

# Bayesian model comparison for rare variant association studies

Guhan Ram Venkataraman<sup>1</sup>, Christopher DeBoever<sup>1</sup>, Yosuke Tanigawa<sup>1</sup>, Matthew Aguirre<sup>1</sup>, Alexander G. Ioannidis<sup>1</sup>, Hakhamanesh Mostafavi<sup>1</sup>, Chris C. A. Spencer<sup>2</sup>, Timothy Poterba<sup>3</sup>, Carlos D. Bustamante<sup>1,4</sup>, Mark J. Daly<sup>3,5</sup>, Matti Pirinen<sup>6,7,8\*</sup>, Manuel A. Rivas<sup>1\*</sup>

1 Department of Biomedical Data Science, Stanford University, Stanford, CA, USA

2 Genomics plc, Oxford, UK

3 Broad Institute of MIT and Harvard, Cambridge, MA, USA

4 Department of Genetics, Stanford University, CA, USA

5 Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA

6 Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland

7 Department of Public Health, University of Helsinki, Helsinki, Finland

8 Department of Mathematics and Statistics, University of Helsinki, Helsinki, Finland

\* [matti.pirinen@helsinki.fi](mailto:matti.pirinen@helsinki.fi)

\* [mrivas@stanford.edu](mailto:mrivas@stanford.edu)

## Abstract

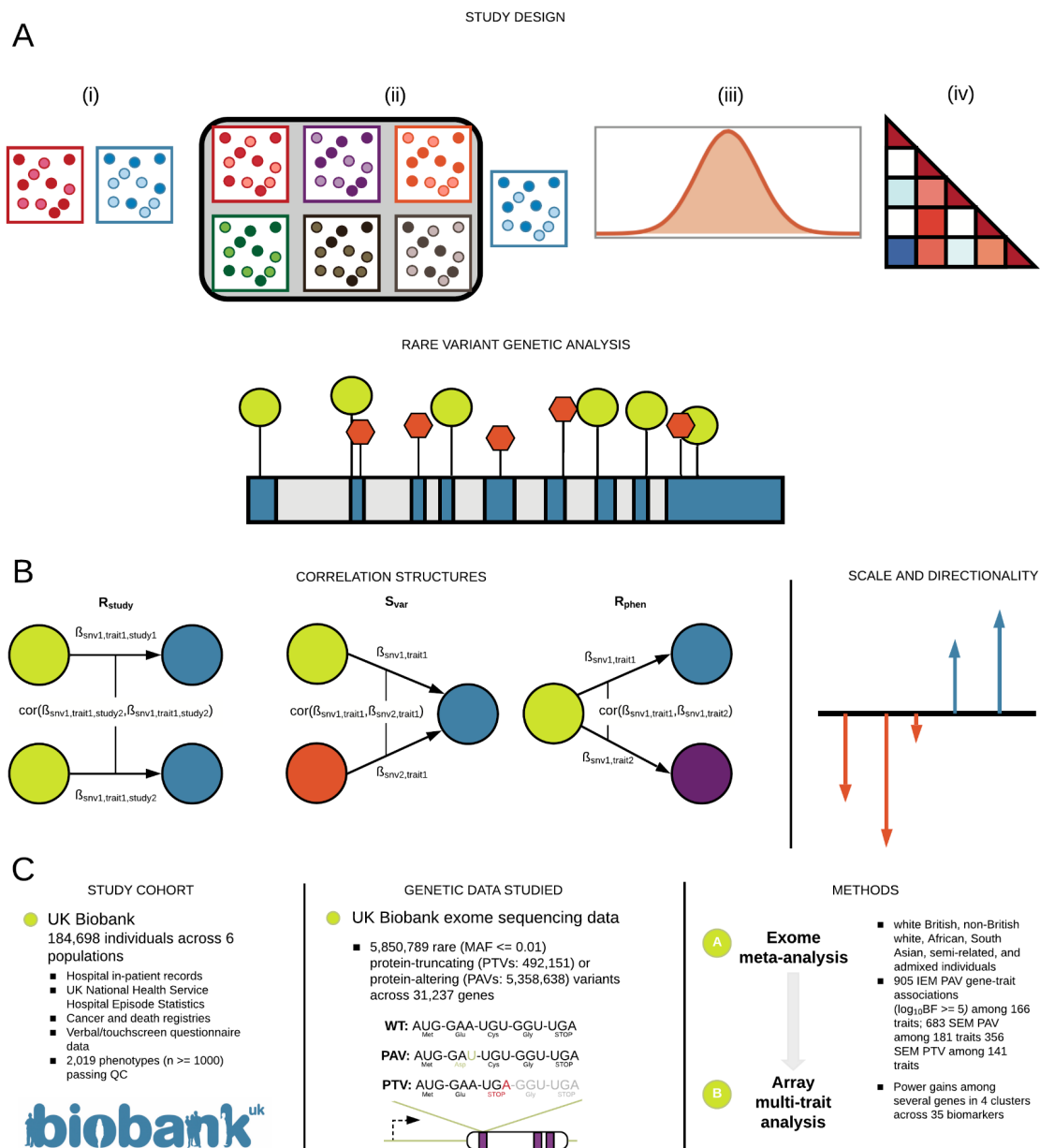
Whole genome sequencing studies applied to large populations or biobanks with extensive phenotyping raise new analytic challenges. The need to consider many variants at a locus or group of genes simultaneously and the potential to study many correlated phenotypes with shared genetic architecture provide opportunities for discovery and inference that are not addressed by the traditional one variant, one phenotype association study. Here, we introduce a Bayesian model comparison approach that we refer to as MRP (Multiple Rare-variants and Phenotypes) for rare-variant association studies that considers correlation, scale, and direction of genetic effects across a group of genetic variants, phenotypes, and studies. The approach requires only summary statistic data. To demonstrate the efficacy of MRP, we apply our method to exome sequencing data (N = 184,698) across 2,019 traits from the UK Biobank, aggregating signals in genes. MRP demonstrates an ability to recover previously-verified signals such as associations between *PCSK9* and LDL cholesterol levels. We additionally find MRP effective in conducting meta-analyses in exome data. Notable non-biomarker findings include associations between *MC1R* and red hair color and skin color, *IL17RA* and monocyte count, *IQGAP2* and mean platelet volume, and *JAK2* and platelet count and crit (mass). Finally, we apply MRP in a multi-phenotype setting; after clustering the 35 biomarker phenotypes based on genetic correlation estimates into four clusters, we find that joint analysis of these phenotypes results in substantial power gains for gene-trait associations, such as in *TNFRSF13B* in one of the clusters containing diabetes and lipid-related traits. Overall, we show that the MRP model comparison approach is able to improve upon useful features from widely-used meta-analysis approaches for rare variant association analyses and prioritize protective modifiers of disease risk.

## Introduction

Sequencing technologies are quickly transforming human genetic studies of complex traits. It is increasingly possible to obtain whole genome sequence data on thousands of samples at manageable costs. As a result, the genome-wide study of rare variants (minor allele frequency [MAF] < 1%) and their contribution to disease susceptibility and phenotype variation is now feasible.<sup>1-4</sup>

In genetic studies of diseases or continuous phenotypes, rare variants are hard to assess individually due to the limited number of observations of each rare variant. Hence, to boost the power to detect a signal, evidence is usually aggregated across variants in blocks. When designing an aggregation method, there are three questions that are usually considered. First, across which biological units should variants be combined (e.g. genes); second, which variants within those units should be included<sup>5</sup>; and third, which statistical model should be used?<sup>6</sup> Given the widespread observations of shared genetic risk factors across distinct diseases, there is also considerable motivation to use gene discovery approaches that leverage the information from multiple phenotypes jointly. In other words, rather than only aggregating variants that may have effects on a single phenotype, we can also bring together sets of phenotypes for which a single variant or set of variants might have effects.

In this paper, we present a Bayesian **M**ultiple **R**are-variants and **P**henotypes (MRP) model comparison approach for identifying rare-variant associations as an alternative to current, widely-used univariate statistical tests. The MRP framework exploits correlation, scale, and/or direction of genetic effects in a broad range of rare-variant association study designs including case-control, multiple diseases and shared controls, a single continuous phenotype, multiple continuous phenotypes or a mixture of case-control and multiple continuous phenotypes (**Figure 1**). MRP makes use of Bayesian model comparison, whereby we compute a Bayes Factor (BF) defined as the ratio of the marginal likelihoods under two models: 1) a null model where all genetic effects are zero; and 2) an alternative model where factors like correlation, scale and direction of genetic effects are considered. For MRP, the BF represents the statistical evidence for a non-zero effect for a particular group of rare variants on the phenotype(s) of interest and can be used as an alternative to p-values from traditional significance testing.



**Figure 1. MRP study overview.** **1A)** MRP is suitable for a broad range of rare variant association study designs, including, from left to right: **i)** case-control, **ii)** multiple diseases with shared controls, **iii)** single quantitative phenotype, and **iv)** mixtures of case-control and quantitative phenotypes. **1B)** Diagram of factors considered in rare variant association analysis including the correlation matrices:  $R_{study}$  (expected correlation of genetic effects among a group of studies),  $S_{var}$  (expected covariance of genetic effects among a group of variants, potentially accounting for annotation of variants), and  $R_{phen}$  (expected correlation of genetic effects among a group of phenotypes). MRP can take into account both scale and direction of effects. **1C)** We focused on 184,698 individuals across 6 ancestry groups in the UK Biobank and analyzed 5,850,789 rare coding variants (492,151 PTVs, 5,358,638 PAVs) in the whole exome sequencing data via single-trait and multi-trait meta-analyses, with a specific focus on 35 biomarker traits.

While many large genetic consortia collect both raw genotype and phenotype data, in practice, sharing of individual genotype and phenotype data across groups is difficult to achieve. To address this, MRP can use summary statistics, such as estimates of effect size and corresponding standard errors from typical single-variant/single-phenotype linear or logistic regressions, as input. Furthermore, we use insights from Liu et al.<sup>7</sup> and Cichonska et al.,<sup>8</sup> which suggest the use of additional summary statistics like covariance estimates across variants and studies, respectively, for the lossless ability to detect gene-based association signals using summary statistics alone.

Aggregation techniques rely on variant annotations to assign variants to groups for analysis. MRP allows for the inclusion of priors on the scale of effect sizes that can be adjusted depending on what type of variants are included in the analysis. For instance, protein truncating variants (PTVs)<sup>9,10</sup> are highly likely to be functional because they often disrupt the normal function of a gene. Additional deleteriousness metrics, such as MPC (which combines subgenic constraints with variant-level data for deleteriousness prediction)<sup>11</sup> and pLI (derived from a comparison of the observed number of PTVs in a sample to the number expected in the absence of fitness effects, i.e. under neutrality, given an estimated mutation rate for the gene)<sup>12</sup>, can further attenuate or accentuate these granular signals. Furthermore, since PTVs typically abolish or severely alter gene function, there is particular interest in identifying protective PTV modifiers of human disease risk that may serve as targets for future therapeutics.<sup>13–15</sup> We therefore demonstrate how the MRP model comparison approach can improve discovery of such protective signals by modeling the direction of genetic effects; this prioritizes variants or genes that are consistent with protecting against disease.

To evaluate the performance of MRP, we use simulations and compare it to other commonly used approaches. Some simple alternatives to MRP include univariate approaches for rare variant association studies including the sequence kernel association test (SKAT)<sup>16</sup>, and the burden test<sup>6</sup>, which are special cases of the MRP model comparison when we assign the prior correlation of genetic effects across different variants to be zero or one, respectively.

We apply MRP to summary statistics computed on a tranche of  $N = 184,698$  exomes for thousands of traits in the UK Biobank for which we have exome data for  $N \geq 1000$  white British individuals, focusing on a meta-analysis context across six UK Biobank subpopulations as defined previously ([Methods](#)).<sup>17</sup> We additionally apply multi-phenotype MRP on clusters of biomarker traits within a single-population context (white British individuals). These analyses show that MRP recovers results from single variant-single phenotype association analyses while increasing the power to detect new rare variant associations, including protective modifiers of disease risk.

# Methods

## Description of MRP

In this section, we provide an overview of the MRP model comparison approach (the [Supplementary Note](#) contains additional details). MRP models GWAS summary statistics as being distributed according to one of two models: the null model, where the effect sizes across all studies for a group of variants and a group of phenotypes is zero, and the alternative model, where effect sizes are distributed according to a multivariate normal distribution with a non-zero mean and/or covariance matrix. MRP compares the evidence between the alternative model and the null model using a Bayes Factor (BF) that is the ratio of the marginal likelihoods under the two models given the observed data.

To define the alternative model, we must specify the prior correlation structure, scale, and direction of the effect sizes. Let  $N$  be the number of individuals and  $K$  the number of phenotype measurements on each individual. Let  $M$  be the number of variants in a testing unit  $G$ , where  $G$  can be, for example, a gene, pathway, or a network. Let  $S$  be the number of studies from which data is obtained — this data may be in the form of **a**) raw genotypes and phenotypes, or **b**) summary statistics including linkage-disequilibrium coefficients, effect sizes, and corresponding standard errors. When considering multiple studies ( $S > 1$ ), multiple rare variants ( $M > 1$ ), and multiple phenotypes ( $K > 1$ ), we define the prior correlation structure of the effect sizes as an  $SMK \times SMK$  matrix,  $U$ . In practice, we define  $U$  as a Kronecker product, of three sub-matrices:

- an  $S \times S$  matrix  $\mathbf{R}_{study}$  containing the correlations of genetic effects among studies that can model the level of heterogeneity in effect sizes between populations<sup>18</sup>;
- an  $M \times M$  matrix  $\mathbf{S}_{var}$  containing the covariances of genetic effects among genetic variants, which may reflect, e.g., the assumption that all the PTVs in a gene may have the same biological consequence<sup>9,10,19</sup> or prior information on scale of the effects obtained through integration of additional functional data<sup>5,20</sup>; by assuming zero correlation of genetic effects, MRP becomes a dispersion test similar to C-alpha<sup>21,22</sup> and SKAT<sup>16</sup>; and
- a  $K \times K$   $\mathbf{R}_{phen}$  matrix containing the correlations of genetic effects among phenotypes, which may be estimated from common variant data.<sup>23–25</sup>

The variance-covariance matrix of the effect size estimates may be obtained from readily available summary statistics such as in-study LD matrices, effect size estimates (or log odds ratios), and the standard errors of the effect size estimates ([Supplementary Note](#)).

MRP allows users to specify priors that reflect knowledge of the variants and phenotypes under study. For instance, we can define an independent effects model (IEM) where the effect sizes of different variants are not correlated at all. In this case,  $\mathbf{S}_{var}$  is the identity matrix, and MRP

behaves similarly to dispersion tests like C-alpha<sup>21,22</sup> and SKAT<sup>16</sup>. We can also define a similar effects model (SEM) by setting every value of  $\mathbf{R}_{var}$  to  $\sim 1$ , where  $\mathbf{R}_{var}$  is the correlation matrix corresponding to covariance matrix  $\mathbf{S}_{var}$ . This model assumes that all variants under consideration have similar effect sizes (with, possibly, differences in scale; like in the burden test). Such a model may be appropriate for PTVs, where each variant completely disrupts the function of the gene leading to a gene knockout. The prior on the scale of effect sizes can be used to denote which variants may have larger effect sizes. For instance, emerging empirical genetic studies have shown that within a gene PTVs may have stronger effects than missense variants.<sup>26</sup> This can be reflected by adjusting the prior variances of effect sizes ( $\sigma$ ) for different categories of variants ([Supplementary Note](#)).

Finally, we can utilize a prior on the expected location and direction of effects to specify alternative models where we seek to identify variants with protective effects against disease. By default, we have assumed that the prior mean, of genetic effects is zero, which makes it possible to analyze a large number of phenotypes without enumerating the prior mean across all phenotypes. To proactively identify genetic variants that are consistent with a protective profile for a disease, we can include a non-zero vector as a prior mean of genetic effects ([Supplementary Note](#)). For this, we can exploit information from Mendelian randomization studies of common variants, such as recent findings where rare protein-truncating loss-of-function variants in *PCSK9* were found to decrease LDL and triglyceride levels and decrease CAD risk<sup>13,27,28</sup> to identify situations where such a prior is warranted.

Applying MRP to variants from a testing unit  $G$  yields a BF for that testing unit that describes the evidence that rare variants in that testing unit have a nonzero effect on the traits used in the model. We can turn this evidence into probability via Bayes rule. Namely, a multiplication of prior-odds of association by BF transforms the prior-odds to posterior-odds. For example, if our prior probability for one particular gene to be associated with a phenotype is  $10^{-4}$ , then an observed BF of  $10^5$  means that our posterior probability of association between the gene and the phenotype is over 90%. Although we see advantages in adopting a Bayesian interpretation for MRP, our approach could also be used in a frequentist context by using BF as a test statistic to compute  $p$ -values ([Supplementary Note](#)).

## UK Biobank Data

### Population definitions

Population	$N_{\text{exome}}$	$N_{\text{array}}$
White British	137,920	337,138
Non-British White	10,432	24,905



African	2,716	6,497
South Asian	3,569	7,885
Semi-Related	18,100	44,632
Admixed	11,961	28,551
Total	184,698	449,608

**Table 1. Number of individuals per population per genotyping platform** (exome/array).

We used a combination of self-reported ancestry (UK Biobank field ID 21000) and principal component analysis to identify six subpopulations in the study: white British, African, South Asian, non-British white, semi-related, and an admixed population. To determine the first four populations, which contain samples not related closer than the 3rd degree, we first used the principal components of the genotyped variants from the UK Biobank and defined thresholds on principal component 1 and principal component 2 and further refined the population definition.<sup>17</sup> Semi-related individuals were grouped as individuals whose genetic data (after passing UK Biobank QC filters; sufficiently low missingness rates; and genetically inferred sex matching reported sex), using a KING relationship table, were between conditional third and conditional second degrees of relatedness to samples in the first four groups. Admixed individuals were grouped as unrelated individuals who were flagged as “used\_in\_pca\_calculation” by the UK Biobank and were not assigned to any of the other populations.

## GWAS Summary Statistics

We performed genome-wide association analysis on 2,019 UK Biobank traits in the six population subgroups as defined above using PLINK v2.00a (20 October 2020). We used the `--glm` Firth fallback option in PLINK to apply an additive-effect model across all sites. Quantitative trait values were rank normalized using the `--pheno-quantile-normalize` flag. We used the following covariates in our analysis: age, sex, array type, and the first ten genetic principal components, where array type is a binary variable that represents whether an individual was genotyped with UK Biobank Axiom Array or UK BiLEVE Axiom Array. For variants that were specific to one array, we did not use array as a covariate.

For the admixed population, we conducted local ancestry-corrected GWAS. We first assembled a reference panel from 1,380 single-ancestry samples in the 1000 Genomes Project,<sup>29</sup> the Human Genome Diversity Project,<sup>30</sup> and the Simons Genome Diversity Project,<sup>31</sup> choosing appropriate ancestry clusters by running ADMIXTURE<sup>32</sup> with the unsupervised setting. Using cross-validation, eight well-supported ancestral population clusters were identified: African, African Hunter-Gatherer, East Asian, European, Native American, Oceanian, South Asian, and West Asian. We then used RFMix v2.03<sup>33</sup> to assign each of the 20,727 windows across the phased genomes to one of these eight ancestry clusters (for all individuals in the UK Biobank).

These local ancestry assignments were subsequently used with PLINK2 as local covariates in the GWAS for the admixed individuals for SNPs within those respective windows. PLINK2 allows for the direct input of the RFMix output (the MSP file, which contains the most likely subpopulation assignment per conditional random field [CRF] point) as local covariates using the “local-cov”, “local-psam”, and “local-haps” flags, the “local-cats0=n” flag (where n is the number of assignments), and the “local-pos-cols=2,1,2,7” flag (for a typical RFMix MSP output file - see <https://www.cog-genomics.org/plink/2.0/assoc>).

## Variant Quality Control and Metadata Generation

For quality control (QC), In total, we ensured that variant-level missingness was less than 10%, that the  $p$ -value for the Hardy-Weinberg equilibrium test (computed within unrelated individuals of white British ancestry) was greater than  $10^{-15}$ , and that the variant was uniquely represented (the “CHROM:POS:REF:ALT” variant string was uniquely identified) in the PLINK dataset file. In total, we removed 195,920 variants that failed to meet all of these criteria, except for 134 variants on the Y chromosome.

For the remainder, we used Variant Effect Predictor (VEP)<sup>34</sup> to annotate the most severe consequence, the gene symbol, and HGVS<sub>p</sub> of each variant in the UK Biobank exome and array data. We calculated minor allele frequencies using PLINK. MPC<sup>11</sup> values (variant-level) and pLI gene memberships<sup>12</sup> were annotated from source. To determine LD independence criteria, we used PLINK’s --indep-pairwise function with a window size of 1000kb, a step size of 1, and an  $r^2$  threshold of 0.1 on those variants that pass QC. As our analyses focused on PTVs and PAVs, we then performed this same LD independence analysis on only these, overriding assignments in the first analysis if necessary. We provide these essential metadata, which are necessary for MRP, in exome

([https://biobankengine.stanford.edu/static/ukb\\_exm\\_oqfe-consequence\\_wb\\_maf\\_gene\\_ld\\_indep\\_mpc\\_pli.tsv.gz](https://biobankengine.stanford.edu/static/ukb_exm_oqfe-consequence_wb_maf_gene_ld_indep_mpc_pli.tsv.gz)) and array

([https://biobankengine.stanford.edu/static/ukb\\_cal-consequence\\_wb\\_maf\\_gene\\_ld\\_indep\\_mpc\\_pli.tsv.gz](https://biobankengine.stanford.edu/static/ukb_cal-consequence_wb_maf_gene_ld_indep_mpc_pli.tsv.gz)) tables, available for direct download via the Global Biobank Engine.<sup>35</sup>

## Applications

For exome applications, we chose variants with  $MAF \leq 1\%$  and that were LD-independent according to the criteria mentioned above. For quantitative traits, we removed variants whose regression effect size had standard error greater than 100, and for binary traits, we removed variants whose regression effect size had standard error greater than 0.2. For array applications, we chose variants with  $MAF \leq 1\%$  and removed variants whose regression effect size had standard error greater than 0.2. While MRP is capable of handling all variant types (e.g. proximal coding and intronic variants), we included only protein-altering variants (PAVs) and protein-truncating variants (PTVs) in both exome and array analyses (exome data features many more PAVs and thus potential for power gain; **Table S1**; **Figure S2**). These sets respectively contain the following consequence annotations:

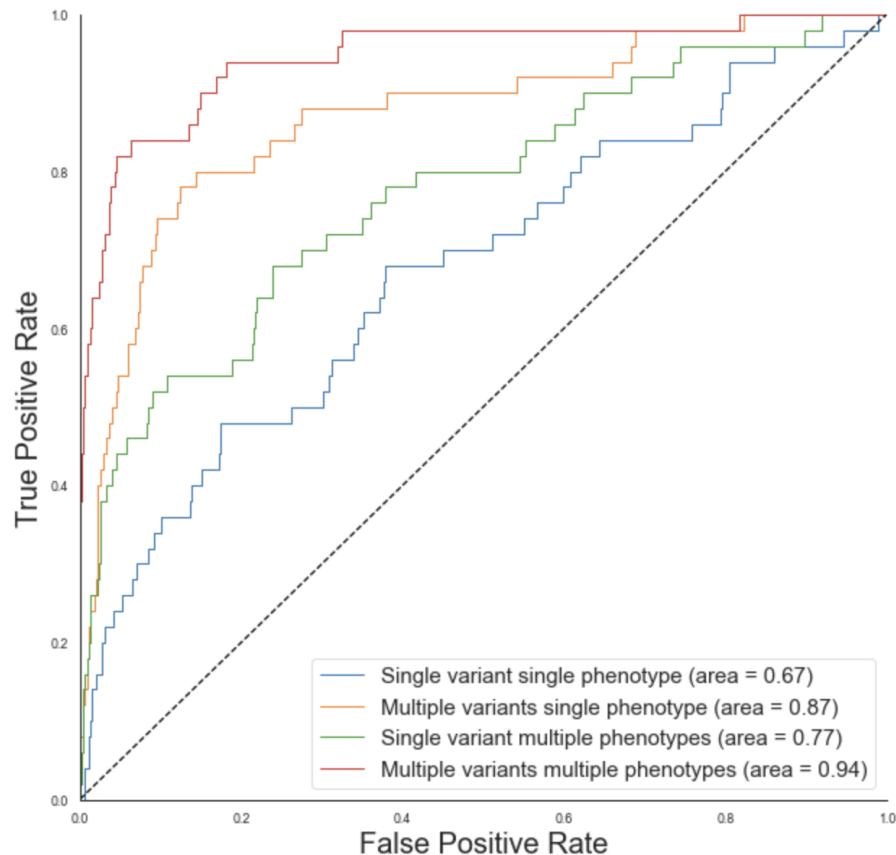


- PAVs: protein\_altering\_variant, inframe\_deletion, inframe\_insertion, splice\_region\_variant, start\_retained\_variant, stop\_retained\_variant, missense\_variant
- PTVs: frameshift\_variant, splice\_acceptor\_variant, splice\_donor\_variant, stop\_gained, start\_lost, stop\_lost

For both quantitative and binary traits, PTVs were assigned a  $\sigma$  (standard deviation of prior on effect size) of 0.2, whereas PAVs were assigned a  $\sigma$  value of 0.05. We also incorporated MPC and pLI deleteriousness metrics into our exome analyses. For those PTVs with a pLI of  $> 0.8$ , we increased  $\sigma$  to 0.5, and for those PAVs with an MPC  $\geq 1$ , we set  $\sigma = 0.05 \times \text{MPC}$ . These adjustments serve to further granularize and weight MRP results in biologically meaningful ways ([Table S3](#); [Figure S4](#)). For the exome meta-analysis, we assumed a similar effects model across studies and an independent effects model across variants.

We also studied how the application of MRP to multiple phenotypes together would potentially boost power to detect rare-variant associations. We calculated pairwise genetic correlation between 35 biomarker phenotypes<sup>17</sup> using LD score regression,<sup>25</sup> and then used the hclust algorithm in the R stats package<sup>36</sup> to generate phenotype clusters. For each of these clusters, using the array data, we performed MRP in the multi-phenotype setting.

## Results



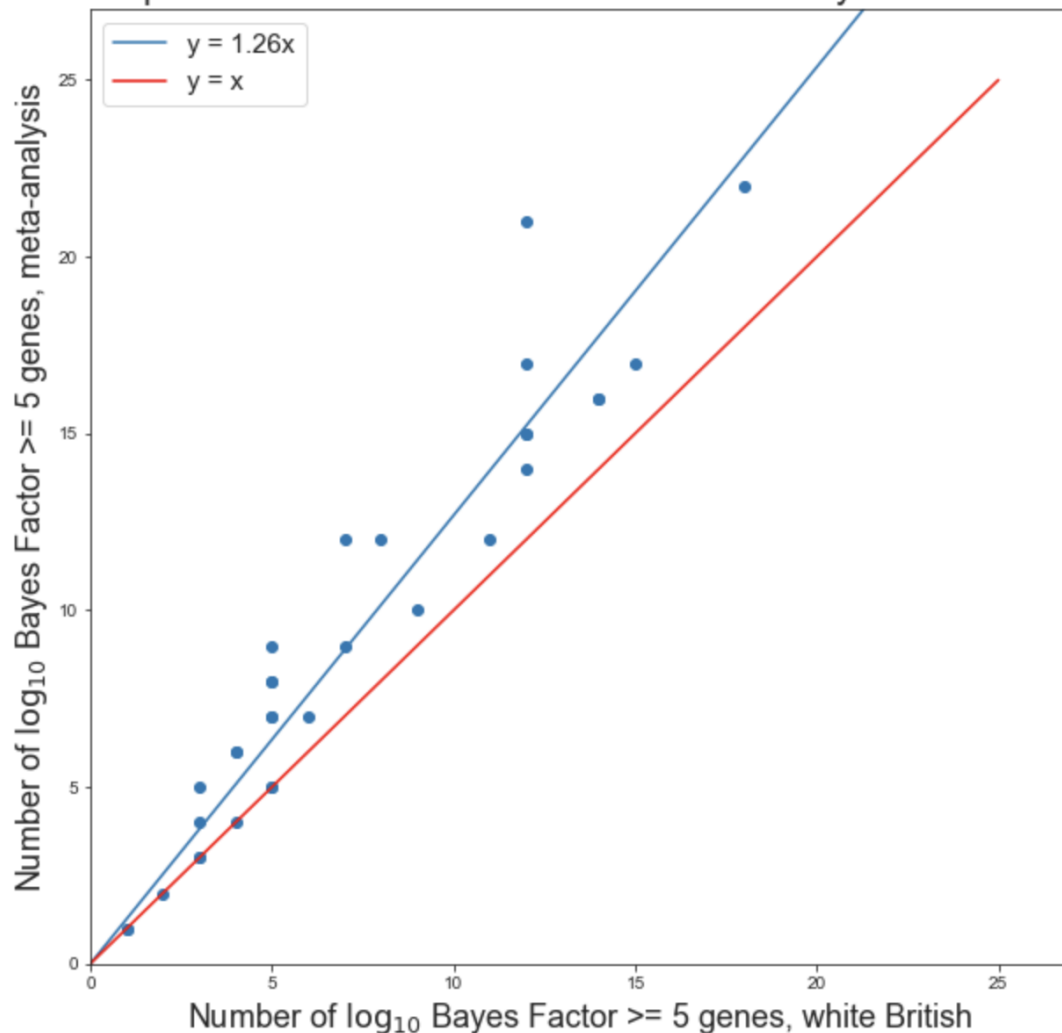
**Figure 2A. From single-variant and single-phenotype to multiple-variant and multiple-phenotype gene discovery.** ROC curves for detecting simulated gene association to any of the phenotypes using single variant/single phenotype association (blue) to multiple-variant and multiple-phenotype association (red).

## Simulations

To study the behavior of MRP going from a single phenotype to multiple phenotypes, we conducted a simulation study where we assumed an allelic architecture consistent to that discovered for *APOC3* in relation to triglycerides, low-density lipoprotein cholesterol (LDL-C), and high-density lipoprotein cholesterol (HDL-C).<sup>37–39</sup> We simulated three continuous phenotypes with a total correlation consistent with that observed for triglycerides, LDL-C, and HDL-C. Furthermore, we introduced effects to four variants consistent with the effects observed in four PTVs (approximately 0.35 standard deviations away from the population mean) and to another four variants consistent with the effects observed for missense variants (approximately 0.2 standard deviations away from the population mean) all with minor allele frequency of

0.05%. The PTV group of variants had the same effects, whereas out of the missense variants, half had positive and the other half had negative effect sizes. The correlation of effects between the group of phenotypes was set to be directionally consistent with the direction of genetic effects observed for lipid phenotypes and PTVs in *APOC3*, i.e. proportional effects for triglycerides and LDL-C, and inversely proportional for LDL-C and HDL-C, and triglycerides and HDL-C. We simulated 1000 genes where 50 of the genes contained non-zero effects on the multivariate phenotype. Given we know which of the 1000 genes contained non-zero effects, we could compute the true positive rates and false positive rates for a given BF threshold. We analyze the data as follows: i) single-variant and single-phenotype, ii) multiple variants and single-phenotype, iii) single-variant and multiple-phenotypes, and iv) multiple-variants and multiple-phenotypes (Figure 2A). We find that in some scenarios, analyzing multiple-variants and multiple-phenotypes jointly improved the ability to detect signals.

### Power comparison between white British and meta-analysis across biomarkers



**Figure 2B. From single to multiple populations.** Scatterplot showing number of genes with  $\log_{10} \text{BF} \geq 5$  for white British population only (x-axis) versus meta-analysis (y-axis) across 35 biomarkers. Assuming that BFs are correctly calibrated in both analyses and that meta-analysis is not inflated compared to white

British-only MRP, suggests a ~26% increase in power when incorporating summary statistics across multiple populations.

## Exome single-phenotype meta-analyses

MRP was used to perform exome meta-analysis on 2019 traits across six UK Biobank populations as described in [Methods](#). Among the best powered and represented traits were a set of 35 biomarkers, the focus of a previous publication<sup>17</sup>. We see the number of  $\log_{10} BF \geq 5$  genes increasing from a single-population to a meta-analysis setting. Since we expect that the meta-analysis over different ancestries cannot be more confounded than the analysis of a single ancestry, we interpret the increase in the number of genes as an increase in the statistical power to detect rare-variant associations ([Figure 2B](#)).

We categorize these biomarkers into six categories as in Sinnott-Armstrong et. al.<sup>17</sup> (Bone and Joint, Cardiovascular, Diabetes, Hormone, Liver, and Renal - [Figure 3](#)) and we recover several known gene-trait associations and discover several others.

Among the “Bone and Joint” biomarkers (alkaline phosphatase, calcium, and vitamin D), we recover associations between *CASR* and calcium<sup>40</sup> and *HAL* and vitamin D<sup>41</sup>. As compared to results from array data as found in Sinnott-Armstrong, et. al.<sup>17</sup>, we also recover exome-specific associations between *ALDH5A1* and alkaline phosphatase<sup>42</sup> and *PDE3B* and vitamin D<sup>41</sup>.

For the “Cardiovascular” phenotypes (apolipoprotein A, apolipoprotein B, C-reactive protein, total cholesterol, HDL cholesterol, LDL cholesterol, lipoprotein A, and triglycerides), MRP recovers array associations between: *PLG*, *LPA*, and lipoprotein A<sup>43</sup>; *APOC3* and triglycerides<sup>44</sup>; *ANGPTL3* and triglycerides<sup>45</sup>; *APOB* and apolipoprotein B<sup>46</sup> and LDL cholesterol<sup>47</sup>; *ABCA1* and apolipoprotein A<sup>44</sup> and HDL cholesterol<sup>47</sup>; *PCSK9* and total cholesterol<sup>48</sup>; and *CRP* and C-reactive protein<sup>49</sup>. Exome-only signals recover associations such as between *ZPR1*<sup>50</sup> and *SIK3*<sup>44</sup> and triglycerides.

In the two diabetes-related phenotypes (glucose and HbA1c), we recover associations between *G6PC2* and glucose<sup>48</sup> as well as *PIEZO1* and HbA1c<sup>51</sup> and an additional exome association between *G6PD* and HbA1c<sup>51</sup>. Hormonal recoveries include those between *SHBG* and SHBG and testosterone levels and *GH1* and IGF-1 levels.

MRP applied to liver-related phenotypes recover known associations between: *UGT* genes and bilirubins<sup>52</sup>; *GOT1* and aspartate aminotransferase<sup>53</sup>; *FCGRT* and albumin<sup>42</sup>; and *GPT* and AST-ALT ratio<sup>49</sup>. In the exome sequencing, we additionally recover associations between *GGT1* and gamma glutamyltransferase<sup>54</sup>, *TMEM236* and aspartate aminotransferase<sup>42</sup>, and *SLCO1B3* and bilirubin.<sup>55</sup>

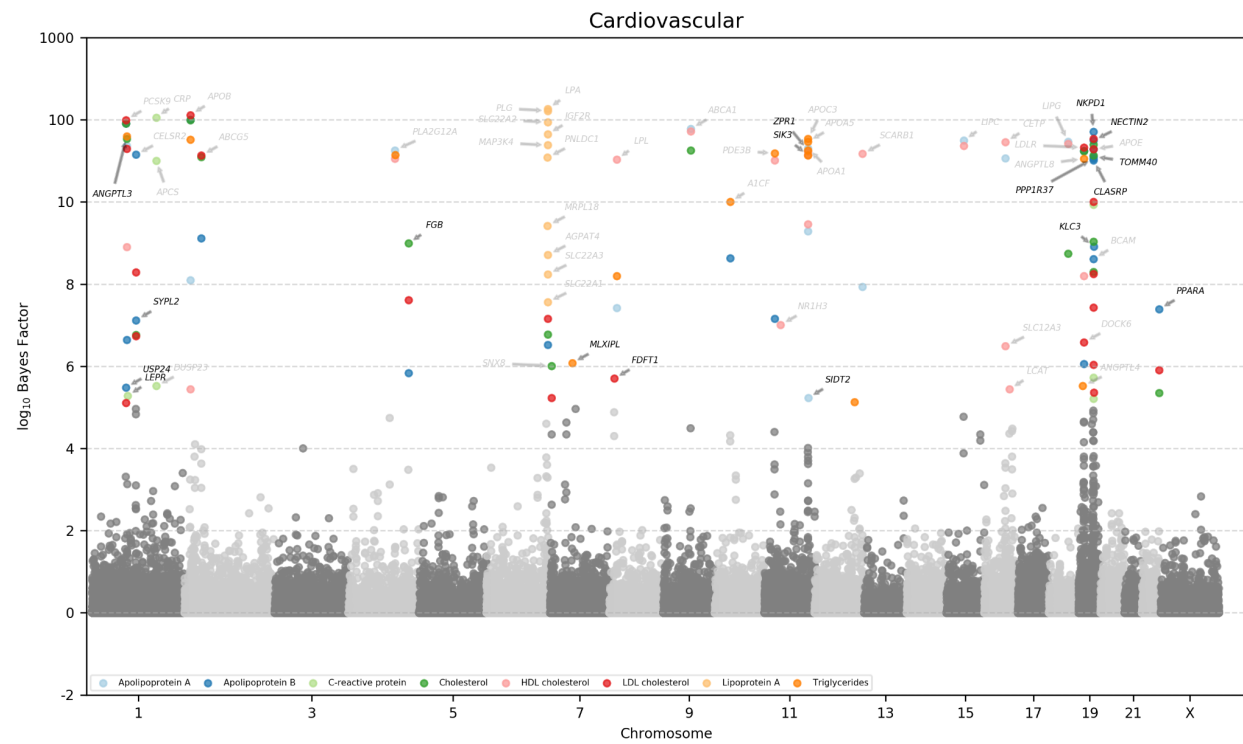
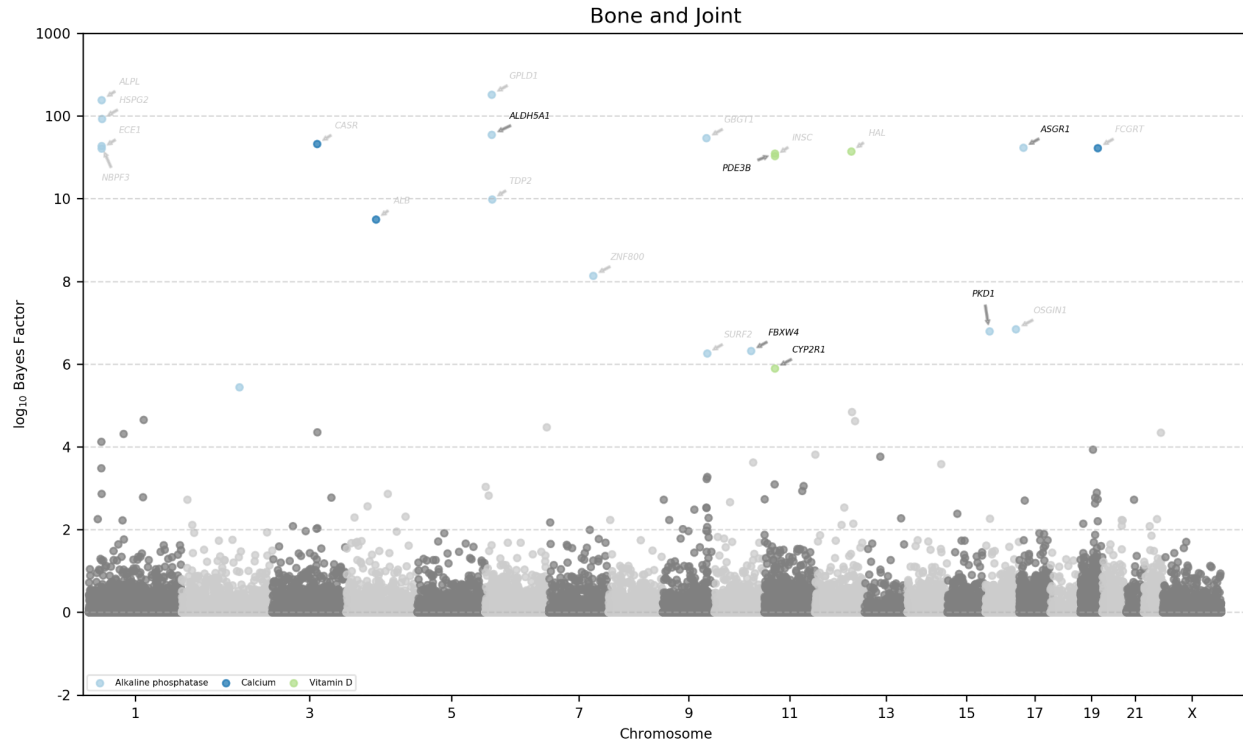
The renal traits similarly feature a mix of array recoveries and exome discoveries. We recover signal between: *SLC22A2* and creatinine<sup>56</sup>; *CST3* and Cystatin C<sup>46</sup>; *COL4A4* and

microalbumin<sup>57</sup>; *TNFRSF13B* and non-albumin protein<sup>42</sup>; *FCGRT* and total protein; *WDR1*, *RASGRP2*, *DRD5*, and urate<sup>58,59</sup>; and *LRP2* and eGFR levels<sup>60</sup>. We additionally discover novel gene-trait associations (not found in the NHGRI-EBI catalog or Open Targets Genetics) across these biomarker categories, including: *GLPD1* and alkaline phosphatase; *NKPD1* and apolipoprotein B; *RENBP/MAP3K15* and Hba1c; *PARPBP* and IGF-1; *NLGN2* and SHBG; *ALB* and albumin; *ALPL* and phosphate; *RBM47* and urea; *ALDH16A1* and urate; *THBD* and Cystatin C; *ITPR3* and phosphate; *SLC22A7* and creatinine; and *FCGR2B* and non-albumin protein.

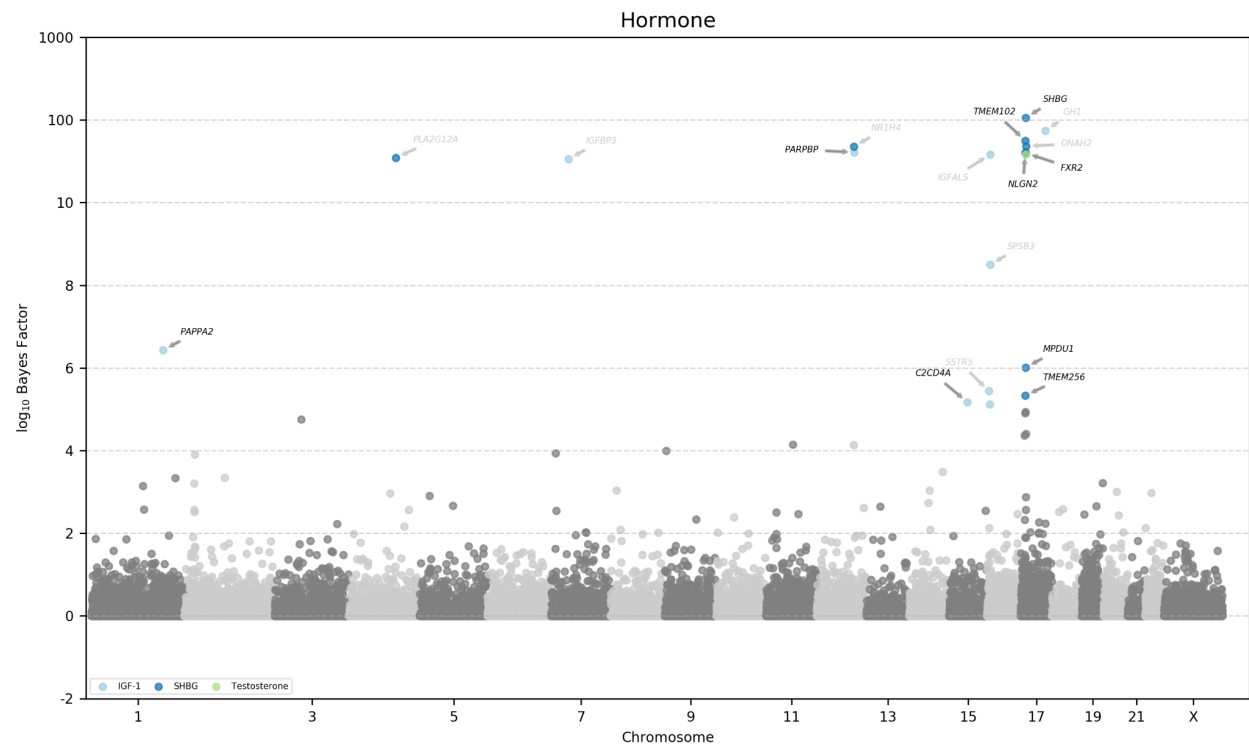
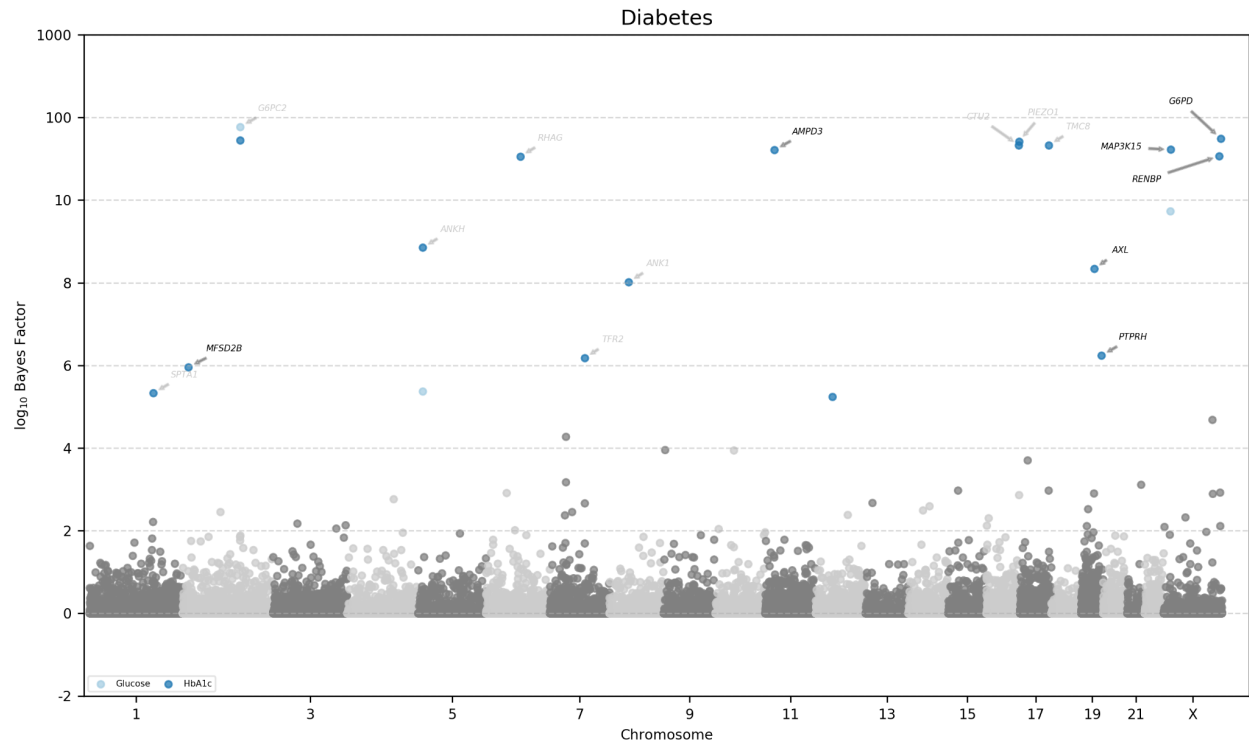
For the 2,019 traits for which MRP was performed, there were also a considerable number of associations found amongst non-biomarker traits. We found associations between: *TUBB1* and platelet distribution width and mean platelet volume<sup>61</sup>; *IL17RA* and monocyte count and percentage<sup>61</sup>; *OCA2/MC1R* and skin color/hair color<sup>62-65</sup>; *IQGAP2* and mean platelet volume<sup>61</sup>; *SLC24A5*, *HERC2*, *TCF25*, *TYR* and skin color<sup>66</sup>; *SH2B3*, *JAK2* and platelet crit<sup>67</sup> and count<sup>68</sup>; *KALRN* and mean platelet volume<sup>61</sup>; *HBB* and mean corpuscular volume<sup>69</sup>, mean corpuscular hemoglobin<sup>70</sup>, and red blood cell count<sup>71</sup>; and *CXCR2* and neutrophil count<sup>61</sup>.

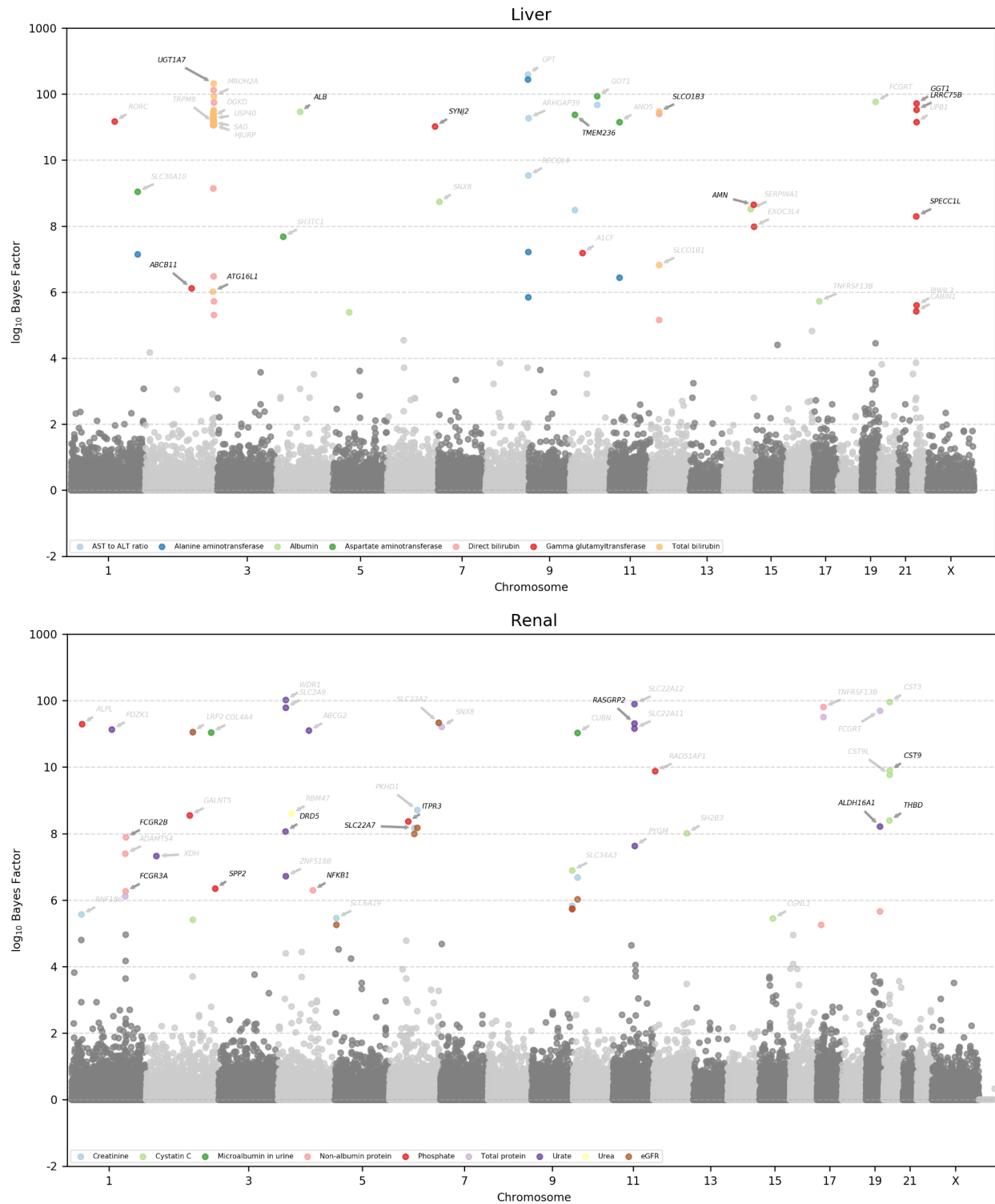
We have published the full set of associations ( $\log_{10} \text{BF} \geq 5$ ) from an independent effects model amongst PAVs, from a similar effects model amongst PAVs, as well as from a similar effects model amongst PTVs on the Global Biobank Engine for exomes ([https://biobankengine.stanford.edu/RIVAS\\_HG38/mrpgene/all](https://biobankengine.stanford.edu/RIVAS_HG38/mrpgene/all)) and array data ([https://biobankengine.stanford.edu/RIVAS\\_HG19/mrpgene/all](https://biobankengine.stanford.edu/RIVAS_HG19/mrpgene/all)).

MRP was implemented using Python (dependencies: pandas v1.1.5, numpy v1.16.4, rpy2 v3.0.4, scipy v1.3.0). The requirements, code, and metadata files can be found at <https://github.com/rivas-lab/mrp>.



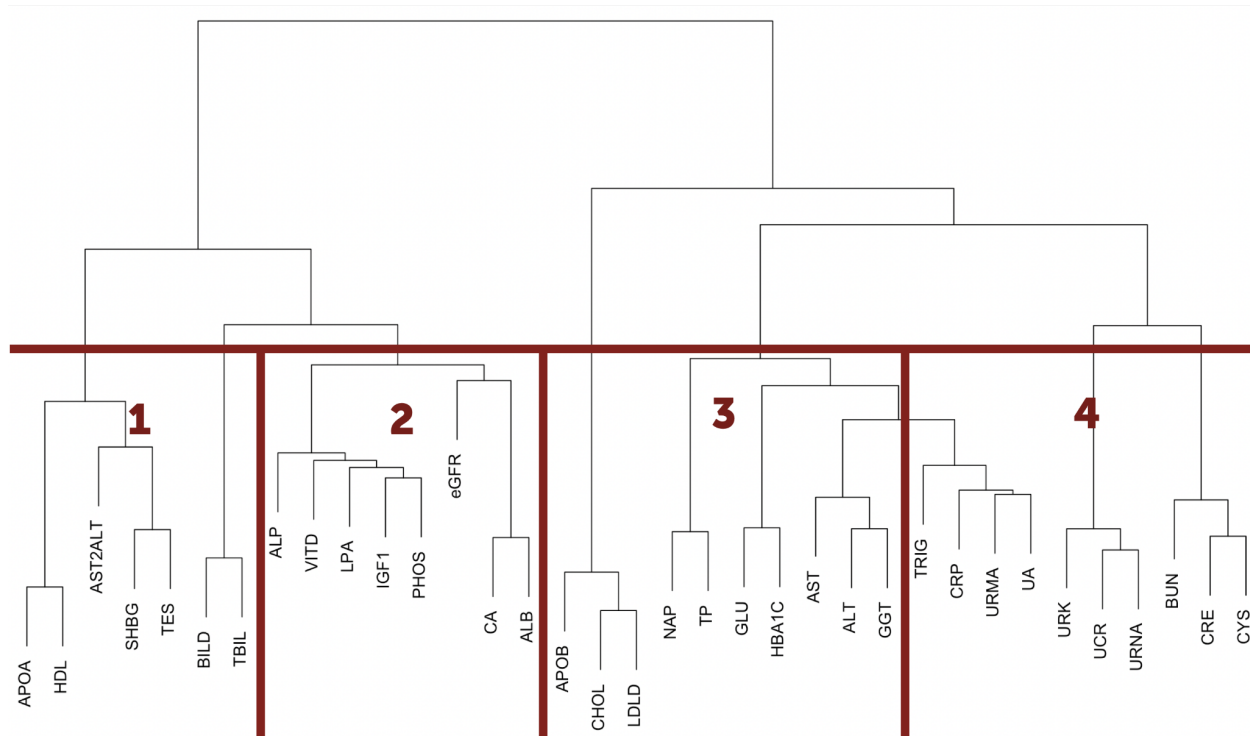






**Figure 3. Manhattan plots showing  $\log_{10}$  BF under an independent effects variant model amongst protein-altering variants for 6 categories across 35 biomarkers.** Scale is logarithmic after  $\log_{10}$  BF  $\geq$  10. Genes found in Sinnott-Armstrong, et.al.<sup>17</sup> are annotated in grey, whereas the other genes are annotated in black.

## Array single-population multi-phenotype analyses

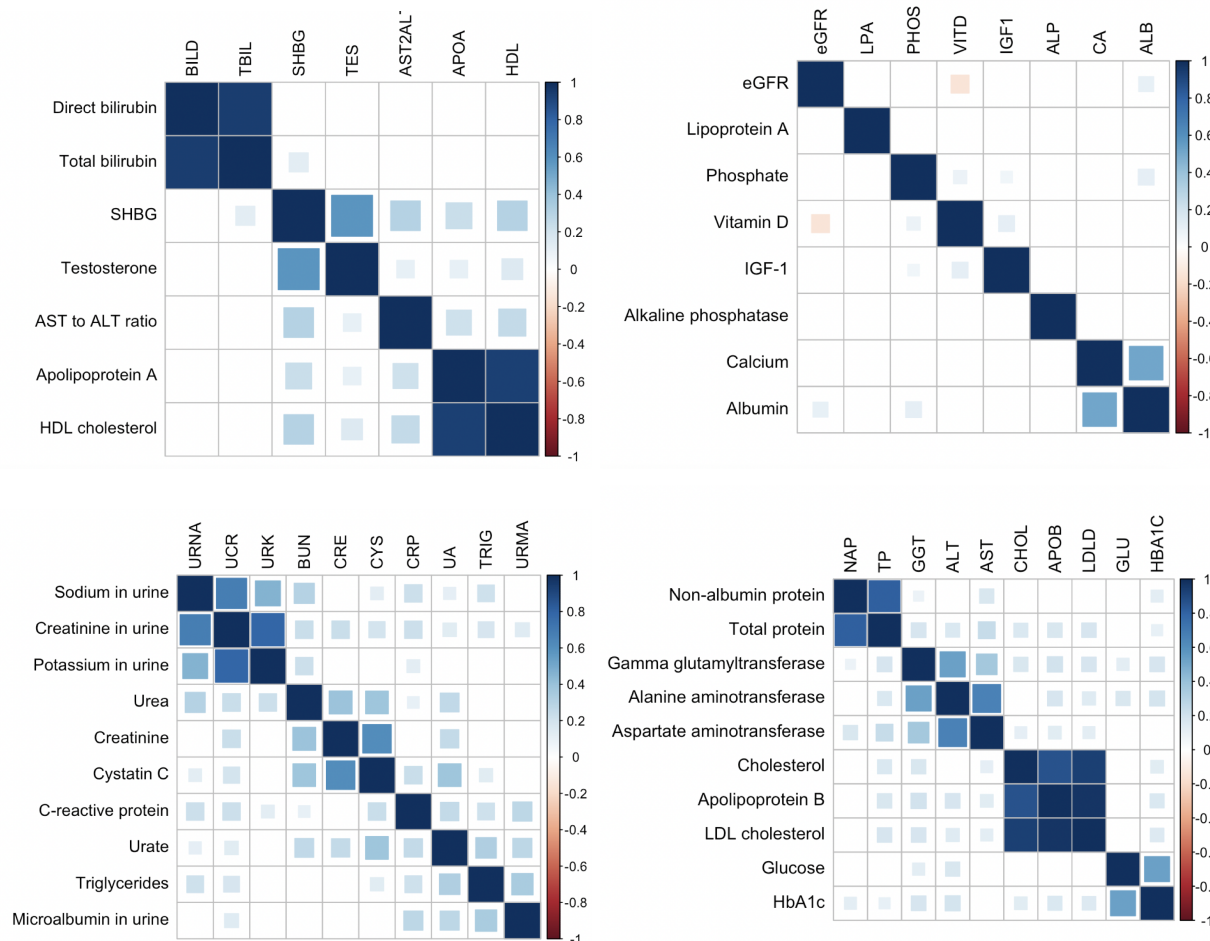


**Figure 4A. Hierarchical clustering dendrogram.** Based on genetic correlation derived from an LD-score regression-based distance matrix between 35 biomarker traits.

In order to demonstrate the effectiveness of MRP to boost signal in a multi-phenotype context, we used LD-score regression<sup>25</sup> to determine genetic correlations between the 35 biomarker traits (Figure S5) that were a focus of a previous paper<sup>17</sup>. This correlation matrix was then used for hierarchical clustering followed by dynamic tree cutting, which formed four clusters of between seven and ten traits each (Figure 4A). We generated the correlation plots as shown in Figure 4B.

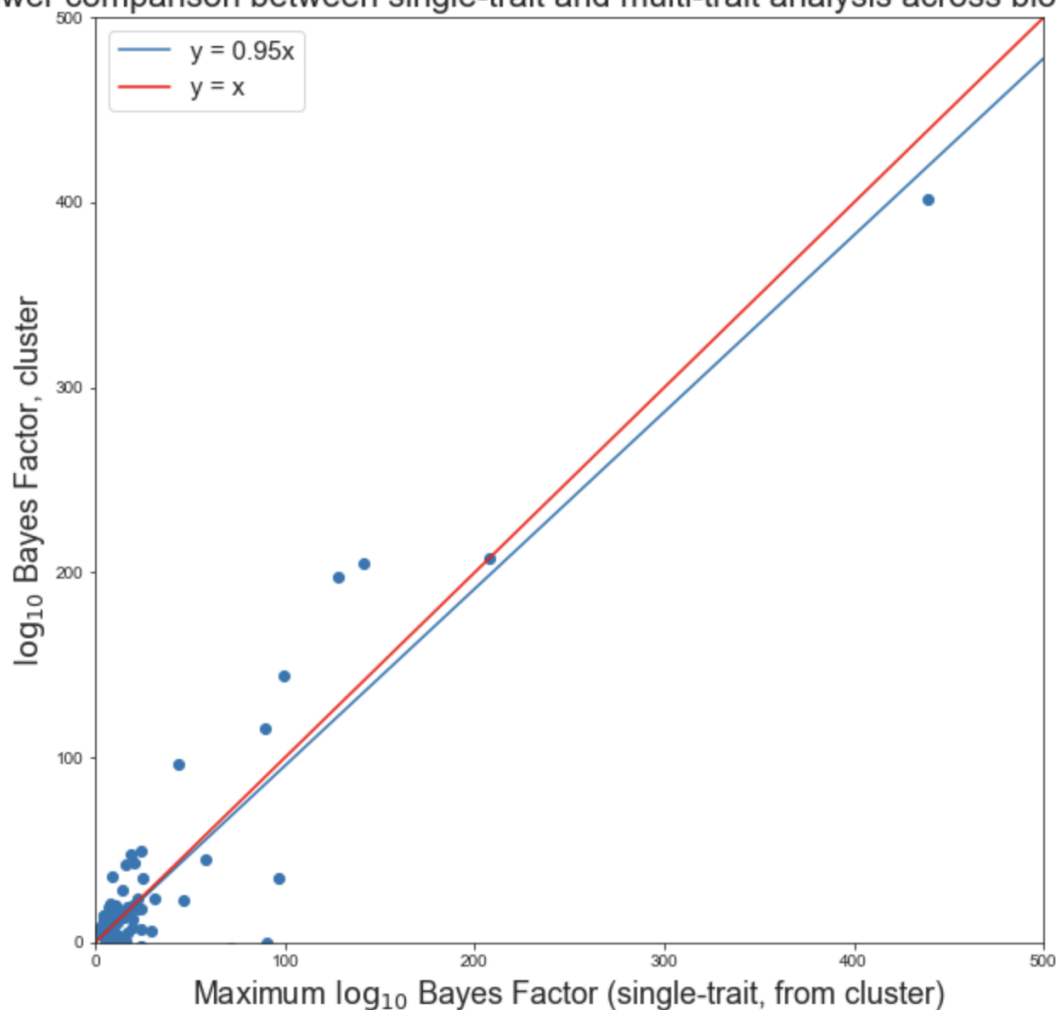
Multi-phenotype MRP results in several substantial power gains throughout the four clusters; one of these clusters is highlighted in Figure 4C. As compared to the maximum  $\log_{10}$  BF from the constituent phenotypes, the multi-phenotype analysis generally fares comparably, while also highlighting clear targets. We found evidence for association between rare coding variants in several genes and the clusters above; *TNFRSF13B* ( $\log_{10}$  BF<sub>multi-trait</sub> = 204.5, max[ $\log_{10}$  BF<sub>single-trait</sub>] = 141.0), *APOB* ( $\log_{10}$  BF<sub>multi-trait</sub> = 197.9, max[ $\log_{10}$  BF<sub>single-trait</sub>] = 128.0), *SNX8* ( $\log_{10}$  BF<sub>multi-trait</sub> = 96.0, max[ $\log_{10}$  BF<sub>single-trait</sub>] = 43.8) receive a boost in  $\log_{10}$  BF of over 50 units for cluster 1 (Alanine aminotransferase, Aspartate aminotransferase, Gamma glutamyltransferase, Glucose, HbA1c, Total protein, Apolipoprotein B, Cholesterol, LDL cholesterol, and Non-albumin protein). Several other genes that are clearly below 5 (in  $\log_{10}$  BF) in the single-trait settings become

above 5 in the joint setting (e.g., *G6PC*;  $\log_{10} \text{BF}_{\text{multi-trait}} = 5.3$ ,  $\max[\log_{10} \text{BF}_{\text{single-trait}}] = 1.3$ ). The *G6PC* gene provides instructions for making the glucose 6-phosphatase enzyme, found on the membrane of the endoplasmic reticulum. The enzyme is expressed in active form in the liver, kidneys, and intestines, and is the main regulator of glucose production in the liver; given the traits included in cluster 1, the increase in power may be biologically relevant<sup>72</sup>. These results demonstrate that MRP can identify biologically meaningful targets that may be missed by standard GWAS approaches.



**Figure 4B. LD-score regression-based genetic correlation plots of candidate clusters.** Derived from the dendrogram in [Figure 4A](#) using a dynamic tree cutting algorithm.

### Power comparison between single-trait and multi-trait analysis across biomarkers



**Figure 4C. Cluster vs. single-trait power analysis.** Power comparison of genes with  $\log_{10} \text{BF} \geq 5$  in either i) any of the single-trait analyses of the traits within the cluster or ii) the multi-trait analysis, for a cluster of biomarkers (Alanine aminotransferase, Aspartate aminotransferase, Gamma glutamyltransferase, Glucose, HbA1c, Total protein, Apolipoprotein B, Cholesterol, LDL cholesterol, and Non-albumin protein). x-axis depicts the maximum  $\log_{10} \text{BF}$  of the gene amongst any of the constituent single-trait analyses, and y-axis depicts the multi-trait result. Multi-trait analyses roughly equal the highest-powered single-trait analyses, while also substantially boosting signal in some genes.

## Discussion

In this study, we developed MRP, a Bayesian model comparison approach that shares information across variants, phenotypes, and studies to identify rare variant associations. We used simulations to verify that jointly considering both variants and phenotypes can improve the ability to detect associations. We also applied the MRP model comparison framework in a meta-analysis setting to exome summary statistics across the UK Biobank, identifying strong evidence for the previously described associations between, for example, *HAL* and vitamin D<sup>41</sup>,

and discovering several novel associations, such as between *GLPD1* and alkaline phosphatase. We made the full results set available on the Global Biobank Engine (<https://biobankengine.stanford.edu/>)<sup>35</sup>. We also leveraged MRP to boost signal in a multi-phenotype setting using the array data (which has many more samples than the exome data), finding genes such as *G6PC* that do not come up in the single-trait context but show strong evidence in the joint analysis. These results demonstrate the ability of the MRP model comparison approach to leverage information across multiple phenotypes and variants to discover rare variant associations.

As genetic data linked to high-dimensional phenotype data is increasingly being made available through biobanks, health systems, and research programs, there is a large need for statistical approaches that can leverage information across different genetic variants, phenotypes, and studies to make strong inferences about disease-associated genes. The approach presented here relies only on summary statistics from marginal association analyses, which can be shared with less privacy concerns compared to raw genotype and phenotype data. Combining joint analysis of variants and phenotypes with meta-analysis across studies offers new opportunities to identify gene-disease associations.

## Author Contributions

M.A.R. and M.P. designed the method and derived all analytical calculations. G.R.V. implemented the method and analyzed the data. G.R.V., M.A.R., M.P., and C.D. wrote the manuscript. G.R.V., M.A.R., M.P., C.C.A.S., C.D., Y.T., M.A., T.P., and H.M. provided quality control analysis, figure edits, and revisions to the manuscript. A.G.I. also provided revisions and feedback on local ancestry-corrected GWAS. C.D.B. and M.J.D. provided critical feedback on methodology.

## Acknowledgements and Funding

This research was conducted using the UK Biobank Resource under application number 24983, “Generating effective therapeutic hypotheses from genomic and hospital linkage data” (<http://www.ukbiobank.ac.uk/wp-content/uploads/2017/06/24983-Dr-Manuel-Rivas.pdf>). Based on the information provided in protocol 44532, the Stanford IRB has determined that the research does not involve human subjects as defined in 45 CFR 46.102(f) or 21 CFR 50.3(g). All participants in the UK Biobank study provided written informed consent (more information is available at <https://www.ukbiobank.ac.uk/2018/02/gdpr/>). Statin adjustment analyses were further conducted via UK Biobank application 7089 using a protocol approved by the Partners HealthCare Institutional Review Board. We thank all the participants in the UK Biobank. We thank members of the Rivas lab for their feedback. M.A.R. is in part supported by the NHGRI of the NIH under award R01HG010140 (M.A.R.) and an NIH Center for Multi- and Trans-ethnic Mapping of Mendelian and Complex Diseases grant (5U01 HG009080). G.R.V. is supported by the National Library of Medicine (NLM) T15 Continuing Education Training Grant. The content is

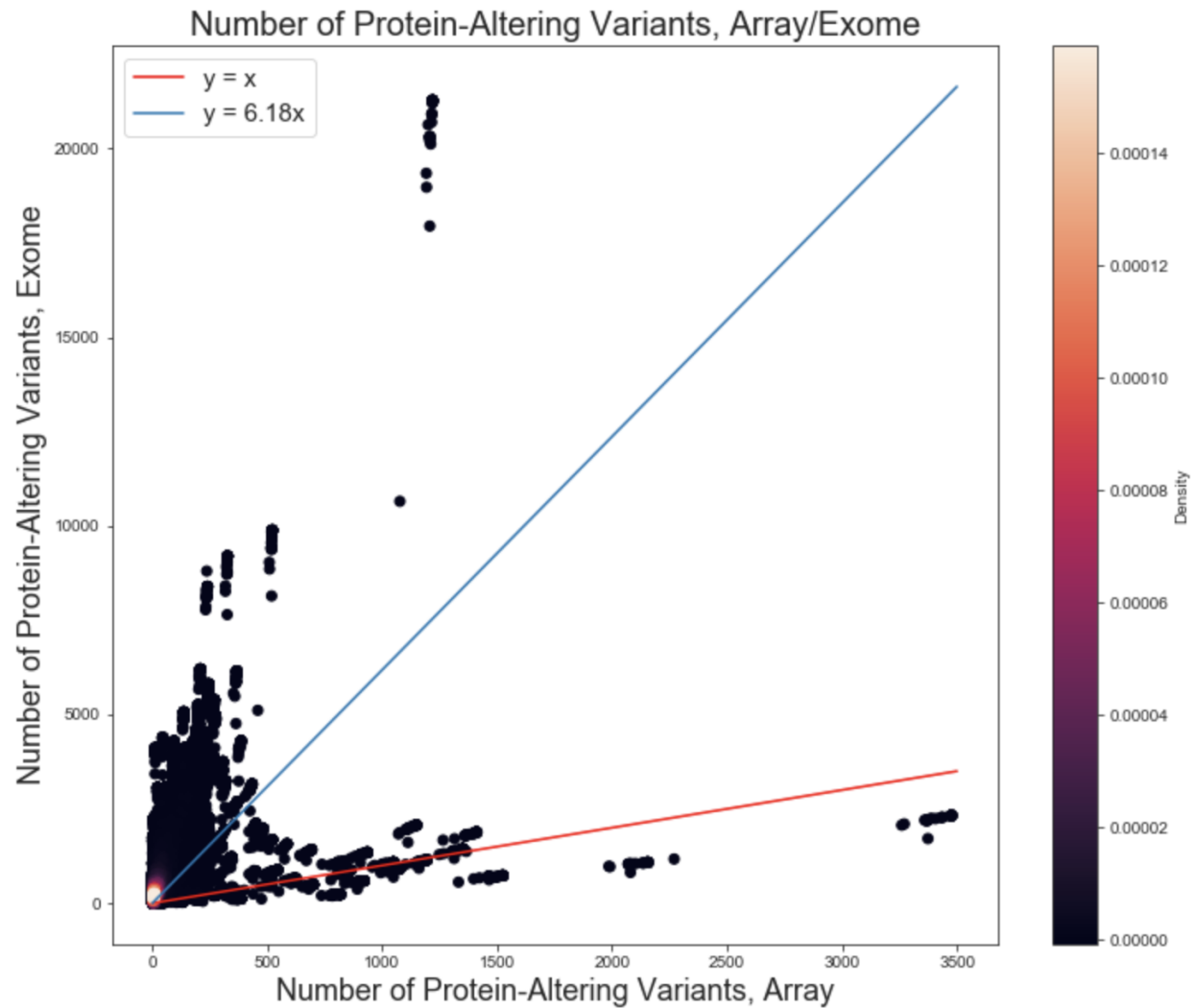


solely the responsibility of the authors and does not necessarily represent the official views of the NIH. Some of the computing for this project was performed on the Sherlock cluster at Stanford University. We would like to thank Stanford University and the Stanford Research Computing Center for providing computational resources and support that contributed to these research results.

## Supplementary Materials

Trait	gene	Number of PAVs, array	log <sub>10</sub> BF, array	Number of PAVs, exome	log <sub>10</sub> BF, exome	log <sub>10</sub> BF Difference
Total bilirubin	<i>UGT1A7</i>	5	1.2	247	213	211.8
Direct bilirubin	<i>UGT1A7</i>	5	0.6	228	133	132.4
Lipoprotein A	<i>PLG</i>	57	38.9	583	165	126.1
SHBG	<i>SHBG</i>	7	2.7	284	114	111.3
LDL cholesterol	<i>PCSK9</i>	94	4.0	759	99	95.0
Total bilirubin	<i>MROH2A</i>	33	4.6	1649	85.8	81.2
Apolipoprotein B	<i>PCSK9</i>	94	3.1	756	80.7	77.6
Cholesterol	<i>PCSK9</i>	94	4.0	760	80.9	76.9
IGF-1	<i>GH1</i>	5	2.1	301	55.1	53.0
Direct bilirubin	<i>MROH2A</i>	33	2.9	1497	55.7	52.8
Gamma glutamyltransferase	<i>GGT1</i>	5	0.008	545	52.1	52.1
Triglycerides	<i>ANGPTL3</i>	7	-0.02	337	39.9	39.9
Cholesterol	<i>ANGPTL3</i>	7	-0.6	337	34.3	34.9
Cholesterol	<i>APC</i>	1409	-34.7	1882	-0.5	34.2
LDL cholesterol	<i>APC</i>	1410	-33.7	1882	-0.5	33.2
Apolipoprotein B	<i>APC</i>	1405	-32.7	1876	-0.7	32.0
Total bilirubin	<i>UGT1A5</i>	12	2.0	225	33	31.0
Albumin	<i>APC</i>	1366	-31.6	1807	-1.2	30.4
Vitamin D	<i>APC</i>	1379	-29.8	1828	0.2	30.0
Creatinine	<i>APC</i>	1411	-31.4	1883	-1.9	29.5

**Table S1. Genes with considerable power gain in exome data as compared to array data.**

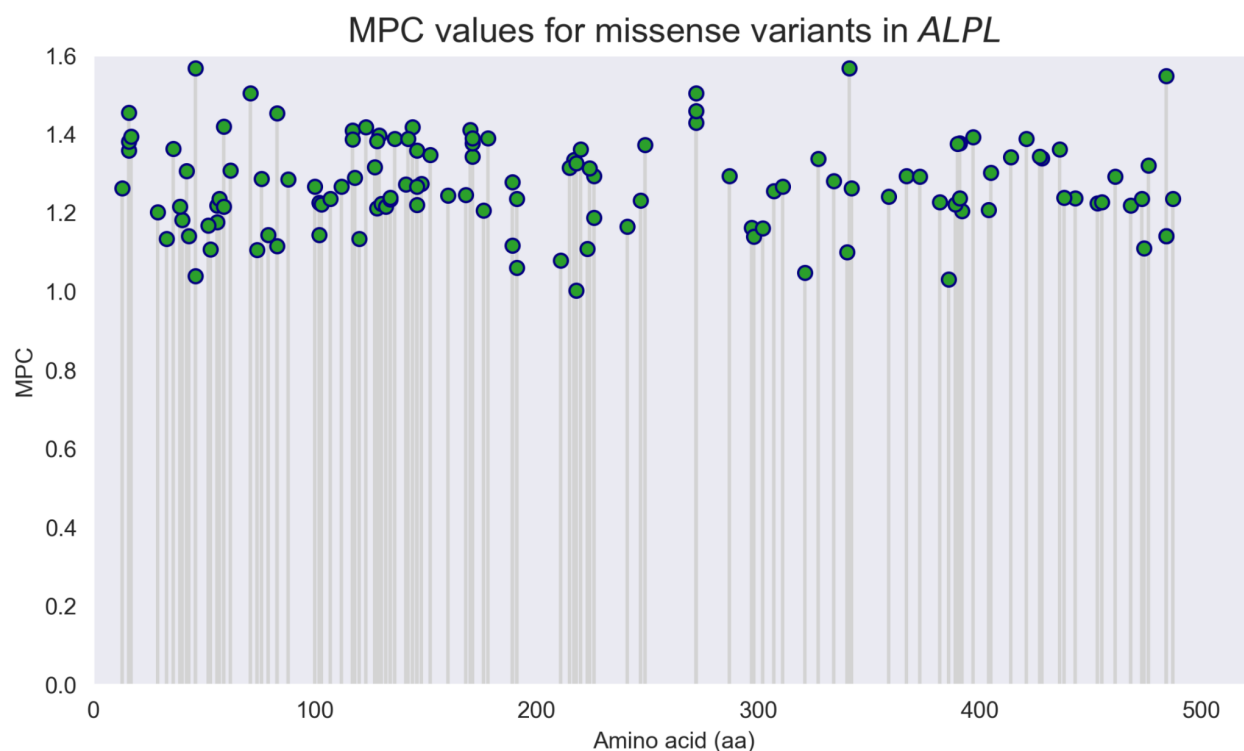


**Figure S2. From array to exome.** Scatterplot showing the increase in number of protein-altering variants in genes used in the analysis when comparing array (x-axis) to exome (y-axis) data. Data is taken from MRP calculations across 35 biomarker traits within the UK Biobank.

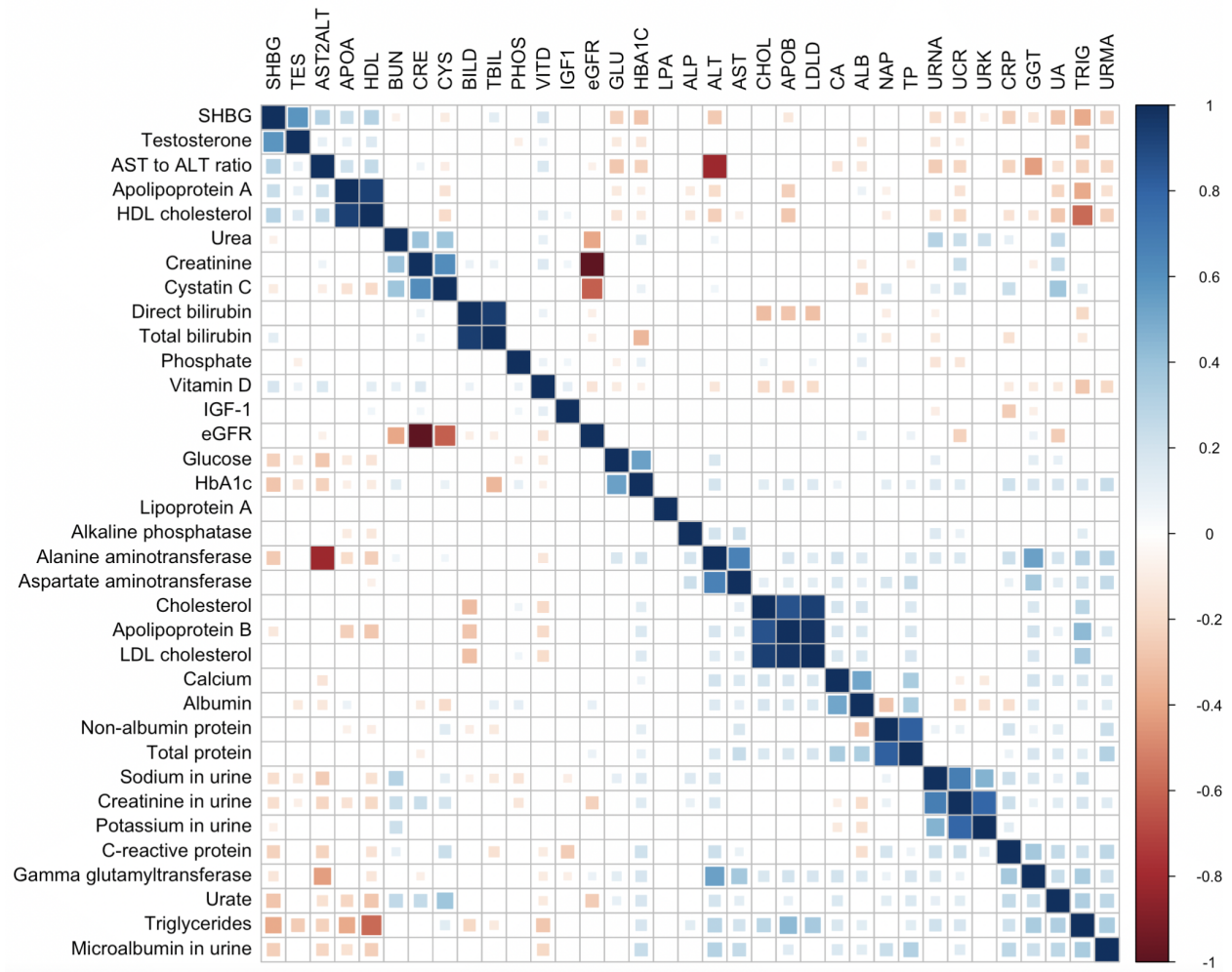
Trait	Gene	Number of PAVs	$\log_{10}$ BF without MPC	Number of MPC-augmented PAVs	Number of pLI-augmented PAVs	$\log_{10}$ BF with MPC	$\log_{10}$ BF Difference
Alkaline phosphatase	<i>ALPL</i>	198	126	93	0	160	34
Lipoprotein A	<i>LPA</i>	512	109	20	0	114	5
Apolipoprotein A	<i>APOA1</i>	102	11.7	30	0	15.7	4
HDL cholesterol	<i>APOA1</i>	103	9.36	30	0	13.2	3.84
Aspartate aminotransferase	<i>SLC30A10</i>	112	3.76	50	6	7.2	3.44
Phosphate	<i>ALPL</i>	192	10.9	91	0	14.3	3.4

Lipoprotein A	<i>IGF2R</i>	763	29.8	153	27	33.1	3.3
HDL cholesterol	<i>SCARB1</i>	220	5.45	66	0	8.29	2.84
Apolipoprotein B	<i>APOE</i>	142	5.48	60	0	8.27	2.79
Alanine aminotransferase	<i>SLC30A10</i>	112	2.94	50	6	5.56	2.62

**Table S3. Power comparison between variant annotation-based MRP and MPC/pLI-augmented MRP analyses across 35 biomarkers.** We see considerable gains in power in several gene/trait combinations.



**Figure S4. *ALPL* gene plot.** Gene plot showing variants for which MPC pathogenicity information was incorporated, resulting in a power gain for *ALPL* gene that encodes alkaline phosphatase; for the Alkaline phosphatase phenotype, the incorporation of this information resulted in a  $\log_{10}$ BF gain of 34 (Table S3).



**Figure S5. LD-score regression-based genetic correlation plots of all 35 biomarkers included in the multi-trait analyses.** The traits are ordered by hierarchical clustering.

## References

1. 1000 Genomes Project Consortium, Abecasis, G.R., Altshuler, D., Auton, A., Brooks, L.D., Durbin, R.M., Gibbs, R.A., Hurles, M.E., and McVean, G.A. (2010). A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061–1073.
2. Nejentsev, S., Walker, N., Riches, D., Egholm, M., and Todd, J.A. (2009). Rare variants of IFIH1, a gene implicated in antiviral responses, protect against type 1 diabetes. *Science* **324**, 387–389.
3. Rivas, M.A., Beaudoin, M., Gardet, A., Stevens, C., Sharma, Y., Zhang, C.K., Boucher, G., Ripke, S., Ellinghaus, D., Burt, N., et al. (2011). Deep resequencing of GWAS loci identifies independent rare variants associated with inflammatory bowel disease. *Nat. Genet.* **43**, 1066–1073.
4. Consortium, T. 1000 G.P., and The 1000 Genomes Project Consortium (2012). An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56–65.
5. Majithia, A.R., Flannick, J., Shahinian, P., Guo, M., Bray, M.-A., Fontanillas, P., Gabriel, S.B., GoT2D Consortium, NHGRI JHS/FHS Allelic Spectrum Project, SIGMA T2D Consortium, et al. (2014). Rare variants in PPARG with decreased activity in adipocyte differentiation are associated with increased risk of type 2 diabetes. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 13127–13132.
6. Lee, S., Abecasis, G.R., Boehnke, M., and Lin, X. (2014). Rare-variant association analysis: study designs and statistical tests. *Am. J. Hum. Genet.* **95**, 5–23.
7. Liu, D.J., Peloso, G.M., Zhan, X., Holmen, O.L., Zawistowski, M., Feng, S., Nikpay, M., Auer, P.L., Goel, A., Zhang, H., et al. (2014). Meta-analysis of gene-level tests for rare variant association. *Nat. Genet.* **46**, 200–204.
8. Cichonska, A., Rousu, J., Marttinen, P., Kangas, A.J., Soininen, P., Lehtimäki, T., Raitakari, O.T., Järvelin, M.-R., Salomaa, V., Ala-Korpela, M., et al. (2016). metaCCA: summary statistics-based multivariate meta-analysis of genome-wide association studies using canonical correlation analysis. *Bioinformatics* **32**, 1981–1989.
9. Rivas, M.A., Pirinen, M., Neville, M.J., Gaulton, K.J., Moutsianas, L., GoT2D Consortium, Lindgren, C.M., Karpe, F., McCarthy, M.I., and Donnelly, P. (2013). Assessing association between protein truncating variants and quantitative traits. *Bioinformatics* **29**, 2419–2426.
10. Rivas, M.A., Pirinen, M., Conrad, D.F., Lek, M., Tsang, E.K., Karczewski, K.J., Maller, J.B., Kukurba, K.R., DeLuca, D.S., Fromer, M., et al. (2015). Human genomics. Effect of predicted protein-truncating genetic variants on the human transcriptome. *Science* **348**, 666–669.
11. Samocha, K.E., Kosmicki, J.A., Karczewski, K.J., O'Donnell-Luria, A.H., Pierce-Hoffman, E., MacArthur, D.G., Neale, B.M., and Daly, M.J. Regional missense constraint improves variant deleteriousness prediction.
12. Fuller, Z.L., Berg, J.J., Mostafavi, H., Sella, G., and Przeworski, M. (2019). Measuring

intolerance to mutation in human genetics. *Nat. Genet.* 51, 772–776.

13. Cohen, J., Pertsemlidis, A., Kotowski, I.K., Graham, R., Garcia, C.K., and Hobbs, H.H. (2005). Low LDL cholesterol in individuals of African descent resulting from frequent nonsense mutations in PCSK9. *Nat. Genet.* 37, 161–165.

14. Cohen, J.C., Boerwinkle, E., Mosley, T.H., Jr, and Hobbs, H.H. (2006). Sequence variations in PCSK9, low LDL, and protection against coronary heart disease. *N. Engl. J. Med.* 354, 1264–1272.

15. Sullivan, D., Olsson, A.G., Scott, R., Kim, J.B., Xue, A., GebSKI, V., Wasserman, S.M., and Stein, E.A. (2012). Effect of a monoclonal antibody to PCSK9 on low-density lipoprotein cholesterol levels in statin-intolerant patients: the GAUSS randomized trial. *JAMA* 308, 2497–2506.

16. Wu, M.C., Lee, S., Cai, T., Li, Y., Boehnke, M., and Lin, X. (2011). Rare-Variant Association Testing for Sequencing Data with the Sequence Kernel Association Test. *The American Journal of Human Genetics* 89, 82–93.

17. Sinnott-Armstrong, N., Tanigawa, Y., Amar, D., Mars, N., Benner, C., Aguirre, M., Venkataraman, G.R., Wainberg, M., Ollila, H.M., Kiiskinen, T., et al. (2021). Genetics of 35 blood and urine biomarkers in the UK Biobank. *Nat. Genet.* 53, 185–194.

18. Band, G., Le, Q.S., Jostins, L., Pirinen, M., Kivinen, K., Jallow, M., Sisay-Joof, F., Bojang, K., Pinder, M., Sirugo, G., et al. (2013). Imputation-based meta-analysis of severe malaria in three African populations. *PLoS Genet.* 9, e1003509.

19. MacArthur, D.G., Balasubramanian, S., Frankish, A., Huang, N., Morris, J., Walter, K., Jostins, L., Habegger, L., Pickrell, J.K., Montgomery, S.B., et al. (2012). A systematic survey of loss-of-function variants in human protein-coding genes. *Science* 335, 823–828.

20. Findlay, G.M., Boyle, E.A., Hause, R.J., Klein, J.C., and Shendure, J. (2014). Saturation editing of genomic regions by multiplex homology-directed repair. *Nature* 513, 120–123.

21. Neale, B.M., Rivas, M.A., Voight, B.F., Altshuler, D., Devlin, B., Orho-Melander, M., Kathiresan, S., Purcell, S.M., Roeder, K., and Daly, M.J. (2011). Testing for an unusual distribution of rare variants. *PLoS Genet.* 7, e1001322.

22. Clarke, G.M., Rivas, M.A., and Morris, A.P. (2013). A flexible approach for the analysis of rare variants allowing for a mixture of effects on binary or quantitative traits. *PLoS Genet.* 9, e1003694.

23. Cotsapas, C., Voight, B.F., Rossin, E., Lage, K., Neale, B.M., Wallace, C., Abecasis, G.R., Barrett, J.C., Behrens, T., Cho, J., et al. (2011). Pervasive sharing of genetic effects in autoimmune disease. *PLoS Genet.* 7, e1002254.

24. Solovieff, N., Cotsapas, C., Lee, P.H., Purcell, S.M., and Smoller, J.W. (2013). Pleiotropy in complex traits: challenges and strategies. *Nat. Rev. Genet.* 14, 483–495.

25. Bulik-Sullivan, B., Finucane, H.K., Anttila, V., Gusev, A., Day, F.R., Loh, P.-R., ReproGen



Consortium, Psychiatric Genomics Consortium, Genetic Consortium for Anorexia Nervosa of the Wellcome Trust Case Control Consortium 3, Duncan, L., et al. (2015). An atlas of genetic correlations across human diseases and traits. *Nat. Genet.* 47, 1236–1241.

26. Do, R., Stitzel, N.O., Won, H.-H., Jørgensen, A.B., Duga, S., Angelica Merlini, P., Kiezun, A., Farrall, M., Goel, A., Zuk, O., et al. (2015). Exome sequencing identifies rare LDLR and APOA5 alleles conferring risk for myocardial infarction. *Nature* 518, 102–106.

27. Do, R., Willer, C.J., Schmidt, E.M., Sengupta, S., Gao, C., Peloso, G.M., Gustafsson, S., Kanoni, S., Ganna, A., Chen, J., et al. (2013). Common variants associated with plasma triglycerides and risk for coronary artery disease. *Nat. Genet.* 45, 1345–1352.

28. (2015). Loss-of-Function Mutations in APOC3, Triglycerides, and Coronary Disease. *N. Engl. J. Med.* 372, 690–690.

29. Consortium, T. 1000 G.P., and The 1000 Genomes Project Consortium (2015). A global reference for human genetic variation. *Nature* 526, 68–74.

30. Bergström, A., McCarthy, S.A., Hui, R., Almarri, M.A., Ayub, Q., Danecek, P., Chen, Y., Felkel, S., Hallast, P., Kamm, J., et al. (2020). Insights into human genetic variation and population history from 929 diverse genomes. *Science* 367,.

31. Mallick, S., Li, H., Lipson, M., Mathieson, I., Gymrek, M., Racimo, F., Zhao, M., Chennagiri, N., Nordenfelt, S., Tandon, A., et al. (2016). The Simons Genome Diversity Project: 300 genomes from 142 diverse populations. *Nature* 538, 201–206.

32. Alexander, D.H., and Lange, K. (2011). Enhancements to the ADMIXTURE algorithm for individual ancestry estimation. *BMC Bioinformatics* 12,.

33. Maples, B.K., Gravel, S., Kenny, E.E., and Bustamante, C.D. (2013). RFMix: a discriminative modeling approach for rapid and robust local-ancestry inference. *Am. J. Hum. Genet.* 93, 278–288.

34. McLaren, W., Gil, L., Hunt, S.E., Riat, H.S., Ritchie, G.R.S., Thormann, A., Flicek, P., and Cunningham, F. (2016). The Ensembl Variant Effect Predictor. *Genome Biol.* 17, 122.

35. McInnes, G., Tanigawa, Y., DeBoever, C., Lavertu, A., Olivieri, J.E., Aguirre, M., and Rivas, M.A. (2019). Global Biobank Engine: enabling genotype-phenotype browsing for biobank summary statistics. *Bioinformatics* 35, 2495–2497.

36. Chambers, J.M., and Murtagh, F. (1985). *Multidimensional Clustering Algorithms* (Springer).

37. Institute, T.T.G.A.H.W.G. of T.E.S.P.N.H.L.A.B., The TG and HDL Working Group of the Exome Sequencing Project, National Heart, Lung, and and Blood Institute (2014). Loss-of-Function Mutations in APOC3, Triglycerides, and Coronary Disease. *New England Journal of Medicine* 371, 22–31.

38. Pollin, T.I., Damcott, C.M., Shen, H., Ott, S.H., Shelton, J., Horenstein, R.B., Post, W., McLenithan, J.C., Bielak, L.F., Peyser, P.A., et al. (2008). A Null Mutation in Human APOC3 Confers a Favorable Plasma Lipid Profile and Apparent Cardioprotection. *Science* 322,

1702–1705.

39. Jørgensen, A.B., Frikke-Schmidt, R., Nordestgaard, B., and Tybjaerg-Hansen, A. (2016). Loss-of-function mutations in APOC3, remnant cholesterol, LDL cholesterol, and risk of ischemic vascular disease. *Atherosclerosis* 252, e251–e252.

40. Kapur, K., Johnson, T., Beckmann, N.D., Sehmi, J., Tanaka, T., Kutalik, Z., Styrkarsdottir, U., Zhang, W., Marek, D., Gudbjartsson, D.F., et al. (2010). Genome-wide meta-analysis for serum calcium identifies significantly associated SNPs near the calcium-sensing receptor (CASR) gene. *PLoS Genet.* 6, e1001035.

41. Manousaki, D., Mitchell, R., Dudding, T., Haworth, S., Harroud, A., Forgetta, V., Shah, R.L., Luan, J., Langenberg, C., Timpson, N.J., et al. (2020). Genome-wide Association Study for Vitamin D Levels Reveals 69 Independent Loci. *The American Journal of Human Genetics* 106, 327–337.

42. Kanai, M., Akiyama, M., Takahashi, A., Matoba, N., Momozawa, Y., Ikeda, M., Iwata, N., Ikegawa, S., Hirata, M., Matsuda, K., et al. (2018). Genetic analysis of quantitative traits in the Japanese population links cell types to complex human diseases. *Nat. Genet.* 50, 390–400.

43. Mack, S., Coassin, S., Rueedi, R., Yousri, N.A., Seppälä, I., Gieger, C., Schönherr, S., Forer, L., Erhart, G., Marques-Vidal, P., et al. (2017). A genome-wide association meta-analysis on lipoprotein (a) concentrations adjusted for apolipoprotein (a) isoforms. *J. Lipid Res.* 58, 1834–1844.

44. Richardson, T.G., Sanderson, E., Palmer, T.M., Ala-Korpela, M., Ference, B.A., Davey Smith, G., and Holmes, M.V. (2020). Evaluating the relationship between circulating lipoprotein lipids and apolipoproteins with risk of coronary heart disease: A multivariable Mendelian randomisation analysis. *PLoS Med.* 17, e1003062.

45. Kilpeläinen, T.O., Bentley, A.R., Noordam, R., Sung, Y.J., Schwander, K., Winkler, T.W., Jakupović, H., Chasman, D.I., Manning, A., Ntalla, I., et al. (2019). Multi-ancestry study of blood lipid levels identifies four loci interacting with physical activity. *Nat. Commun.* 10, 376.

46. Suhre, K., Arnold, M., Bhagwat, A.M., Cotton, R.J., Engelke, R., Raffler, J., Sarwath, H., Thareja, G., Wahl, A., DeLisle, R.K., et al. (2017). Connecting genetic risk to disease end points through the human blood plasma proteome. *Nat. Commun.* 8, 14357.

47. Hoffmann, T.J., Theusch, E., Haldar, T., Ranatunga, D.K., Jorgenson, E., Medina, M.W., Kvale, M.N., Kwok, P.-Y., Schaefer, C., Krauss, R.M., et al. (2018). A large electronic-health-record-based genome-wide study of serum lipids. *Nat. Genet.* 50, 401–413.

48. Wojcik, G.L., Graff, M., Nishimura, K.K., Tao, R., Haessler, J., Gignoux, C.R., Highland, H.M., Patel, Y.M., Sorokin, E.P., Avery, C.L., et al. (2019). Genetic analyses of diverse populations improves discovery for complex traits. *Nature* 570, 514–518.

49. Nielsen, J.B., Rom, O., Surakka, I., Graham, S.E., Zhou, W., Roychowdhury, T., Fritsche, L.G., Gagliano Taliun, S.A., Sidore, C., Liu, Y., et al. (2020). Loss-of-function genomic variants highlight potential therapeutic targets for cardiovascular disease. *Nat. Commun.* 11, 6417.

50. Bentley, A.R., Sung, Y.J., Brown, M.R., Winkler, T.W., Kraja, A.T., Ntalla, I., Schwander, K., Chasman, D.I., Lim, E., Deng, X., et al. (2019). Multi-ancestry genome-wide gene-smoking interaction study of 387,272 individuals identifies new loci associated with serum lipids. *Nat. Genet.* *51*, 636–648.
51. Wheeler, E., Leong, A., Liu, C.-T., Hivert, M.-F., Strawbridge, R.J., Podmore, C., Li, M., Yao, J., Sim, X., Hong, J., et al. (2017). Impact of common genetic determinants of Hemoglobin A1c on type 2 diabetes risk and diagnosis in ancestrally diverse populations: A transethnic genome-wide meta-analysis. *PLoS Med.* *14*, e1002383.
52. Gurdasani, D., Carstensen, T., Fatumo, S., Chen, G., Franklin, C.S., Prado-Martinez, J., Bouman, H., Abascal, F., Haber, M., Tachmazidou, I., et al. (2019). Uganda Genome Resource Enables Insights into Population History and Genomic Discovery in Africa. *Cell* *179*, 984–1002.e36.
53. Moon, S., Kim, Y.J., Han, S., Hwang, M.Y., Shin, D.M., Park, M.Y., Lu, Y., Yoon, K., Jang, H.-M., Kim, Y.K., et al. (2019). The Korea Biobank Array: Design and Identification of Coding Variants Associated with Blood Biochemical Traits. *Scientific Reports* *9*,.
54. Seo, J.Y., Lee, J.-E., Chung, G.E., Shin, E., Kwak, M.-S., Yang, J.I., and Yim, J.Y. (2020). A genome-wide association study on liver enzymes in Korean population. *PLoS One* *15*, e0229374.
55. Kang, T.-W., Kim, H.-J., Ju, H., Kim, J.-H., Jeon, Y.-J., Lee, H.-C., Kim, K.-K., Kim, J.-W., Lee, S., Kim, J.Y., et al. (2010). Genome-wide association of serum bilirubin levels in Korean population. *Hum. Mol. Genet.* *19*, 3672–3678.
56. Graham, S.E., Nielsen, J.B., Zawistowski, M., Zhou, W., Fritsche, L.G., Gabrielsen, M.E., Skogholt, A.H., Surakka, I., Hornsby, W.E., Fermin, D., et al. (2019). Sex-specific and pleiotropic effects underlying kidney function identified from GWAS meta-analysis. *Nat. Commun.* *10*, 1847.
57. Casanova, F., Tyrrell, J., Beaumont, R.N., Ji, Y., Jones, S.E., Hattersley, A.T., Weedon, M.N., Murray, A., Shore, A.C., Frayling, T.M., et al. (2019). A genome-wide association study implicates multiple mechanisms influencing raised urinary albumin-creatinine ratio. *Hum. Mol. Genet.* *28*, 4197–4207.
58. Gill, D., Cameron, A.C., Burgess, S., Li, X., Doherty, D.J., Karhunen, V., Abdul-Rahim, A.H., Taylor-Rowan, M., Zuber, V., Tsao, P.S., et al. (2021). Urate, Blood Pressure, and Cardiovascular Disease: Evidence From Mendelian Randomization and Meta-Analysis of Clinical Trials. *Hypertension* *77*, 383–392.
59. Kamatani, Y., Matsuda, K., Okada, Y., Kubo, M., Hosono, N., Daigo, Y., Nakamura, Y., and Kamatani, N. (2010). Genome-wide association study of hematological and biochemical traits in a Japanese population. *Nature Genetics* *42*, 210–215.
60. Gorski, M., van der Most, P.J., Teumer, A., Chu, A.Y., Li, M., Mijatovic, V., Nolte, I.M., Cocca, M., Taliun, D., Gomez, F., et al. (2017). 1000 Genomes-based meta-analysis identifies 10 novel loci for kidney function. *Sci. Rep.* *7*, 45040.
61. Astle, W.J., Elding, H., Jiang, T., Allen, D., Ruklisa, D., Mann, A.L., Mead, D., Bouman, H.,

- Riveros-Mckay, F., Kostadima, M.A., et al. (2016). The Allelic Landscape of Human Blood Cell Trait Variation and Links to Common Complex Disease. *Cell* 167, 1415–1429.e19.
62. Morgan, M.D., Pairo-Castineira, E., Rawlik, K., Canela-Xandri, O., Rees, J., Sims, D., Tenesa, A., and Jackson, I.J. (2018). Genome-wide study of hair colour in UK Biobank explains most of the SNP heritability. *Nat. Commun.* 9, 5271.
63. Adhikari, K., Mendoza-Revilla, J., Sohail, A., Fuentes-Guajardo, M., Lampert, J., Chacón-Duque, J.C., Hurtado, M., Villegas, V., Granja, V., Acuña-Alonzo, V., et al. (2019). A GWAS in Latin Americans highlights the convergent evolution of lighter skin pigmentation in Eurasia. *Nat. Commun.* 10, 358.
64. Liu, F., Visser, M., Duffy, D.L., Hysi, P.G., Jacobs, L.C., Lao, O., Zhong, K., Walsh, S., Chaitanya, L., Wollstein, A., et al. (2015). Genetics of skin color variation in Europeans: genome-wide association studies with functional follow-up. *Hum. Genet.* 134, 823–835.
65. Hysi, P.G., Valdes, A.M., Liu, F., Furlotte, N.A., Evans, D.M., Bataille, V., Visconti, A., Hemani, G., McMahon, G., Ring, S.M., et al. (2018). Genome-wide association meta-analysis of individuals of European ancestry identifies new loci explaining a substantial fraction of hair color variation and heritability. *Nat. Genet.* 50, 652–656.
66. Lona-Durazo, F., Hernandez-Pacheco, N., Fan, S., Zhang, T., Choi, J., Kovacs, M.A., Loftus, S.K., Le, P., Edwards, M., Fortes-Lima, C.A., et al. (2019). Meta-analysis of GWA studies provides new insights on the genetic architecture of skin pigmentation in recently admixed populations. *BMC Genet.* 20, 59.
67. Vuckovic, D., Bao, E.L., Akbari, P., Lareau, C.A., Mousas, A., Jiang, T., Chen, M.-H., Raffield, L.M., Tardaguila, M., Huffman, J.E., et al. (2020). The Polygenic and Monogenic Basis of Blood Traits and Diseases. *Cell* 182, 1214–1231.e11.
68. Gieger, C., Radhakrishnan, A., Cvejic, A., Tang, W., Porcu, E., Pistis, G., Serbanovic-Canic, J., Elling, U., Goodall, A.H., Labrune, Y., et al. (2011). New gene functions in megakaryopoiesis and platelet formation. *Nature* 480, 201–208.
69. Hodonsky, C.J., Jain, D., Schick, U.M., Morrison, J.V., Brown, L., McHugh, C.P., Schurmann, C., Chen, D.D., Liu, Y.M., Auer, P.L., et al. (2017). Genome-wide association study of red blood cell traits in Hispanics/Latinos: The Hispanic Community Health Study/Study of Latinos. *PLoS Genet.* 13, e1006760.
70. Chen, M.-H., Raffield, L.M., Mousas, A., Sakaue, S., Huffman, J.E., Moscati, A., Trivedi, B., Jiang, T., Akbari, P., Vuckovic, D., et al. (2020). Trans-ethnic and Ancestry-Specific Blood-Cell Genetics in 746,667 Individuals from 5 Global Populations. *Cell* 182, 1198–1213.e14.
71. Hodonsky, C.J., Baldassari, A.R., Bien, S.A., Raffield, L.M., Highland, H.M., Sitlani, C.M., Wojcik, G.L., Tao, R., Graff, M., Tang, W., et al. (2020). Ancestry-specific associations identified in genome-wide combined-phenotype study of red blood cell traits emphasize benefits of diversity in genomics. *BMC Genomics* 21, 228.
72. Hutton, J.C., and O'Brien, R.M. (2009). Glucose-6-phosphatase Catalytic Subunit Gene

Family. *Journal of Biological Chemistry* 284, 29241–29245.