# Detection and replication of epistasis influencing transcription in humans

Gibran Hemani[1,2,*], Konstantin Shakhbazov[1,2], Harm-Jan Westra[3], Tonu Esko[4,5,6], Anjali K Henders[7], Allan F McRae[1,2], Jian Yang[2], Greg Gibson[8], Nicholas G Martin[7], Andres Metspalu[4], Lude Franke[3], Grant W Montgomery[7,+], Peter M Visscher[1,2,+], and Joseph E Powell[1,2,+]

[1]University of Queensland Diamantina Institute, University of Queensland, Princess Alexandra Hospital, Brisbane, Queensland, Australia. [2]Queensland Brain Institute, University of Queensland, Brisbane, QLD, Australia. [3]Department of Genetics, University Medical Center Groningen, University of Groningen, Hanzeplein 1, Groningen, the Netherlands. [4]Estonian Genome Center, University of Tartu, Tartu, 51010, Estonia. [5]Medical and Population Genetics, Broad Institute, Cambridge, MA, 02142, US. [6]Divisions of Endocrinology, Children's Hospital, Boston, MA, 02115, US. [7]Queensland Institute of Medical Research, Brisbane, Queensland, Australia. [8]School of Biology and Centre for Integrative Genomics, Georgia Institute of Technology, Atlanta, Georgia United States of America. [+]These authors contributed equally. [*]Corresponding author: g.hemani@uq.edu.au

1

**Abstract**

Epistasis is the phenomenon whereby a polymorphism's effect on a trait depends on other polymorphisms present in the genome. The extent to which epistasis influences complex traits[1] and contributes to their variation[2,3] is a fundamental question in evolution and human genetics. Though epistasis has been demonstrated in artificial gene manipulation studies in model organisms,[4,5] and examples have been reported in other species,[6] few convincing examples with independent replication exist for epistasis amongst natural polymorphisms in human traits.[7,8] Its absence from empirical findings may simply be due to its low incidence in the genetic control of complex traits,[2,3] but an alternative view is that it has previously been too technically challenging to detect due to statistical power and computational issues.[9] Here we show that, using advanced computation techniques[10] and a gene expression study design, many instances of epistasis are found between common single nucleotide polymorphisms (SNPs). In a cohort of 846 individuals with data on 7339 gene expression levels in peripheral blood, we found 501 significant pairwise epistatic interactions between common SNPs acting on the expression levels of 238 genes ($p < 2.91 \times 10^{-16}$). Replication of these interactions in two independent data sets[11,12] showed both concordance of direction of epistatic effects ($p = 5.56 \times 10^{-31}$) and enrichment of interaction $p$-values, with 30 being significant at a conservative threshold of $p < 0.05/501$. There was evidence of functional enrichment for the interacting SNPs, for instance 44 of the genetic interactions are located within 5Mb of regions of known physical chromosome interactions[13] ($p = 1.8 \times 10^{-10}$). Epistatic networks of three SNPs or more influence the expression levels of 129 genes, whereby one *cis*-acting SNP is modulated by several *trans*-acting SNPs. For example MBNL1 is influenced by an additive effect at rs13069559 which itself is masked by *trans*-SNPs on 14 different chromosomes, with nearly identical genotype-phenotype (GP) maps for each *cis-trans* interaction. This study presents the first evidence for multiple instances of segregating common polymorphisms interacting to influence human traits.

# Main text

In the genetic analysis of complex traits it is usual for SNP effects to be estimated using an additive model where they are assumed to contribute independently and cumulatively to the mean of a trait. This framework has been successful in identifying thousands of associations.[14] But to date, though its contribution to phenotypic variance is frequently the subject of debate,[1–3] there is little empirical exploration of the role that epistasis plays in the architecture of complex traits in humans.[7,8] Beyond the prism of human association studies there is evidence for epistasis, not only at the molecular scale from artificially induced mutations[4] but also at the evolutionary scale in fitness adaptation[15] and speciation.[16]

Methods are now available to overcome the computational problems involved in searching for epistasis, but its detection still remains problematic due to re-

duced statistical power. For example increased dependence on linkage disequilibrium (LD) between causal SNPs and observed SNPs,[17, 18] increased model complexity in fitting interaction terms,[19] and more extreme significance thresholds to account for increased multiple testing[9] all make it more difficult to detect epistasis in comparison to additive effects. Thus, when combined with small genetic effect sizes, as is expected in most complex traits of interest,[14] the power to detect epistasis diminishes rapidly. There are two simple ways to overcome this problem. One is by using extremely large sample sizes;[20] another is by analysing traits that are likely to have large effect sizes among common variants. Because our focus was to ascertain the extent to which instances of epistasis arises from natural genetic variation we designed a study around the latter approach and searched for epistatic genetic effects that influence gene expression levels. Transcription levels can be measured for thousands of genes and like most complex diseases, these expression traits are typically heritable.[21] But unlike complex diseases, genetic associations with gene expression commonly have very large effect sizes that explain large proportions of the genetic variance,[22] making them good candidates to search for epistasis, should it exist.

In our discovery dataset (Brisbane Systems Genetics Study, BSGS[23]) of 846 individuals genotyped at 528,509 SNPs, we used a two stage approach to identify genetic interactions. First, we exhaustively test every pair of SNPs for pairwise effects against each of 7339 expression traits in peripheral blood (family-wise error rate of 5% corresponding to a significance threshold of $p < 2.91 \times 10^{-16}$, Methods). Second, we filtered the SNP pairs from stage 1 on LD and genotype class counts, and tested the remaining pairwise effects for significant interaction terms and used a Bonferroni correction for multiple testing (estimated type 1 error rate $0.05 \le \alpha \le 0.14$, Methods, Supplementary Figure S1). Using this design we identified 501 putative genetic interactions influencing the expression levels of 238 genes (Supplementary Table S1). We used strict quality control measures to avoid statistical associations being driven by technical artifacts (Methods). However it remains possible that unexplained technical artifacts may have led to the significant discovery interactions. Of the 501 discovery interactions, 434 had available data and passed filtering (Methods) in two independent replication datasets, Fehrmann[12] and the Estonian Genomics Centre University of Tartu (EGCUT),[11] in which we saw convincing evidence for replication. We used the summary statistics from the replication datasets to perform a meta analysis to obtain an independent $p$-value for the putative interactions, and 30 were significant after applying a Bonferroni correction for multiple testing (5% significance threshold $p < 0.05/501$, Table 1). To quantify the similarity of GP maps between the independent datasets (Figure 1) we decomposed the genetic effects of each of the SNP pairs into orthogonal additive, dominance and epistatic effects ($A1$, $A2$, $D1$, $D2$, $A \times A$, $A \times D$, $D \times A$, $D \times D$) and tested for concordance of the sign of the most signicant effect (Supplementary Table S3, Methods). Sign concordance between the discovery and both replication datasets was observed in 22 out of the 30 significantly replicated interactions (expected value = 7.5 under the null hypothesis of no interactions, $p = 3.76 \times 10^{-8}$).

In addition, using the meta analysis from the replication samples only, we

3

observed that 316 of the remaining 404 discovery SNP pairs had replication interaction $p$-values more extreme than the 2.5% confidence interval of the quantile-quantile plot against the null hypothesis of no interactions where $p$-values are assumed to be uniformly distributed ($p << 1.0 \times 10^{-16}$, Figure 2 and Supplementary Figure S2). Concordance of the direction of the effect of the largest variance component was also highly significant ($p = 5.71 \times 10^{-31}$, Supplementary Table S3). The congruence of the epistatic networks in discovery and replication datasets is shown in Figure 3, demonstrating that these complex genetic patterns are common even across independent datasets. A further replication was attempted using the Centre for Health Discovery and Wellbeing (CHDWB) dataset,[24] but only 20 of the SNP pairs passed filtering because the sample size was small ($n = 139$), and likely due to insufficient power we found no evidence for replication (Supplementary Figure S6).

It should be noted that although it is a necessary step to establish the veracity of the interactions from the discovery set, replication of epistasis is difficult in practice. For example, LD between causal variants and observed markers plays an important role. Not only is the dependence on LD much greater for epistatic effects than for additive effects (Supplementary Figure S7), but when estimating epistatic variance it is more sensitive to changes in LD between observed SNPs and causal variants between independent samples when compared to additive effects (Supplementary Figure S8). This has a direct effect on statistical power for replication. The sampling variance of LD $r$ leads to the ascertainment of marker associations with higher sample $r$ in the discovery stage in comparison to the replication stage. However, the average decrease in $\hat{r}^x$ in replication samples becomes larger as $x$ increases (Methods, Supplementary Figure S9). For example, the decrease in $\hat{r}^8$ (which is proportional to the power of deteting $D \times D$ effects), is on average three fold greater than the decrease in $\hat{r}^2$ (which is proportional to the power of detecting additive effects).

Though seldom the focus of association studies, SNPs with known main effects are often tested for additive $\times$ additive genetic interactions,[9] but our analysis shows that this is unlikely to be the most effective strategy for its detection. The majority of our discovery interactions comprised of one SNP that was significantly associated with the gene expression level in the discovery dataset, and one SNP that had no previous association[22] (439 out of 501, Methods). Only nine interactions were between SNPs that both had known main effects while 64 were between SNPs that had no known main effects. Additionally, we observed that the largest epistatic variance component for the 501 interactions was equally divided amongst additive $\times$ additive, additive $\times$ dominance, dominance $\times$ additive and dominance $\times$ dominance at the discovery stage ($p = 0.22$ for departure from expectation). This is not surprising because the patterns of epistasis used for statistical decomposition (*i.e.* $A \times A$, $A \times D$, $D \times A$, $D \times D$) are simply convenient orthogonal parameterisations of a two locus model, and are not intended to model biological function.[25]

Of the discovery interactions, 26 were *cis-cis* acting (within 1Mb of the transcription start site, mean distance between SNPs was 0.53Mb), 462 were *cis-trans*-acting, and 13 were *trans-trans*-acting. We observed a wide range of

significant GP maps (Figure 1) but the most common pattern of epistasis that we detected involved a *trans*-SNP masking the effect of an additive *cis*-SNP. For example, MBNL1 (involved in RNA modification and regulation of splicing[26]) has a *cis* effect at rs13069559 which in turn is controlled by 13 *trans*-SNPs and one *cis*-SNP that each exhibit a masking pattern, such that when the *trans*-SNP is homozygous for the masking allele the decreasing allele of the *cis*-SNP no longer has an effect (Supplementary Figure S10). Each of these interactions has evidence for replication in at least one dataset and six are significantly replicated at the Bonferroni level (Supplementary Figure S3). We see similar epistatic networks involving multiple (eight or more) *trans*-acting SNPs for other gene expresson levels too, for example TMEM149 (Supplementary Figure S11), NAPRT1 (Supplementary Figure S12), TRAPPC5 (Supplementary Figure S13), and CAST (Supplementary Figure S14). We observed that from pedigree analysis these five gene expression phenotypes had non-additive variance component estimates within the 95th percentile of the 17,994 gene expression phenotypes that were analysed previously[22] (Supplementary Table S2, Methods).

In total the 501 interactions comprised 781 unique SNPs, which we analysed for functional enrichment (Methods). We tested the SNPs for cell-type specific overlap with transcriptionally active chromatin regions, tagged by histone-3-lysine-4,tri-methylation (H3K4me3) chromatin marks, in 34 cell types[27] (Supplementary Figure S5). There was significant enrichment for *cis*-acting SNPs in haematopoietic cell types only ($p < 1 \times 10^{-4}$ for the three tissues with the strongest enrichment after adjusting for multiple testing). However *trans*-acting SNPs did not show any tissue specific enrichment ($p > 0.1$ for all tissues). This difference between *cis* and *trans* SNPs suggests different roles in epistatic interactions where tissue specificity is provided by the *cis* SNPs. There is also enrichment for *cis*-SNPs to be localised in regions with regulatory genomic features as measured by chromatin states[28] (Supplementary Figure S4).

We also demonstrate physical organisation of interacting loci within the cell, suggesting a mechanism by which biological function can lead to epistatic genetic variance. It has been shown that different chromosomal regions spatially colocalise in the cell through chromatin interactions.[13] We cross-referenced our epistatic SNPs with a map of chromosome interacting regions ($n = 96,139$) in K562 blood cell lines[29] (Methods) and found that 44 epistatic interactions mapped to within 5Mb ($p < 1.8 \times 10^{-10}$), (Supplementary Figure S15). Interaction of distant loci may occur through physical proximity in transcriptional factories that organise across different chromosome regions and can regulate transcription of related genes.[30, 31]

Though we present many instances of epistasis, quantifying its relative importance to complex traits in humans remains an open question. In this study we are able to identify 238 gene expression traits with at least one significant interaction given our experiment-wide threshold, where the minimum estimated variance explained by the epistatic effects of any interaction was 2.1% of phenotypic variance. Taking results from our previously published eQTL[23] we calculated that 1848 of the 7339 gene expression levels analysed were influenced by additive effects where the estimated additive variance of a locus was 2.1% or

5

greater. Thus, we can infer that the number of instances of large additive effects is significantly greater than the number of instances of large epistatic effects.

In terms of their contribution to complex traits a more important metric might be the proportion of the variance that the epistatic loci explain.[2] Ideally one would approach this question from a whole genome perspective[32] but this is intractable for non-additive variance components. Nevertheless, some inference can be made from the ascertained effects in these analyses and it is evident that estimated additive variance is overall a larger component than estimated epistatic variance, as has been argued previously.[2,3] Taking all additive effects detected in Powell *et al* (2012) that have additive variance explaining 2.1% or greater of phenotypic variance, we calculated that the proportion of total phenotypic variance of all 7339 gene expression levels explained by additive effects alone was 2.16%. By contrast, the estimated epistatic variance from the interacting SNPs detected in this study on average explain a total of 0.22% of phenotypic variance, approximately ten times lower than the estimated additive variance. There are several caveats to this comparison. Firstly, the ratio of additive to epistatic variance may differ at different minimum variance thresholds, and our estimate is determined by the threshold used. Secondly, the power of a 1 *d.f.* test exceeds that of an 8 *d.f.* test. Thirdly, the non-additive variance at causal variants is expected to be underestimated by observed SNPs in comparison to estimates for additive variance. This is due to differences in the rate of decay of the estimate of the genetic variance of the causal SNPs as LD decreases with the observed SNPs. And forthly, the extent of winner's curse in estimation of effect sizes may differ between the the two studies.

Overall, we have demonstrated that it is possible to identify and replicate epistasis in complex traits amongst common human variants, despite the relative contribution of pairwise epistasis to phenotypic variation being small. The bioinformatic analysis of the significant epistatic loci suggests that there are a large number of possible mechanisms that can lead to non-additive genetic variation. Further research into such epistatic effects may provide a useful framework for understanding molecular mechanisms and complex trait variation in greater detail. With computational techniques and data now widely available the search for epistasis in larger datasets for traits of broader interest is warranted.

## Methods Summary

We searched for pairwise epistasis exhaustively in the BSGS discovery dataset,[23] which comprises 846 individuals who are genotyped at 528,509 autosomal SNPs. Each individual had gene expression levels measured in peripheral blood at 47,323 probes. Only the probes that passed quality control and had significant expression in $\geq$ 90% of individuals were used in the analysis (7,339 probes representing 6,158 RefSeq genes). Recent hardware and software[10] advances that use graphics processing units (GPUs) made it possible to perform the $1.03 \times 10^{15}$ statistical tests to complete this analysis. We used permutation analysis[33] to calculate an experiment-wide significance threshold of $T_e = 2.91 \times 10^{-16}$ at the 5% family-wise error rate (FWER). SNP pairs were modelled for

full genetic effects, including marginal additive and dominance at both SNPs plus four interaction terms. Though we could have used a less complex model to improve statistical efficiency, we deemed it important to be agnostic about the type of epistasis that might exist, and therefore chose not to over-parameterise the test.[18, 19] Because there are many large marginal effects present in these data it was necessary to perform several filtering steps to exclude SNP pairs that were significant due to marginal effects alone. All SNP pairs with LD $r^2 > 0.1$ and $D'^2 > 0.1$ were removed to minimise the possibility of haplotype effects. All SNP pairs were required to have at least five data points in all nine genotype classes. If multiple SNP pairs were present on the same chromosomes for a particular expression trait then only the sentinel SNP pair was retained. Finally, a nested test contrasting the full genetic model against the marginal additive and dominance model was performed for each remaining SNP pair (Methods), resulting in 501 significant interactions after Bonferroni correction for multiple testing of the filtered SNPs. The 501 significant SNP pairs were carried forward for replication in two independent datasets that used the same expression assays for analysing transcription in peripheral blood, the Fehrmann dataset[12] ($n = 1240$) and the Estonian Genome Centre University of the University of Tartu (EGCUT) dataset[11] ($n = 891$). Of these, 434 passed filtering in both replication datasets. A meta analysis on the interaction $p$-values from each replication dataset was performed to provide an overall replication statistic for each putative interaction.

## Acknowledgements

7

# Tables

Table 1: Epistatic interactions significant at the Bonferroni level in two replication sets

|    | Gene (chr.) | SNP 1 (chr.) | SNP 2 (chr.) | BSGS[2] | Fehrmann[3] | EGCUT[3] | Meta[4] |
|----|-------------|--------------|--------------|---------|-------------|----------|---------|
| 1  | ADK (10) | rs2395095 (10) | rs10824092 (10) | 6.69[1] | 18.33[1] | 21.21[1] | 39.82[1] |
| 2  | ATP13A1 (19) | rs4284750 (19) | rs873870 (19) | 5.30 | 12.18 | 3.25 | 14.23 |
| 3  | C21ORF57 (21) | rs9978658 (21) | rs11701361 (21) | 9.42 | 6.08 | 16.36 | 21.67 |
| 4  | CSTB (21) | rs9979356 (21) | rs3761385 (21) | 11.99 | 25.20 | 16.72 | 42.27 |
| 5  | CTSC (11) | rs7930237 (11) | rs556895 (11) | 7.16 | 18.76 | 15.06 | 33.53 |
| 6  | FN3KRP (17) | rs898095 (17) | rs9892064 (17) | 16.16 | 28.24 | 29.39 | 59.95 |
| 7  | GAA (17) | rs11150847 (17) | rs12602462 (17) | 13.91 | 19.98 | 12.99 | 32.60 |
| 8  | HNRPH1 (5) | rs6894268 (5) | rs4700810 (5) | 15.38 | 8.55 | 3.01 | 10.37 |
| 9  | LAX1 (1) | rs1891432 (1) | rs10900520 (1) | 19.16 | 18.60 | 11.22 | 29.24 |
| 10 | MBNL1 (3) | rs16864367 (3) | rs13079208 (3) | 13.49 | 16.25 | 24.74 | 41.56 |
| 11 | MBNL1 (3) | rs7710738 (5) | rs13069559 (3) | 7.92 | 2.55 | 7.89 | 9.28 |
| 12 | MBNL1 (3) | rs2030926 (6) | rs13069559 (3) | 7.10 | 0.91 | 5.80 | 5.53 |
| 13 | MBNL1 (3) | rs2614467 (14) | rs13069559 (3) | 5.74 | 4.13 | 2.22 | 5.30 |
| 14 | MBNL1 (3) | rs218671 (17) | rs13069559 (3) | 7.63 | 0.62 | 5.82 | 5.23 |
| 15 | MBNL1 (3) | rs11981513 (7) | rs13069559 (3) | 7.71 | 0.43 | 5.36 | 4.58 |
| 16 | MBP (18) | rs8092433 (18) | rs4890876 (18) | 5.40 | 7.06 | 21.91 | 28.73 |
| 17 | NAPRT1 (8) | rs2123758 (8) | rs3889129 (8) | 8.45 | 15.12 | 16.08 | 30.77 |
| 18 | NCL (2) | rs7563453 (2) | rs4973397 (2) | 7.31 | 7.51 | 6.33 | 12.70 |
| 19 | PRMT2 (21) | rs2839372 (21) | rs11701058 (21) | 4.81 | 0.69 | 4.47 | 4.06 |
| 20 | RPL13 (16) | rs352935 (16) | rs2965817 (16) | 4.98 | 3.79 | 14.41 | 17.24 |
| 21 | SNORD14A (11) | rs2634462 (11) | rs6486334 (11) | 7.31 | 13.11 | 10.96 | 23.22 |
| 22 | TMEM149 (19) | rs807491 (19) | rs7254601 (19) | 12.16 | 81.55 | 45.78 | 145.78 |
| 23 | TMEM149 (19) | rs8106959 (19) | rs6926382 (6) | 5.80 | 3.06 | 8.80 | 10.72 |
| 24 | TMEM149 (19) | rs8106959 (19) | rs914940 (1) | 6.22 | 3.36 | 6.96 | 9.20 |
| 25 | TMEM149 (19) | rs8106959 (19) | rs2351458 (4) | 7.30 | 0.04 | 9.61 | 8.00 |
| 26 | TMEM149 (19) | rs8106959 (19) | rs6718480 (2) | 8.55 | 3.31 | 5.15 | 7.36 |
| 27 | TMEM149 (19) | rs8106959 (19) | rs1843357 (8) | 6.21 | 3.72 | 3.33 | 6.00 |
| 28 | TMEM149 (19) | rs8106959 (19) | rs9509428 (13) | 9.44 | 0.10 | 5.75 | 4.47 |
| 29 | TRA2A (7) | rs7776572 (7) | rs11770192 (7) | 8.23 | 3.19 | 1.89 | 4.09 |
| 30 | VASP (19) | rs1264226 (19) | rs2276470 (19) | 5.09 | 0.94 | 5.14 | 4.95 |

[1] $-\log_{10} p$-values for 4 *d.f.* interaction tests

[2] Discovery dataset

[3] Independent replication dataset

[4] Meta analysis of interaction terms between replication datasets only

# Figures

Figure 1: **Replication of GP maps in two independent populations**
The GP maps for each epistatic interaction that is significant at the Bonferroni
level in both replication datasets are shown. Each GP map consists of nine
tiles where each tile represents the expression level for that two-locus genotype
class. Phenotypes are for gene transcript levels (dark coloured tiles = high
expression, light coloured tiles = low expression). Columns of GP maps are for
each independent dataset. Rows of GP maps are for each of 30 significantly
replicated interactions at the Bonferroni level, corresponding to the rows in
Table 1. There is a clear trend of the GP maps replicating across all three
datasets.

Figure 2: **Q-Q plots of interaction $p$-values from replication datasets**
The top panel shows all 434 discovery SNPs that were tested for interactions.
Observed $p$-values ($y$-axis, $-\log_{10}$ scale) are plotted against the expected $p$-
values ($x$-axis, $-\log_{10}$ scale). The multiple testing correction threshold for
significance following Bonferroni correction is denoted by a dotted line. The
bottom panel shows the same data as the top panel but excluding the 30 inter-
actions that were significant at the Bonferroni level in the replication datasets.
The shaded grey area represents the 5% confidence interval for the expected
distribution of $p$-values. Dark blue points represent $p$-values that exceed the
confidence interval, light blue are within the confidence interval.

Figure 3: **Discovery and replication of epistatic networks** All 434 putative
genetic interactions (edges) with data common to discovery and replication sets
is shown, where black nodes represent SNPs and red nodes represent traits
(gene expression probes). Three hundred and forty-five interactions had $p$-values
exceeding the 2.5% confidence interval following meta analysis of the replication
data The remaining 89 interactions that did not replicate are depicted in grey.
It is evident that a large proportion of the complex networks identified in the
discovery set also exist in independent populations. An interactive version of
this graph can be found here: `http://kn3in.github.io/detecting_epi/`

# References

[1] Carlborg, O. & Haley, C. S. Epistasis: too often neglected in complex trait studies? *Nature Reviews Genetics* **5**, 618–25 (2004).

[2] Hill, W. G., Goddard, M. E. & Visscher, P. M. Data and Theory Point to Mainly Additive Genetic Variance for Complex Traits. *PLoS Genetics* **4** (2008).

[3] Crow, J. F. On epistasis: why it is unimportant in polygenic directional selection. *Philosophical transactions of the Royal Society of London. Series B, Biological sciences* **365**, 1241–4 (2010).

[4] Costanzo, M. *et al.* The genetic landscape of a cell. *Science (New York, N.Y.)* **327**, 425–31 (2010).

[5] Bloom, J. S., Ehrenreich, I. M., Loo, W. T., Lite, T.-L. V. o. & Kruglyak, L. Finding the sources of missing heritability in a yeast cross. *Nature* 1–6 (2013).

[6] Carlborg, O., Jacobsson, L., Ahgren, P., Siegel, P. & Andersson, L. Epistasis and the release of genetic variation during long-term selection. *Nature Genetics* **38**, 418–420 (2006).

[7] Strange, A. *et al.* A genome-wide association study identifies new psoriasis susceptibility loci and an interaction between HLA-C and ERAP1. *Nature Genetics* **42**, 985–90 (2010).

[8] Evans, D. M. *et al.* Interaction between ERAP1 and HLA-B27 in ankylosing spondylitis implicates peptide handling in the mechanism for HLA-B27 in disease susceptibility. *Nature Genetics* **43** (2011).

[9] Cordell, H. J. Detecting gene-gene interactions that underlie human diseases. *Nature Reviews Genetics* **10**, 392–404 (2009).

[10] 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).

[11] Metspalu, A. The Estonian Genome Project. *Drug Development Research* **62**, 97–101 (2004).

[12] Fehrmann, R. S. N. *et al.* Trans-eQTLs reveal that independent genetic variants associated with a complex phenotype converge on intermediate genes, with a major role for the HLA. *PLoS genetics* **7**, e1002197 (2011).

[13] Lieberman-Aiden, E. *et al.* Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science (New York, N.Y.)* **326**, 289–93 (2009).

[14] Visscher, P. M., Brown, M. a., McCarthy, M. I. & Yang, J. Five years of GWAS discovery. *American journal of human genetics* **90**, 7–24 (2012).

[15] Weinreich, D. M., Delaney, N. F., Depristo, M. a. & Hartl, D. L. Darwinian evolution can follow only very few mutational paths to fitter proteins. *Science (New York, N.Y.)* **312**, 111–4 (2006).

[16] Breen, M. S., Kemena, C., Vlasov, P. K., Notredame, C. & Kondrashov, F. a. Epistasis as the primary factor in molecular evolution. *Nature* **490**, 535–538 (2012).

[17] Weir, B. S. Linkage disequilibrium and association mapping. *Annual review of genomics and human genetics* **9**, 129–42 (2008).

[18] Hemani, G., Knott, S. & Haley, C. An Evolutionary Perspective on Epistasis and the Missing Heritability. *PLoS Genetics* **9**, e1003295 (2013).

[19] Marchini, J., Donnelly, P. & Cardon, L. R. Genome-wide strategies for detecting multiple loci that influence complex diseases. *Nature Genetics* **37**, 413–417 (2005).

[20] Lango Allen, H. *et al.* Hundreds of variants clustered in genomic loci and biological pathways affect human height. *Nature* **467**, 832–8 (2010).

[21] Schadt, E. *et al.* Genetics of gene expression surveyed in maize, mouse and man. *Nature* **422**, 297–302 (2003).

[22] 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).

[23] Powell, J. E. *et al.* The Brisbane Systems Genetics Study: genetical genomics meets complex trait genetics. *PloS one* **7**, e35430 (2012).

[24] Preininger, M. *et al.* Blood-informative transcripts define nine common axes of peripheral blood gene expression. *PLoS genetics* **9**, e1003362 (2013).

[25] Cockerham, C. C. An extension of the concept of partitioning hereditary variance for analysis of covariances among relatives when epistasis is present. *Genetics* **39**, 859–882 (1954).

[26] Ho, T. H. *et al.* Muscleblind proteins regulate alternative splicing. *The EMBO journal* **23**, 3103–12 (2004).

[27] Trynka, G. *et al.* Chromatin marks identify critical cell types for fine mapping complex trait variants. *Nature genetics* **45**, 124–30 (2013).

[28] Hoffman, M., Buske, O., Wang, J. & Weng, Z. Unsupervised pattern discovery in human chromatin structure through genomic segmentation. *Nature Methods* **9**, 473–476 (2012).

[386] [29] Lan, X. *et al.* Integration of Hi-C and ChIP-seq data reveals distinct types of chromatin linkages. *Nucleic acids research* **40**, 7690–704 (2012).

[388] [30] Osborne, C. S. *et al.* Active genes dynamically colocalize to shared sites of ongoing transcription. *Nature genetics* **36**, 1065–71 (2004).

[390] [31] Rieder, D., Trajanoski, Z. & McNally, J. G. Transcription factories. *Frontiers in genetics* **3**, 221 (2012).

[392] [32] Visscher, P. M., Hill, W. G. & Wray, N. R. Heritability in the genomics era–concepts and misconceptions. *Nature Reviews Genetics* **9**, 255–66 (2008).

[394] [33] Churchill, G. A. & Doerge, R. W. Empirical threshold values for quantitative trait mapping. *Genetics* **138**, 963–71 (1994).