-represented transcription
 472 factor binding sites in co-expressed genes. *Nucleic acids research* **33**, 3154–64
 473 (2005).