

# Extensive epistasis within the MHC contributes to the genetic architecture of celiac disease

Benjamin Goudey<sup>1,2</sup>, Gad Abraham<sup>3</sup>, Eder Kikianty<sup>4</sup>, Qiao Wang<sup>1,2</sup>, Dave Rawlinson<sup>1,2</sup>, Fan Shi<sup>1,2</sup>, Izhak Haviv<sup>5</sup>, Linda Stern<sup>2</sup>, Adam Kowalczyk<sup>1,2,\*</sup>, Michael Inouye<sup>3,\*</sup>

<sup>1</sup>NICTA Victoria Research Lab, The University of Melbourne, Parkville, Victoria 3010, Australia

<sup>2</sup>Department of Computing and Information Systems, The University of Melbourne, Parkville, Victoria 3010, Australia

<sup>3</sup>Medical Systems Biology, Department of Pathology and Department of Microbiology & Immunology, The University of Melbourne, Parkville, Victoria 3010, Australia

<sup>4</sup>Department of Mathematics, University of Johannesburg, PO Box 524, Auckland Park 2006, South Africa

<sup>5</sup>Bar Ilan University, Safed, Israel

\*These authors contributed equally

Correspondence should be addressed to Michael Inouye ([minouye@unimelb.edu.au](mailto:minouye@unimelb.edu.au)) and Adam Kowalczyk ([kowa@unimelb.edu.au](mailto:kowa@unimelb.edu.au))

## Abstract

Epistasis has long been thought to contribute to the genetic aetiology of complex diseases, yet few robust epistatic interactions in humans have been detected. We have conducted exhaustive genome-wide scans for pairwise epistasis in five independent celiac disease (CeD) case-control studies, using a rapid model-free approach to examine over 500 billion SNP pairs in total. We found extensive epistasis within the MHC region with 7,270 statistically significant pairs achieving stringent replication criteria across multiple studies. These robust epistatic pairs partially tagged CeD risk HLA haplotypes, and replicable evidence for epistatic SNPs outside the MHC was not observed. Both within and between European populations, we observed striking consistency of epistatic models and epistatic model distribution, thus providing empirical estimates of their frequencies in a complex disease. Within the UK population, models of CeD comprised of both epistatic and additive single-SNP effects increased explained CeD variance by approximately 1% over those of single SNPs. Further analysis showed that additive SNP effects tag epistatic effects (and vice versa), sometimes involving SNPs separated by a megabase or more. These findings show that the genetic architecture of CeD consists of overlapping additive and epistatic components, indicating that the genetic architecture of CeD, and potentially other common autoimmune diseases, is more complex than previously thought.

## Author Summary

There are few bona fide examples of interactions between genetic variants (epistasis) which affect human disease risk. Here, we assess multiple genome-wide genotyped case-control datasets to investigate the role that epistasis plays in celiac disease, a common immune-mediated illness. We find thousands of replicable, statistically significant pairs of SNPs exhibiting epistasis and, interestingly, all of these fall within the well-known Major Histocompatibility Complex (MHC) region on chromosome 6. We investigate the underlying distribution of epistatic models and further assess the amount of celiac disease variance that can be explained by epistatic pairs, single SNPs and a combination thereof. Our results indicate that there is a substantial amount of shared disease variance between single SNPs and epistatic pairs, but also that a combination of the effects gives a better model of disease. These findings support powerful and routine epistasis scans for the next generation of genome-wide association studies and indicate that the genetic architecture of celiac disease, and potentially other immune-mediated diseases, is more complex than currently appreciated.

## Introduction

The limited success of genome-wide association studies (GWAS) to identify common variants that substantially explain the heritability of many complex human diseases and traits has led researchers to explore other potential sources of heritability, including the low/rare allele frequency spectrum as well as epistatic interactions between genetic variants [1,2]. Many studies are now leveraging high-throughput sequencing with initial findings beginning to elucidate the effects of low frequency alleles [3-6]. However, the characterization of the epistatic component of complex human disease has been limited, despite the availability of a multitude of statistical approaches for epistasis detection [7-13]. Large-scale systematic research into epistatic interactions has been hampered by several computational and statistical challenges mainly stemming from the huge number of variables that need to be considered in the analysis (>100 billion pairs for even a small SNP array), the subsequent stringent statistical corrections necessary to avoid being swamped by large number of false positive results, and the requirement of large sample size in order to achieve adequate statistical power.

The strongest evidence for wide-ranging epistasis has so far come from model organisms [14,15], and recent evidence has demonstrated that epistasis is pervasive across species and is a major factor in constraining amino acid substitutions [16]. Motivated by the hypothesis that epistasis is commonplace in humans as well, recent studies have begun providing evidence for the existence of epistatic interactions in several human diseases, including psoriasis [17], multiple sclerosis [18], type 1 diabetes [19], and ankylosing spondylitis [20]. While these studies have been crucial in demonstrating that epistasis does indeed occur in human disease, several questions remain including how wide-ranging epistatic effects are, how well epistatic pairs replicate in other datasets, how the discovered epistatic effects can be characterized in terms of previously hypothesized models of interaction [21,22], and how much (if at all) epistasis contributes to disease heritability [23].

Celiac disease (CeD) is a complex human disease characterized by an autoimmune response to dietary gluten. CeD has a strong heritable component largely concentrated in the MHC region, due to its dependence on the HLA-DQ2/DQ8 heterodimers encoded by the HLA class II genes *HLA-DQA1* and *HLA-DQB1* [24]. The genetic basis of CeD in terms of individual SNP associations has been well-characterized in several GWAS [25-28], including the additional albeit smaller contribution of non-HLA variants to disease risk [29]. The success of GWAS for common variants in CeD has recently been emphasized by the development of a genomic risk score that could prove relevant in the diagnostic pathway of CeD [30]. Autoimmune diseases have so far yielded the most convincing evidence for epistatic associations, potentially due to power considerations since these diseases usually tend to depend on common variants of moderate to large effect within the MHC. Given these findings in conjunction with recent observations that rare coding variants may play only a negligible role in common autoimmune diseases [3], we sought to determine whether robust epistasis is detectable in CeD and whether it accounts for some of the unexplained disease heritability.

Here, we present the first large-scale exhaustive study of pairwise epistasis in celiac disease. Leveraging GWIS, a highly efficient approach for epistasis detection [31], we conduct genome-wide scans for all epistatic pairs across five separate CeD case/control datasets of European descent, finding thousands of statistically significant pairs despite stringent multiple testing corrections. Next, we show a high degree of concordance of these interactions across the datasets, demonstrating that they are highly robust and replicable. We characterize the common epistatic models found and compare them to previously proposed theoretical models. Finally, we examine the issue of whether epistatic pairs add more predictive power and explain more disease variation than do single SNPs.

## Results

Datasets are summarized in **Table 1**, these include five independent, previously published GWAS datasets of CeD with individuals genotyped from four different European ethnicities: United Kingdom (UK1 and UK2), Finland (FIN), The Netherlands (NL) and Italy (IT) [26,27]. To limit the impact of genotyping error and other sources of non-biological variation, we implemented three stages of validation and quality control (QC): (i) standard QC within each dataset, (ii) independent exhaustive epistatic scans within each of the five datasets, and (iii) derivation of a validated list of epistatic interactions based on UK1. The study workflow is shown in **Figure 1**.

### *Exhaustive epistatic scans and replication*

For each dataset, we implemented stringent sample and SNP level quality control (**Methods**), and then conducted an exhaustive analysis of all possible SNP pairs using the GWIS methodology [31]. Each pair was tested using the GSS statistic, which determines whether a pair of SNPs in combination provides significantly more discrimination of cases and controls than either SNP individually (**Methods**). Forty-five billion pairs were evaluated in the UK1 study (Illumina Hap300/Hap550) and 133 billion SNP pairs were evaluated in each of the four remaining cohorts (Illumina 670Quad and/or 1.2M-DuoCustom). Given this multiple testing burden, we adopted stringent Bonferroni-corrected significance levels of  $P = 1.1 \times 10^{-12}$  for the UK1 and  $P = 3.75 \times 10^{-13}$  for the remaining datasets.

To further ensure that the downstream results were robust to technological artefact and population stratification, we took two additional steps: (a) utilizing the raw genotype intensity data available for UK1 for independent SNP cluster plot inspection (performed by Karen A. Hunt, QMUL), and (b) replicating the epistatic interactions of the SNPs passing cluster plot inspection, where replication is defined as a SNP pair exhibiting Bonferroni-adjusted significance both in UK1 and in at least one additional study. Using these criteria, we found that 7,270 SNP pairs (654 unique SNPs) from the UK1 dataset passed both (a) and (b) above. We denote these pairs as 'validated epistatic pairs' (VEPs) below. **Table 2** presents the top 10 validated pairs, after pruning redundant pairs (pairs of pairs with at least two SNPs in perfect LD); the full list of VEPs is given in **Supplementary Table 1**. Notably, all VEPs fulfilling these robustness criteria were within the MHC.

More than 128,000 unique pairs achieved Bonferroni-adjusted significance across all five studies, with the vast majority lying within the extended MHC region of chr6 (**Figure 2** and **Supplementary Table 2**). Of the 131 epistatic pairs outside the MHC that were significant in at least one study, none passed Bonferroni-adjusted significance in at least one other study and were thus deemed not replicated. As expected, the number and strength of epistatic interactions increased as sample size increased. Interestingly, some of the strongest epistatic interactions tended to be in close proximity and in moderate LD, though only 1% of pairs had  $r^2 > 0.5$  (**Supplementary Figure 1**). The heatmaps in **Figure 2** also showed that epistasis was widely distributed with distances of  $>1\text{Mb}$  common between epistatic pairs. Across all studies, epistatic interactions were consistently located in and around HLA class II genes, however further examination of the VEPs found that many of the top pairs were proximal to HLA class III genes,  $>1\text{Mb}$  upstream of *HLA-DQA1* and *HLA-DQB1*, the strongest known risk loci (**Supplementary Figure 2**).

The extent of replication of the epistatic pairs was apparent from the high degree of similarity in the rankings when pairs were sorted by GSS significance (**Figure 3a**), with  $\sim 70\text{-}80\%$  overlap between the UK1 and UK2 datasets extending all the way to the top 10,000 pairs, and  $40\text{-}60\%$  overlap with the pairs found in the NL and FIN datasets. Such high degrees of overlap have essentially zero probability of occurring by chance ( $P < 10^{-600}$  for  $\sim 80\%$  overlap between the UK1 and UK2 top 50 pairs, hypergeometric test). The pairs found in the IT dataset showed lower levels of consistency with those detected in the UK1 dataset but overall were still much more than expected by chance

with ~30% overlap at ~30,000 pairs ( $P < 10^{-1000}$ ).

### ***Empirical epistatic model distributions***

The epistatic model provides insight into how disease risk is distributed across the nine pairwise genotype combinations. Following the conventions of Li and Reich [21], we discretized the models for the VEPs to use fully-penetrant values where each genotype combination implies a susceptibility or protective effect on disease (**Methods**), simplifying the comparison of models between different SNP pairs.

To establish model consistency, we first replicated the most frequent full penetrance VEP models in the other datasets (**Figure 4**). When considering the distribution of epistatic models we found striking consistency of the UK1 models with those from UK2 and the other Northern European populations (Finnish and Dutch) (**Figure 4**). Only four models from the possible 50 classes [21] occurred with >5% frequency in the Northern European studies, and there was substantial variation in epistatic model as a function of the strength of the interaction. Amongst all VEPs in UK1, the four models corresponded to the threshold model (T; 34.7% frequency), jointly dominant-dominant model (DD; 31.1%), jointly recessive-dominant model (RD; 17.9%), and modifying effect model (Mod; 14.73%) [21,32]. The DD and RD models are considered multiplicative, the Mod model is conditionally dominant (i.e. one variant behaves like a dominant model if the other variant takes a certain genotype), and the T model is recessive. The T model was the most frequent model, especially amongst the strongest pairs.

To our knowledge, the frequencies of these epistatic models have not previously been determined in a complex human disease. Interestingly, despite the consistency of MHC epistasis, the VEPs showed noticeable differences in epistatic model distribution in the IT population. This was in contrast to the other Northern European populations but consistent with the different ranking in GSS significance observed above. In the IT population, the distribution of models was altered such that there was a more even distribution. The four most frequent models were still the T model (13.6%), DD (13.45%), modifying effects (13.45%), and RD model (11.55%). But, we also observed that many of the strongest pairs within the IT cohort followed the M86 model, though M86 represented only a small proportion of models overall (0.98%). In IT, the remaining 50% of the VEP models overall consisted of many low-frequency models.

The cause(s) of these differences is unclear. While cryptic technical factors cannot be ruled out at this stage, it may be the case that there is population specific epistatic variation that follows the known North/South European genetic gradient [33].

### ***Contribution of epistatic pairs to celiac disease heritability***

We next sought to estimate the CeD heritability explained by the VEPs and single SNPs. The GSS test selects each epistatic pair based on it being more predictive than either of its constituent individual SNPs. However, the procedure does not *a priori* guarantee that a given pair is a better predictor of disease than all other individual SNPs not included in the pair; this requires a further step to determine which pairs and/or single SNPs provide the most predictive power overall. This task is further complicated by the fact that, like linkage disequilibrium for individual SNPs, many pairs are highly correlated and thus may not add substantial predictive power after accounting for the most predictive pair. These issues can naturally be addressed within the framework of a multivariable model, accounting for all SNPs and/or pairs at once. Hence, to better assess the contribution of epistatic pairs to CeD prediction and thus heritability explained, we employed L1 penalized linear support vector machines (SVM, see **Methods**), an approach which models all variables concurrently (SNPs and/or pairs) and which has been previously shown to be particularly suited for maximizing predictive ability from SNPs in CeD and other autoimmune diseases [30,34]. While we find that VEPs, as expected, are associated with the *HLA-DQA1* and *HLA-DQB1* risk haplotypes (**Supplementary Figure 3**), we have previously found that additive models of single



SNPs explain substantially more CeD variance than haplotype-based models [30]. We therefore employ the former to estimate the gain in heritability here.

We assessed CeD variance explained by constructing three separate models: (a) genome-wide single SNPs only, using the 290,277 SNPs present across all datasets, (b) the VEPs only, i.e. the 7270 VEPs encoded as 65,430 indicator variables, and (c) a 'combined' model of both single SNPs and VEPs together. The models were evaluated in cross-validation on the UK1 dataset, and the best models in terms of Area Under the Curve (AUC) were then taken forward for external validation in the other four datasets without further modification.

In UK1 cross-validation, the VEPs and combined models led to an increase in maximum AUC of 0.5% over single SNPs alone (0.883 to 0.888), corresponding to an additional ~1.5% in explained CeD variance, from 32.6% to 34.1% (**Table 3**). External validation of these models showed that the VEPs and combined models showed similar and significant gains in AUC over the single SNPs for the UK2 and IT dataset at +1% for IT ( $P=0.0163$ ) and +0.9% for UK2 ( $P=0.0066$ ), but the differences in FIN and NL were smaller and not significant. In external validation, the VEPs model was highly predictive yet slightly less predictive than that based on single SNPs, with the combined model yielding the highest AUC. The increased sample size of a combined UK1 and UK2 dataset in cross-validation did not yield better AUCs nor corresponding CeD variance (**Supplementary Table 3**).

## Discussion

This study has shown the robust presence of epistasis in celiac disease. Epistatic interactions were observed within the extended MHC, most strongly between neighbouring SNPs in low to moderate LD, indicating that these interactions may play a role in segregating specific haplotype classes. We have shown that these epistatic SNP pairs strongly replicate across cohorts in terms of significance, ranking, and epistatic model. To our knowledge, this level of epistatic signal strength, number of epistatic pairs, and degree of replication has not been previously shown in a complex human disease.

Despite observations that epistatic interactions between SNPs within a locus are enriched for batch effects and poorly clustered genotype clouds [37], the stringent quality control and extensive replication of the analyses in this study indicate that these SNPs are largely bona fide epistatic pairs. When considering those pairs not achieving the Bonferroni significance criteria for replication, a large number of epistatic pairs were still highly statistically associated with CeD consistently across datasets, indicating that our estimates of epistasis may be conservative. For validated epistatic pairs (VEPs), we found that much of strongest epistatic signal is over 1MB upstream of the well-known *HLA-DQA1* and *HLA-DQB1* risk loci, suggesting a potentially important contribution of HLA class III genes. We also performed a large-scale empirical characterization of the epistatic models underlying the interactions in CeD, with the majority of the VEPs approximately following the threshold model, and a smaller number following dominant-dominant, dominant-recessive, and recessive-recessive models. Further, these patterns were found to be strongly consistent across most of the datasets.

We have previously found that penalized predictive models based on individual SNPs similar to those used here are able to extract more predictive ability from the MHC region than models based on coarse-grain HLA types [30]. Here, we have found that combined models of both epistatic SNP pairs and single SNPs achieve slightly improved accuracy over models created with single SNPs alone, and that models of only epistatic SNP pairs explained similar amounts of CeD variance as single SNPs. Examining this redundancy more closely, the epistatic SNP pairs are highly correlated with single SNPs that are usually located near one of the pair (see **Supplementary Figure 7**). This

correlation between single SNPs and combinations of SNPs appears to have been previously hinted at in a study by de Bakker et al examining the effectiveness of SNPs to tag HLA genotypes, where groups of SNPs were found to be more highly correlated with HLA genes than single SNPs [38]. The shared information between these single SNP and epistatic effects implies that determining the causal signal will be more difficult than previously thought. Just as the redundancy between single SNPs in LD has affected the resolution of causal genetic variants, our findings indicate that a similar, though currently unexplored, sharing of information may exist between epistatic variants and single variants. Such an observation is supported by previous literature [39] and may help to explain some of the controversy around epistatic versus additive genetic effects.

Celiac disease has a strong HLA signal, is highly heritable and is thought to conform to the Common-Disease, Common-Variant (CDCV) model [24]. Yet within this 'model disease' [25], our results suggest the presence of a previously unexplored level of complexity. Given their similar disease etiologies [40], we predict that these observations may hold true for other autoimmune/inflammatory diseases and other diseases that approximate the CDCV model. It is less likely that these observations affect our understanding of complex diseases that are unlikely to approximate CDCV, such as coronary artery disease, though it has been proposed that epistasis plays a role for these types of conditions as well [1].

The limitations of the first generation GWAS approach to explain missing heritability has led to the development and application of more complex approaches to resolve this problem, yet success has been elusive. Recent results suggest that rare variants add little to known heritability for a number of autoimmune diseases including celiac disease [26]. The predictive models generated in this work indicate that while epistatic pairs have substantial predictive power, their overall explained heritability is not substantially more than that for additive effects. Combined models of epistatic and additive effects are likely to constitute the best solution, however it is unlikely that these alone will resolve missing heritability.

These findings have implications for how next generation GWAS should be analysed and interpreted. While epistatic analyses have increasingly been advocated [27], this study demonstrates the usefulness of such an approach alongside that of traditional genome-wide analysis of additive effects. Many challenges remain in conducting this type of analysis. While we found strong epistasis within the MHC, future advances in statistical methods could uncover additional epistasis with weaker effects or involving rare variants, and it is currently unknown how weaker and rare variant epistatic effects interact with additive effects in humans. A main challenge of genetic association studies, the inference of genetic architecture, may very well be complicated by the shared information between epistatic and additive effects and it may be that targeted perturbation experiments will be required to identify the true causal signal.

## Methods

### *Quality control*

A range of quality control measures were applied to all datasets to limit the impact of genotyping error. For all datasets, we removed non-autosomal SNPs, SNPs with MAF <1%, missingness >1% and those deviating from Hardy Weinberg Equilibrium in controls with  $P < 5 \times 10^{-6}$ . Samples were removed if data missingness was >1%. Cryptic relatedness was also stringently assessed by examining all pairs of samples using identity-by-descent in PLINK, and removing one of the samples if  $\hat{\pi} > 0.05$ . The cryptic relatedness filter removed 17 samples within the UK1 cohort that related to other UK1 samples, and 1208 samples from the UK2 cohort which were either related to other UK2 samples or UK1 samples. Dataset sizes in Table 1 are reported after the quality control

steps above. Significant epistatic SNP pairs were further assessed by manually inspecting the genotyping cluster plots of both SNPs in the UK1 cohort. Intensity data for the other studies was not available thus epistatic pairs discovered in these datasets were not classified as robust and were not used in heritability estimates, however the consistency of the statistics and epistatic model across independent datasets indicated that many likely represent bona fide epistasis. Cluster plot inspection removed 115 SNPs with poor genotyping assays.

### *Statistical tests for epistasis*

Here, we briefly describe the intuition behind the Gain in Sensitivity and Specificity (GSS) test we employ to detect epistasis, and later we present several approximations we employ in analyzing epistasis across datasets. The test has been presented in detail in [31] and, currently, a web server implementing the GSS test is at <http://bioinformatics.research.nicta.com.au/software/gwis/>.

There is a long history of discussion around the exact definition of epistasis, or gene-gene interaction [41]. Here, we use a definition that is closely aligned with the multifactor dimensionality reduction (MDR) family of gene-gene interaction methods [42]: an epistatic interaction is defined as a significant improvement of a SNP-pair in classifying cases from controls over what is possible using each SNP individually. There are two main differences between our approach and similar approaches for detecting epistasis [8,22]. First, our approach is “model-free”, as it makes no assumptions about the way in which genotypes combine to affect the phenotype [7,43], but considers all possible pairwise interactions for each pair, making it potentially more powerful to detect unknown epistatic forms, as empirical knowledge about epistasis in humans is currently lacking. Second, instead of measuring the deviation from additive effects (for example, using a likelihood ratio test), our approach focuses on the utility of the test in case/control classification, quantified using the receiver-operating characteristic (ROC) curve, and measuring the deviation in the curve from that induced by the additive model.

The main principle behind the GSS is quantification of the gain in predictive power afforded by a putative epistatic pair over and above the predictive power due to each of its constituent SNPs. The difference in predictive power is assessed in terms of the ROC curves induced by the pair and each of the SNPs. The ROC curve is formed by considering each possible genotype (or pair of genotypes), and measuring the sensitivity (true positive rate, TPR) and specificity (1 – false positive rate, FPR) at that point, and ordering them in decreasing order by the ratio TPR/FPR; hence the curve is piecewise linear. Since the two ROC curves induced by the individual SNPs may intersect, we represent them using a convex hull, which is the best ROC curve that can be produced by any linear combination of the two individual SNPs, and represents a conservative estimate of the predictive power of the individual SNPs. The GSS then assigns a p-value to each point in the pair’s ROC curve, based on the probability of observing a combination of genotypes with a higher or equal TPR and a lower or equal FPR, under the null hypothesis that the true TPR and FPR reside below the convex hull. We employ a highly efficient minimax-based implementation, maximizing the probability for each point on the ROC curve (worst case scenario) against all points of the convex hull, and returning the minimum probability over all points [31]; this is done using an exact procedure rather than relying on approximations based on the normal distribution. Finally, the best p-value is assigned as the overall p-value for the pair, allowing the pairs to be ranked and corrected for multiple testing as is standard practice in GWAS. Those SNPs that are significant after multiple testing correction are deemed significant epistatic pairs.

Analogously to odds ratios used for analyses of single SNPs, we can estimate odds ratios for epistatic pairs based on the GSS statistic

$$OR_{GSS} = \frac{(\pi_{\{0,HR\}})(\pi_{\{1,LR\}})}{(\pi_{\{1,HR\}})(\pi_{\{0,LR\}})},$$



where  $\pi_{\{i,j\}}$  denotes the proportion of samples with phenotype  $i$ , 0 for cases and 1 for controls, and carrying genotype combinations which are marked as  $j$  with *HR* (high risk) indicating genotypes which are associated by GSS with cases and *LR* (low risk) indicating genotypes which are associated with controls. By relying on the model-free GSS approach, this odds ratio can be seen as deriving the specific model maximizing the level of improvement over that of the individual SNPs in the pair.

### *Approximate representation of the epistatic models*

While the GSS approach is the basis for detecting epistatic pairs, the models it produces can be hard to visually interpret and categorize into broad groups. To simplify interpretation, we approximate the models for the statistically significant pairs found via GSS using two representations: balanced penetrance models and full penetrance models.

#### *Balanced penetrance models*

Following Li and Reich [21] we employ the penetrance, that is, the probability of disease given the genotype, estimated from the data for each of the nine genotype combinations as (number of cases with combination)/(number of individuals with combination). Representing the epistatic model in terms of penetrance allows us to clearly see which genotype combinations contribute more to disease risk (or conversely, may be protective).

One limitation of the penetrance is that it is typically considered in isolation of the disease background rate (the prevalence), which may be misleading when comparing penetrance levels across datasets with widely varying proportions of cases. For example, a penetrance of 50% for a given SNP would be considered very high in a dataset consisting of 1% cases and 99% controls, but no better than random guessing in datasets with 50%/50% cases and controls. Hence, we employ a standardization to ensure that the penetrance is comparable across datasets, termed *balanced sample penetrance*, and defined as

$$P_{balanced} = \frac{p_{1v}}{p_{1v} + p_{0v}},$$

where  $p_{iv}$  refers to the proportional frequency of genotype  $v$  in class  $i$ , where controls are 0 and cases are 1 (0=controls, 1 = cases). The balanced sample penetrance ranges between 0 and 1, where 0 means that the genotype only occurs in controls, 1 means that the genotype only occurs in cases and 0.5 means the genotype occurs evenly between the two classes. Balanced penetrance can be related to either standard penetrance or relative risk in the data via monotonic transformations. The definition is easily extended to the case of pair of SNPs. The only difference is the use of the  $3 \times 3 = 9$  possible genotype combinations from each SNP-pair rather than the 3-value set of genotypes from an individual SNP.

#### *Simplification to full penetrance models*

The balanced-penetrance epistatic models provide fine-grained insight into the relative effects of each genotype combination. In addition, we employ a coarse-grain approach where these values are discretized into binary values (0/1), so called “fully penetrant” models, an approach analogous to that of Li and Reich [21]. These binary models forgo some detail but make it easier to categorize epistatic models into broad classes based on their patterns of interaction, such as the classic XOR pattern [8] or the threshold model [22]. Swapping major and minor alleles, and swapping the SNP ordering in the contingency table, can reduce the number of fully penetrant models. Unlike Li and Reich, we do not swap the high and low risk status, as we are interested in distinguishing between protective and deleterious combinations. Furthermore, Li and Reich also excluded models with all high or low risk genotypes. Such models can not exist within the set we are analyzing as they would show no association with disease. Li and Reich were able to show that there are only 51 possible fully penetrant disease models after accounting for symmetries. However, as we do not swap risk

status, there will be 100 possible full-penetrance models that can appear within the analysis conducted here [22].

Given that some genotype combinations in certain SNP pairs are rare, there may be insufficient evidence to determine whether they have a substantial effect on disease risk. As such, we have used a simple heuristic for such entries, denoting all cells with a frequency below 2% in both cases and controls as ‘low risk’. Experiments with this threshold revealed that altering this cutoff between 0% and 7% made little difference to the overall distribution of our models.

### ***The predictive models***

We employed a sparse support vector machine (SVM) implemented in SparSNP [44]. This is a multivariable linear model where the degree of sparsity (number of variables being assigned a non-zero weight) is tuned via penalization. The model is induced by minimizing the L1-penalized squared hinge loss

$$(\beta^*, \beta_0^*) = \arg \min_{\beta, \beta_0} \frac{1}{2N} \sum_{i=1}^N \max\{0, 1 - y_i(x_i^T \beta + \beta_0)\}^2 + \lambda \sum_{j=1}^p |\beta_j|$$

where  $\beta$  and  $\beta_0$  are the model weights and the intercept, respectively,  $N$  is the number of samples,  $p$  is the number of variables (SNPs and/or encoded pairs),  $x_i$  is the  $i$ th vector of  $p$  variables (genotypes and/or encoded pairs),  $y_i$  is the  $i$ th case/control status  $\{+1, -1\}$ , and  $\lambda \geq 0$  is the L1 penalty. To find the optimal penalty, we used a grid of 100 penalty values within 10 replications of 10-fold cross-validation, and found the model/s that maximized the average area under the receiver-operating characteristic curve (AUC). For models based on single SNPs, we used minor allele dosage  $\{0, 1, 2\}$  encoding of the genotypes. For models based on SNP pairs, the standard dosage model is not applicable; hence, we transformed the variable representing each pair (encoded by integers 1 to 9) to 9 indicator variables using the Python library scikit-learn [45], using a consistent encoding scheme across all datasets. The indicator variables were then analyzed in the same way as single SNPs. Results were analyzed in R [46] with the packages ROCR [47] and pROC [48], and plotted using the ggplot2 [49] package.

### ***Evaluation of predictive ability and explained disease variance***

To maximize the number of SNPs available for analysis, we imputed SNPs in the UK2, FIN, NL, and IT dataset to match those that were in the UK1 dataset but not in former, using IMPUTE v2.3.0 [50]. Post QC this left 290,277 SNPs common to all five datasets. Together with  $9 \times 7270$  pairs = 65,430 indicator variables, this led to a total of 355,707 variables in the combined singles+pairs dataset. Models trained in cross-validation on the UK1 dataset were then applied without any further tuning to the four other datasets, and the external-validation AUC for these models was then estimated within the validation datasets. To derive the proportion of phenotypic variance explained by the model (on the liability scale), we used the method of Wray et al. [51], assuming a population prevalence of 1%.

## **Acknowledgements**

MI was supported by an NHMRC early career fellowship 637400. MI and GA were supported by University of Melbourne funding. BG, EK, QW, DR, FS, IH and AK were supported by National ICT Australia (NICTA). NICTA is funded by the Australian Government’s Department of Communications, Information Technology and the Arts, the Australian Research Council through Backing Australia’s Ability, and the ICT Centre of Excellence programs.

We thank the investigators of the van Heel et al., 2007 and Dubois et al., 2010 papers (David van

Heel and Cisca Wijmenga) for providing the celiac disease data. We thank Karen A. Hunt (QMUL) for performing cluster plot inspection on the UK1 data. We also thank Rami Mukhtar for useful technical advice regarding implementation of algorithms used here and Andrew Kowalczyk and Leon Gor for assistance with development of software utilised for this work. We also thank Armita Zarnegar for assistance with processing of data and John Markham, Justin Bedo and Geoff Macintyre for insightful discussions and comments.

## Tables

**Table 1: Datasets**

		SNPs <sup>a</sup>	Celiac cases		Controls		Ref
			Samples <sup>a</sup>	Platform <sup>b</sup>	Samples <sup>a</sup>	Platform <sup>b</sup>	
UK1	UK	301546	763	Illumina Hap300v1-1	1420	Illumina Hap550	(van Heel, et al., 2007)
UK2	UK	515413	1826	Illumina670-Quad	3777	Illumina 1.2M-Duo	(Dubois, et al., 2010)
FIN	Finland	513952	647	Illumina670-Quad	1829	Illumina 610-Quad	(Dubois, et al., 2010)
NL	Netherlands	515169	803	Illumina670-Quad	846	Illumina 670-Quad	(Dubois, et al., 2010)
IT	Italy	515641	497	Illumina670-Quad	543	Illumina 670-Quad	(Dubois, et al., 2010)
Overlapping SNPs		286938					

- The number of samples/SNPs is reported after quality control procedures were applied.
- All platforms contain a common set of Hap300 markers; the Hap550 and 610-Quad contain a common set of Hap550 markers.

**Table 2: Top 10 epistatic signals detected in UK1 cohort within the extended MHC region and their properties in the remaining four cohorts**

Rank	SNP	Chr	Position (bp) <sup>f</sup>	UK1 <sup>univariate</sup> MAF <sup>c</sup>	$\chi^2$	UK1 LD <sup>d</sup>	GSS <sup>a</sup>	OR <sup>b</sup>	UK2 GSS <sup>a</sup>	OR <sup>b</sup>	FIN GSS <sup>a</sup>	OR <sup>b</sup>	NL GSS <sup>a</sup>	OR <sup>b</sup>	IT GSS <sup>a</sup>	OR <sup>b</sup>
1	rs2260000	6	31701455	0.28	40.9	0.68	58.2	14.2	108.3	10.8	95.4	20.5	27.0	10.1	12.6	6.7
	rs805262	6	31736712	0.47	24.7											
2	rs2535315	6	31160106	0.30	16.3	0.50	51.4	8.9	98.3	7.4	74.8	14.2	30.8	7.7	6.6	5.7
	rs2517452	6	31168141	0.43	29.5											
3	rs805303	6	31724345	0.47	69.3	0.61	50.5	14.3	97.6	10.7	77.0	20.1	29.1	16.1	9.8	7.0
	rs805274	6	31773173	0.23	6.9											
4	rs2269426	6	32184477	0.31	31.3	0.23	48.3	11.7	90.9	9.6	56.2	12.2	25.0	11.5	8.5	3.8
	rs394657	6	32295001	0.48	54.9											
5	rs9357152	6	32772938	0.19	47.8	0.13	46.8	9.3	94.1	8.0	17.2	7.7	13.3	6.3	5.7	3.4
	rs9276644	6	32853021	0.42	38.7											
6	rs241440	6	32905339	0.20	20.2	0.58	46.2	7.72	104.3	7.4	80.1	15.9	28.1	7.0	12.3	4.0
	rs241437	6	32905662	0.45	25.7											
7	rs3117098	6	32466491	0.24	14.5	0.59	44.7	10.3	92.2	8.8	43.0	12.0	23.3	7.5	9.2	3.1
	rs6932542	6	32488240	0.45	42.7											
8	rs2395488	6	31553888	0.42	56.2	0.47	44.5	11.9	95.1	9.8	74.9	18.9	29.5	8.9	14.9	5.8
	rs2523647	6	31557757	0.16	12.5											
9	rs2269426	6	32184477	0.31	31.3	0.42	42.1	12.4	71.4	8.9	47.8	10.0	24.4	11.9	10.2	5.4
	rs2269423	6	32253685	0.32	35.7											
10	rs7192	6	32519624	0.47	43.1	0.57	40.6	10.3	81.8	8.9	47.3	12.0	24.5	7.9	7.3	3.6
	rs2395182	6	32521295	0.17	15.7											

- GSS indicates the  $-\log_{10}(\text{p-value})$  of improvement of the pair over each of the SNPs involved measured by the GSS filter described further in the Methods section
- Odds Ratios are calculated directly from the GSS rather than via logistic regression, discussed further in Methods.
- Minor Allele Frequency measured in the Control samples in the UK1 cohort
- $r^2$  was taken from HAPMAP release 2 using the CEU population
- Each signal represents the strongest of any pairs that show an  $r^2 > 0.7$  with both of the SNPs in the pair
- SNP positions were extracted from build 36  
X2 indicates  $\log_{10}(\text{p-value})$  for the standard  $\chi^2$  test of association ( $\chi^2$  statistics with 2 degrees of freedom).



**Table 3: Variance explained by models with additive and epistatic genetic effects**

		Single SNPs		Single SNPs + Pairs		Pairs	
		Variance explained	AUC (95% CI)	Variance explained	AUC (95% CI)	Variance explained	AUC (95% CI)
<b>Cross validation</b>	UK1 290K SNPs (best model)	0.326	0.882 [0.880, 0.883]	0.341	0.888 [0.886, 0.889]	0.342	0.888 [0.887, 0.889]
<b>External validation</b>	UK2	0.269	0.855 [0.844, 0.865]	0.279	0.860 [0.850, 0.870]	0.249	0.845 [0.834, 0.855]
	Finn	0.325	0.880 [0.865, 0.896]	0.330	0.884 [0.870, 0.898]	0.297	0.868 [0.853, 0.884]
	IT	0.267	0.853 [0.830, 0.876]	0.290	0.867 [0.844, 0.889]	0.243	0.841 [0.817, 0.865]
	NL	0.275	0.858 [0.839, 0.876]	0.273	0.858 [0.840, 0.876]	0.254	0.847 [0.828, 0.866]

Predictive power of single SNPs and pairs in cross-validation and in external validation, using SparSNP models. Models were optimized on the UK1 dataset (n=2183 samples) in cross-validation (290K SNPs), and tested without modification on the other datasets. The proportion of heritability explained (on the liability scale) assumes a population prevalence of 1%. The 95% CI for AUC in UK1 was computed over the 10x10 cross-validation, and in external validation was computed using DeLong's method (R package pROC). Two-sided DeLong significance tests for AUC of single SNPs+pairs difference from AUC of single SNPs: UK2 P=0.006592, FIN P=0.3409, IT P=0.01626, NL P=0.8966.

## Figure Legends

### Figure 1: Study workflow

### Figure 2: Epistatic interactions within the extended MHC region

SNP pairs within 30KB of each other are shown as a single point on each heatmap. The colour of each point represents the most significant  $-\log_{10}(\text{P-value})$  returned by the GSS statistic for SNPs pairs within each point. The  $-\log_{10}(\text{P-value})$  is capped at 30 to increase contrast of lower values. The distribution of higher values in these datasets is shown in **Supplementary Figure 4**. The differences in the number of significant pairs detected in each cohort are clearly associated with the relative power of each study.

### Figure 3: Replication of epistatic pairs and corresponding epistatic models between datasets and populations

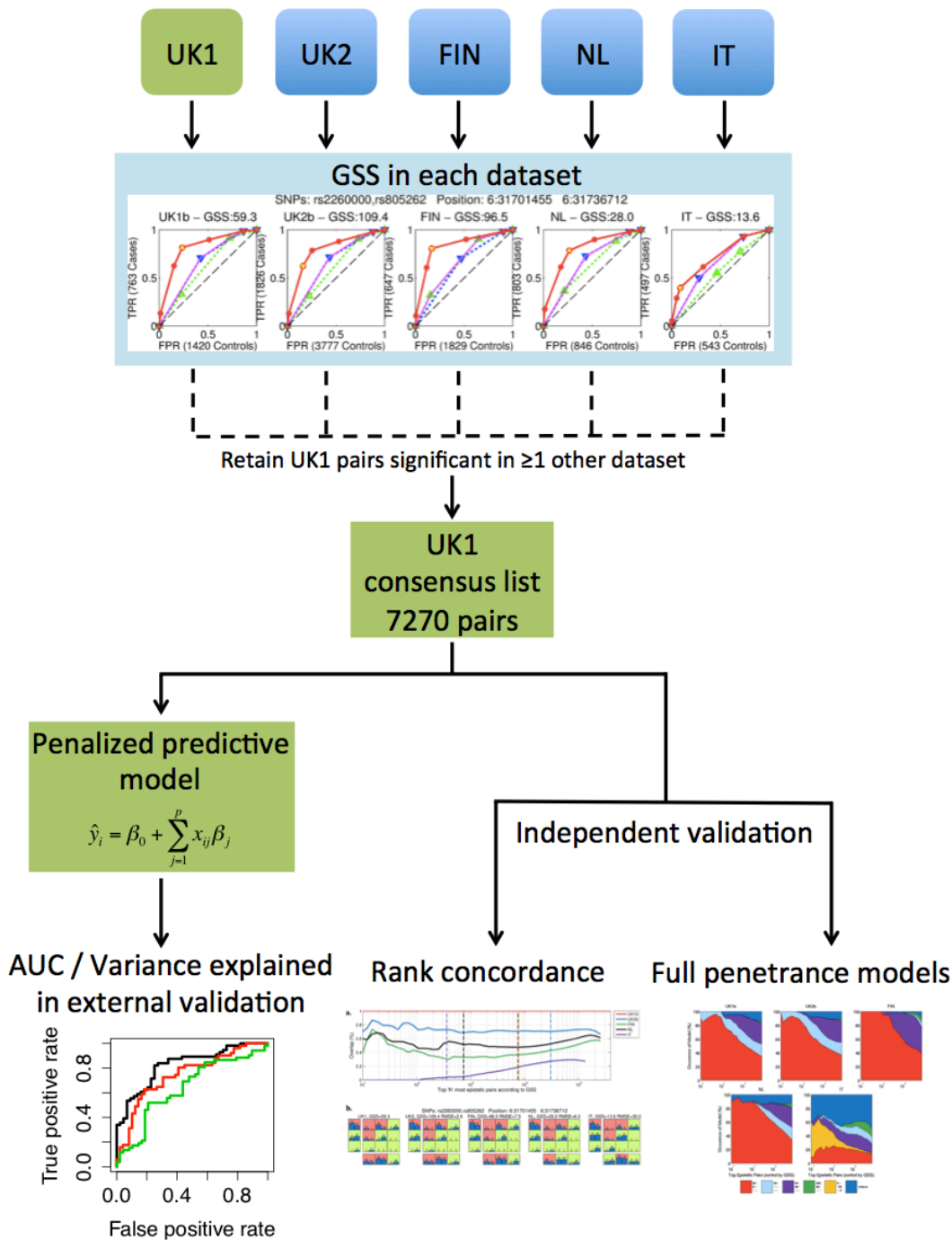
Panel **(a)** shows the overlap of significant epistatic pairs as a percentage between UK1 and remaining cohorts in order of decreasing GSS significance. Vertical dotted lines indicate the Bonferroni-adjusted significance for each study. Panel **(b)** shows the occurrence of genotype combinations for the top pair from UK1. Colouring of cells provides an indication of the epistatic model occurring in each cohort, explained further in the Methods section. Further examples are shown in **Supplementary Figure 5**.

### Figure 4: Variation in epistatic models within and between populations

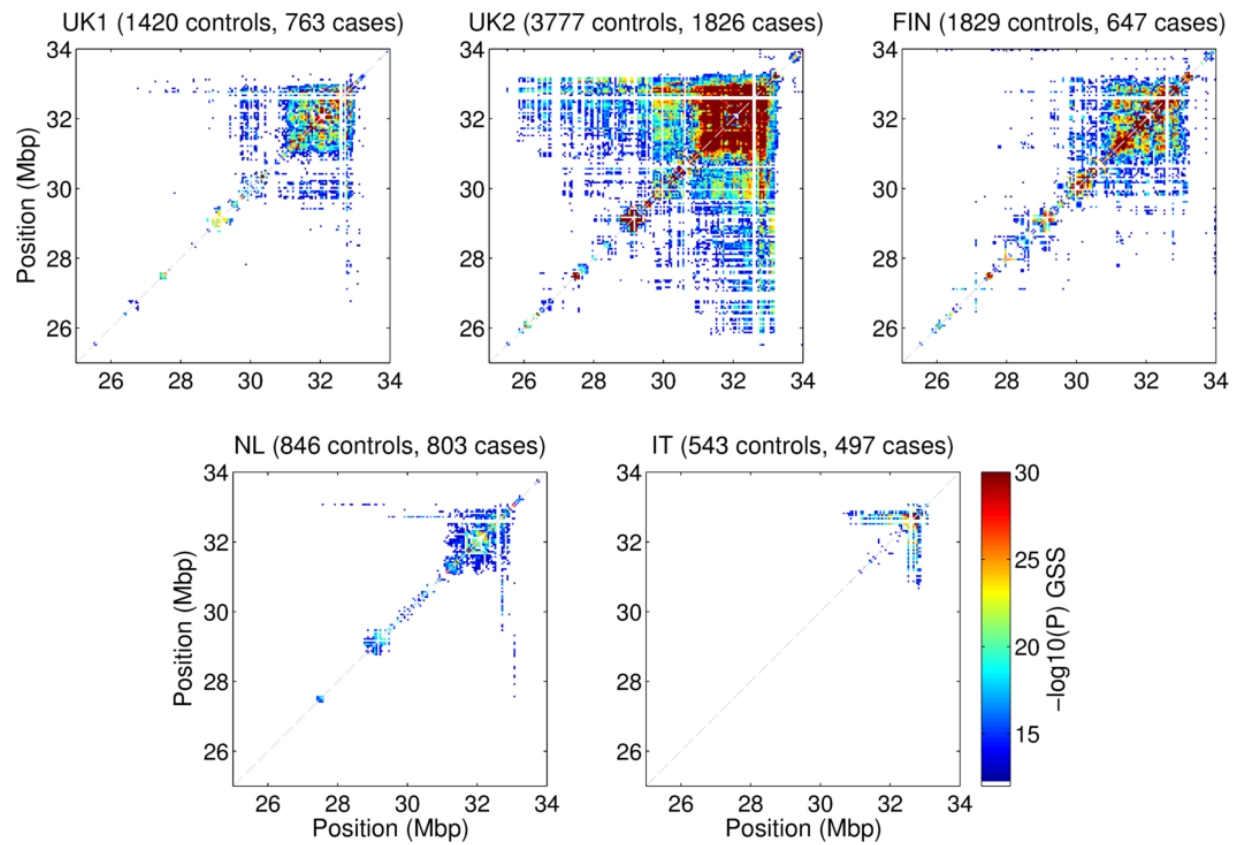
Distribution of epistatic models in different studies as increasing less significant SNP pairs are examined, where the models were selected based on the UK1 dataset. Different colours represent a different subset of epistatic models. The “other” group represents the set of models that occur less than 5% of the time. Models have been simplified using the rules provided in (Li & Reich, 2000). **Supplementary Figure 6** examines the distribution of models for pairs where at least one genotype combination does not occur in the data.

## Figures

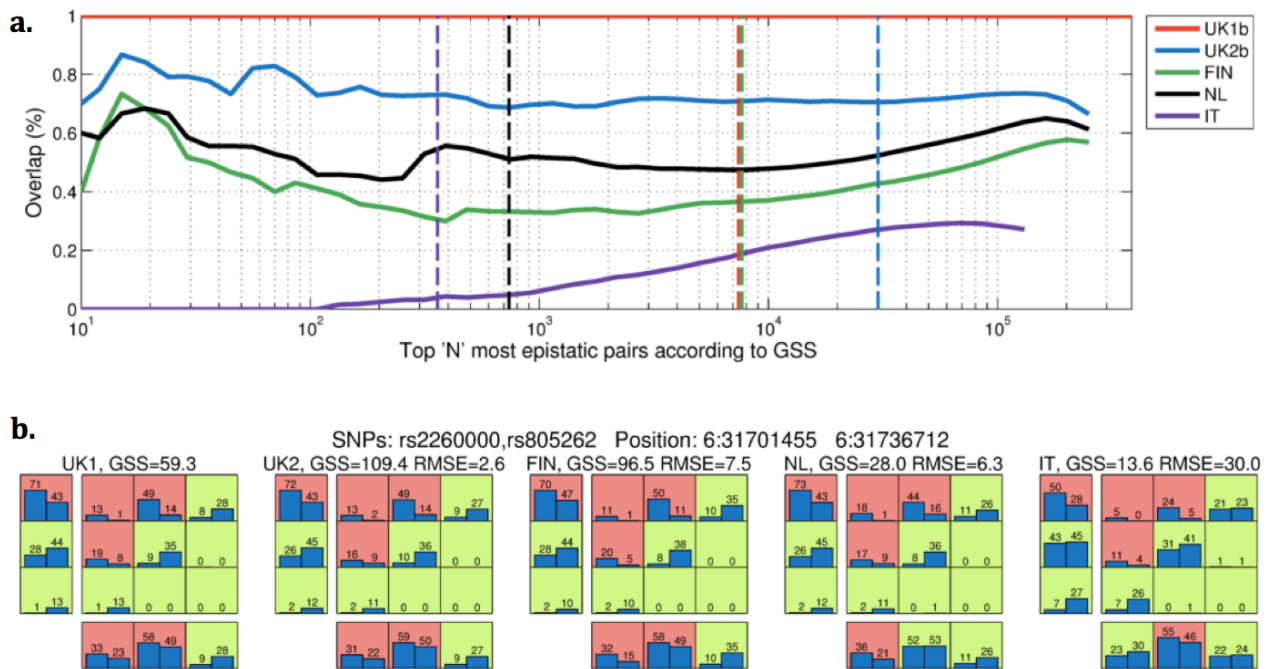
**Figure 1: Study workflow**



**Figure 2: Epistatic interactions within the extended MHC region**

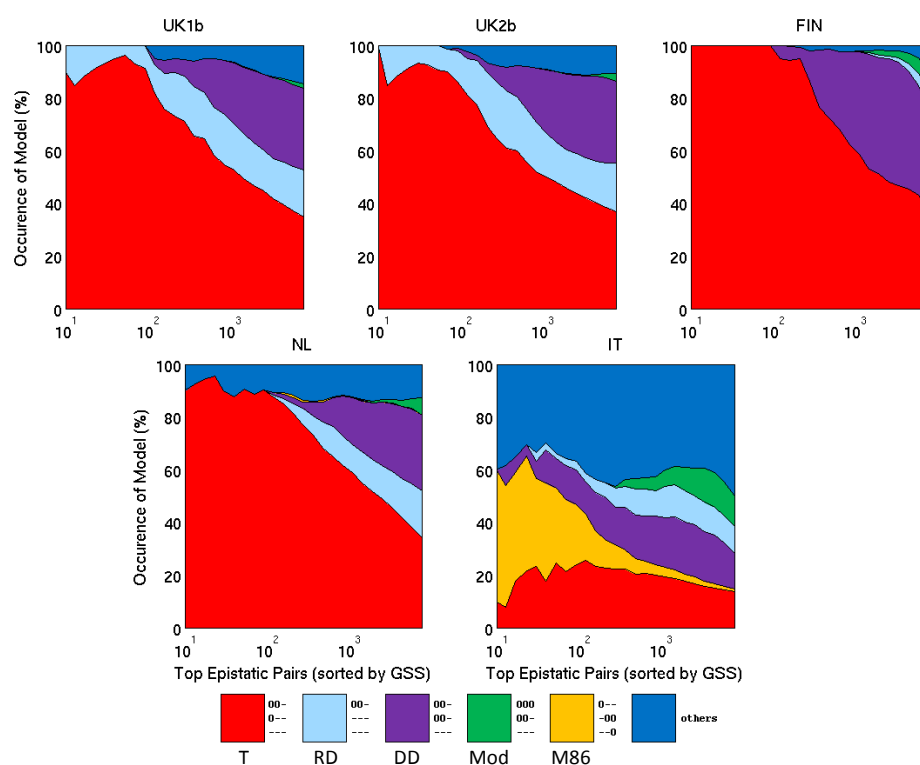


**Figure 3: Replication of epistatic pairs and corresponding epistatic models between datasets and populations**





**Figure 4: Variation in epistatic models within and between populations**



# References

1. Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, et al. (2009) Finding the missing heritability of complex diseases. *Nature* 461: 747-753.
2. Eichler EE, Flint J, Gibson G, Kong A, Leal SM, et al. (2010) Missing heritability and strategies for finding the underlying causes of complex disease. *Nat Rev Genet* 11: 446-450.
3. Hunt KA, Mistry V, Bockett NA, Ahmad T, Ban M, et al. (2013) Negligible impact of rare autoimmune-locus coding-region variants on missing heritability. *Nature* 498: 232-235.
4. Huyghe JR, Jackson AU, Fogarty MP, Buchkovich ML, Stancakova A, et al. (2013) Exome array analysis identifies new loci and low-frequency variants influencing insulin processing and secretion. *Nat Genet* 45: 197-201.
5. Jonsson T, Atwal JK, Steinberg S, Snaedal J, Jonsson PV, et al. (2012) A mutation in APP protects against Alzheimer's disease and age-related cognitive decline. *Nature* 488: 96-99.
6. Stykarsdottir U, Thorleifsson G, Sulem P, Gudbjartsson DF, Sigurdsson A, et al. (2013) Nonsense mutation in the LGR4 gene is associated with several human diseases and other traits. *Nature* 497: 517-520.
7. Cordell HJ (2002) Epistasis: what it means, what it doesn't mean, and statistical methods to detect it in humans. *Hum Mol Genet* 11: 2463-2468.
8. Wan X, Yang C, Yang Q, Xue H, Fan X, et al. (2010) BOOST: A fast approach to detecting gene-gene interactions in genome-wide case-control studies. *Am J Hum Genet* 87: 325-340.
9. Hu X, Liu Q, Zhang Z, Li Z, Wang S, et al. (2010) SHEsisEpi, a GPU-enhanced genome-wide SNP-SNP interaction scanning algorithm, efficiently reveals the risk genetic epistasis in bipolar disorder. *Cell Res* 20: 854-857.
10. Hemani G, Theodoridis A, Wei W, Haley C (2011) EpiGPU: exhaustive pairwise epistasis scans parallelized on consumer level graphics cards. *Bioinformatics* 27: 1462-1465.
11. Prabhu S, Pe'er I (2012) Ultrafast genome-wide scan for SNP-SNP interactions in common complex disease. *Genome Res* 22: 2230-2240.
12. Moore JH, Gilbert JC, Tsai CT, Chiang FT, Holden T, et al. (2006) A flexible computational framework for detecting, characterizing, and interpreting statistical patterns of epistasis in genetic studies of human disease susceptibility. *J Theor Biol* 241: 252-261.
13. Zhang Y, Liu JS (2007) Bayesian inference of epistatic interactions in case-control studies. *Nat Genet* 39: 1167-1173.
14. Carlborg O, Haley CS (2004) Epistasis: too often neglected in complex trait studies? *Nat Rev Genet* 5: 618-625.
15. Corbett-Detig RB, Zhou J, Clark AG, Hartl DL, Ayroles JF (2013) Genetic incompatibilities are widespread within species. *Nature*.
16. Breen MS, Kemena C, Vlasov PK, Notredame C, Kondrashov FA (2012) Epistasis as the primary factor in molecular evolution. *Nature* 490: 535-538.
17. Genetic Analysis of Psoriasis C, the Wellcome Trust Case Control C, Strange A, Capon F, Spencer CC, et al. (2010) A genome-wide association study identifies new psoriasis susceptibility loci and an interaction between HLA-C and ERAP1. *Nat Genet* 42: 985-990.
18. Lincoln MR, Ramagopalan SV, Chao MJ, Herrera BM, Deluca GC, et al. (2009) Epistasis among HLA-DRB1, HLA-DQA1, and HLA-DQB1 loci determines multiple sclerosis susceptibility. *Proc Natl Acad Sci U S A* 106: 7542-7547.
19. Smyth DJ, Plagnol V, Walker NM, Cooper JD, Downes K, et al. (2008) Shared and distinct genetic variants in type 1 diabetes and celiac disease. *The New England Journal of Medicine* 359: 2767-2777.
20. Evans DM, Spencer CC, Pointon JJ, Su Z, Harvey D, et al. (2011) Interaction between ERAP1 and HLA-B27 in ankylosing spondylitis implicates peptide handling in the mechanism for HLA-B27 in disease susceptibility. *Nat Genet* 43: 761-767.
21. Li W, Reich J (2000) A complete enumeration and classification of two-locus disease models.

Hum Hered 50: 334-349.

22. Marchini J, Donnelly P, Cardon LR (2005) Genome-wide strategies for detecting multiple loci that influence complex diseases. *Nat Genet* 37: 413-417.
23. Hill WG, Goddard ME, Visscher PM (2008) Data and theory point to mainly additive genetic variance for complex traits. *PLoS Genet* 4: e1000008.
24. van Heel DA, Hunt K, Greco L, Wijmenga C (2005) Genetics in coeliac disease. *Best Pract Res Clin Gastroenterol* 19: 323-339.
25. Trynka G, Hunt Ka, Bockett Na, Romanos J, Mistry V, et al. (2011) Dense genotyping identifies and localizes multiple common and rare variant association signals in celiac disease. *Nature Genetics* 43: 1193-1201.
26. van Heel DA, Franke L, Hunt KA, Gwilliam R, Zhernakova A, et al. (2007) A genome-wide association study for celiac disease identifies risk variants in the region harboring IL2 and IL21. *Nature Genetics* 39: 827-829.
27. Dubois PCA, Trynka G, Franke L, Hunt Ka, Romanos J, et al. (2010) Multiple common variants for celiac disease influencing immune gene expression. *Nature Genetics* 42: 295-302.
28. Hunt KA, Zhernakova A, Turner G, Heap GA, Franke L, et al. (2008) Newly identified genetic risk variants for celiac disease related to the immune response. *Nat Genet* 40: 395-402.
29. Romanos J, Rosen A, Kumar V, Trynka G, Franke L, et al. (2013) Improving coeliac disease risk prediction by testing non-HLA variants additional to HLA variants. *Gut*.
30. Abraham G, Tye-Din JA, Bhalala OG, Kowalczyk A, Zobel J, et al. (2013) Accurate and robust genomic prediction of celiac disease using statistical learning. *PLoS Genet*.
31. Goudey B, Rawlinson D, Wang Q, Shi F, Ferra H, et al. (2013) GWIS--model-free, fast and exhaustive search for epistatic interactions in case-control GWAS. *BMC Genomics* 14 Suppl 3: S10.
32. Neuman RJ, Rice JP (1992) Two-locus models of disease. *Genet Epidemiol* 9: 347-365.
33. Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, et al. (2008) Genes mirror geography within Europe. *Nature* 456: 98-101.
34. Abraham G, Kowalczyk A, Zobel J, Inouye M (2013) Performance and robustness of penalized and unpenalized methods for genetic prediction of complex human disease. *Genetic Epidemiology* 37: 184-195.
35. Monsuur AJ, de Bakker PI, Zhernakova A, Pinto D, Verduijn W, et al. (2008) Effective detection of human leukocyte antigen risk alleles in celiac disease using tag single nucleotide polymorphisms. *PLoS ONE* 3: e2270.
36. Zheng X, Shen J, Cox C, Wakefield JC, Ehm MG, et al. (2013) HIBAG – HLA Genotype Imputation with Attribute Bagging. *Pharmacogenomics Journal Advance online publication*.
37. Lee SH, Nyholt DR, Macgregor S, Henders AK, Zondervan KT, et al. (2010) A simple and fast two-locus quality control test to detect false positives due to batch effects in genome-wide association studies. *Genet Epidemiol* 34: 854-862.
38. de Bakker PI, McVean G, Sabeti PC, Miretti MM, Green T, et al. (2006) A high-resolution HLA and SNP haplotype map for disease association studies in the extended human MHC. *Nat Genet* 38: 1166-1172.
39. Mackay TF (2014) Epistasis and quantitative traits: using model organisms to study gene-gene interactions. *Nat Rev Genet* 15: 22-33.
40. Cotsapas C, Voight BF, Rossin E, Lage K, Neale BM, et al. (2011) Pervasive sharing of genetic effects in autoimmune disease. *PLoS Genetics* 7: e1002254.
41. Phillips PC (2008) Epistasis--the essential role of gene interactions in the structure and evolution of genetic systems. *Nat Rev Genet* 9: 855-867.
42. Moore JH, Williams SM (2009) Epistasis and its implications for personal genetics. *Am J Hum Genet* 85: 309-320.
43. Szymczak S, Biernacka JM, Cordell HJ, Gonzalez-Recio O, Konig IR, et al. (2009) Machine learning in genome-wide association studies. *Genet Epidemiol* 33 Suppl 1: S51-57.
44. Abraham G, Kowalczyk A, Zobel J, Inouye M (2012) SparSNP: Fast and memory-efficient analysis of all SNPs for phenotype prediction. *BMC Bioinformatics* 13: 88.

45. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, et al. (2011) Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12: 2825-2830.
46. R Core Team (2012) R: A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing.
47. Sing T, Sander O, Beerenwinkel N, Lengauer T (2005) ROCR: visualizing classifier performance in R. *Bioinformatics* 21: 3940-3941.
48. Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, et al. (2011) pROC: an open-source package for R and S+ to analyze and compare ROC curves. *BMC Bioinformatics* 12: 77.
49. Wickham H (2009) ggplot2: elegant graphics for data analysis. New York: Springer.
50. Howie BN, Donnelly P, Marchini J (2009) A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genet* 5: e1000529.
51. Wray NR, Yang J, Goddard ME, Visscher PM (2010) The genetic interpretation of area under the ROC curve in genomic profiling. *PLoS Genetics* 6: e1000864.