

Bayesian model comparison for rare variant association studies of multiple phenotypes

Christopher DeBoever¹, Matthew Aguirre¹, Yosuke Tanigawa¹, Timothy Poterba², Mark J. Daly^{2,3},
Matti Pirinen^{4*}, and Manuel A. Rivas^{1†}

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 model comparison approach we refer to as MRP for rare variant association studies that considers correlation, scale, and location of genetic effects across a group of genetic variants, phenotypes, and studies. In so doing we consider the use of summary statistic data to 1) apply standard univariate and multivariate gene-based meta-analysis models, 2) models for assessing heterogeneity of genetic effects, which may be used in practice for downstream quality control, and 3) models for identifying protective protein-truncating variants, which can expedite drug discovery. Through simulation studies, we demonstrate that the proposed model comparison approach may improve ability to detect rare variant association signals. Finally, we demonstrate its use with rare variant data combined with asthma diagnosis and haematological and spirometry measures for individuals in the UK Biobank. We show that we are able to retain useful features from widely-used meta-analysis approaches and prioritize protective modifiers of disease risk.

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

²Broad Institute of MIT and Harvard, Cambridge, MA, USA

³Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, UK

⁴Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA

⁵Oxford Centre for Diabetes, Endocrinology and Metabolism, Oxford, UK

⁶Institute for Molecular Medicine Finland, University of Helsinki, Finland

*matti.pirinen@helsinki.fi

†mrivas@stanford.edu

1 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 feasible¹⁻⁴. In genetic studies of diseases or continuous phenotypes, rare variants are hard to assess individually due to the limited number of copies of each rare variant. Hence, to boost the ability to detect a signal, evidence is usually ‘aggregated’ across variants. When designing an ‘aggregation’ method, there are three questions that are usually considered. First, across which biological units should variants be combined; second, which variants mapping within those units should be included⁵; and third, which statistical model should be used⁶? Given the widespread observations of shared genetic risk factors across distinct diseases, there is also considerable motivation to proactively modify gene discovery approaches accordingly. 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 sets of variants might have effects.

In this paper, we present a Bayesian multiple rare variants and phenotypes (MRP) model comparison approach for identifying rare variant associations as an alternative to current widely-used statistical tests. The MRP framework, as implemented, can exploit correlation, scale, or location (direction) of genetic effects in a broad range of rare variant association study designs including: case-control; multiple diseases and shared controls; single continuous phenotype; mul-

48 tiple continuous phenotypes; or a mixture of case-control and multiple continuous phenotypes
49 (Fig. 1).

50 In Bayesian model comparison, we compute a Bayes Factor (BF) defined as the ratio of
51 the marginal likelihoods of the observed data under two models: 1) a pre-specified null where
52 all genetic effects are zero; and 2) an alternative model where factors like correlation, scale, or
53 location of genetic effects are considered.

54 In practice, sharing of individual genotype and phenotype data across groups in large genetic
55 consortium studies is difficult to achieve. We present a derivation where summary statistics, such
56 as estimates of effect size and the corresponding standard error from single variant and single phe-
57 notype analysis, are the input data. Furthermore, we use insights from Liu et al.⁷ and Cichonska et
58 al.⁸ who suggest the use of additional summary statistics, like covariance estimates across variants
59 and studies, respectively, that would enable lossless ability to detect gene-based association signals
60 using summary statistics alone.

61 Genome variant annotation is critical for the interpretation of genetic findings. In the MRP
62 model comparison we present approaches for including priors on spread of effect sizes. An ex-
63 ample of an important class of variants where functional prediction is much more straightforward
64 is the class of variants predicted to truncate the protein product, commonly referred to as protein
65 truncating variants or PTV^{9,10}, and in this paper we place special emphasis on the interpretation of
66 PTVs.

Given the translational potential of identifying protective modifiers of human disease risk we characterize how model comparison may be used to improve discovery of such protective signals. We present an extension allowing for modeling of location (direction) of genetic effects, which may be useful for prioritizing variants or genes that are consistent with a protective profile of disease^{11–13}.

To evaluate the performance of MRP and to study its behavior 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)¹⁴, and the burden test, which we show 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. Furthermore, we show that MRP is able to detect the presence of heterogeneity of effects, which in some circumstances, such as those where effects are observed to be heterogeneous across studies, may indicate the presence of technical sources of error.

We also apply MRP to two groups of related phenotypes. First, we apply MRP to asthma (HC382: the corresponding phenotype label in Global Biobank Engine [<https://biobankengine.stanford.edu>]), eosinophil count (INI30150), Forced Expiratory Volume in 1-second (FEV1, INI3063), and Forced Vital Capacity (FVC, INI3062) in the UK Biobank dataset. We next apply MRP to glaucoma (HCX276), intra-ocular pressure (INI5263), and corneal resistance factor (INI5265). We show that we are able to retain useful features from single variant and single phenotype meta-analysis, rare variant analysis, and prioritize protective modifiers of disease risk.

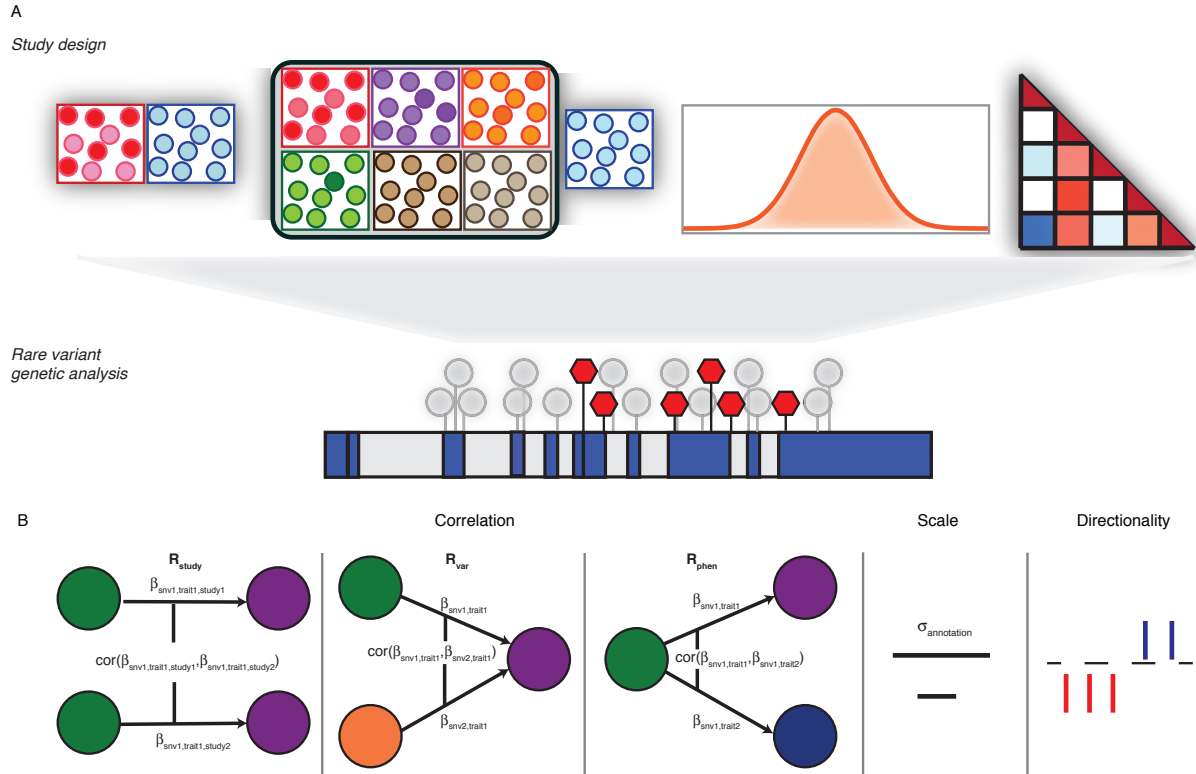


Figure 1 Schematic overview of MRP. A. MRP is suitable for a broad range of rare variant association study designs including (from left to right): i) case-control, ii) multiple diseases with shared controls, iii) single quantitative phenotype, and iv) mixture of case-control and quantitative phenotypes. B. 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), R_{var} (expected correlation of genetic effects among a group of variants), and R_{phen} (expected correlation of genetic effects among a group of phenotypes); the scale parameter for genetic variant annotation; and the location of genetic effects, which may be used to prioritize or identify protective modifiers of disease risk.

2 Results

2.1 Methods overview

A complete description of MRP, along with an analytical derivation, is given in the Online Methods and Appendix. Briefly, we describe the approach.

2.1.1 Model comparison

In Bayesian model comparison, we compare the null and alternative models via the Bayes Factor (BF) defined as the ratio of the marginal likelihoods of the observed data under two models. Here, the null model is that the effect sizes obtained across all studies for a group of variants and a group of phenotypes is zero. To complete the specification of the alternative model we consider the correlation structure, the scale, and the location of the effect sizes (which we find to be informative for alternative models where we seek to identify variants with protective effects).

Let N be the number of individuals and K the number of quantitative phenotype measurements on each individual, for example, here in the applications we focus on asthma, eosinophil count, Forced Expiratory Volume in 1-second, and Forced Vital Capacity, in which case $K = 4$. 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 where data is obtained from - this data may be in the form of raw genotypes and phenotypes or summary statistics including linkage-disequilibrium,

effect sizes (or odds ratio), and standard error of the effect size.

In the multivariate setting we define the prior correlation structure of the effect sizes, denoted by a SMK-by-SMK matrix \mathbf{U} . In practice, we define \mathbf{U} as a Kronecker product, an operation of matrices of arbitrary size, of three sub-matrices: an S-by-S matrix $\mathbf{R}_{\text{study}}$ being the correlations of genetic effects among studies where different values can be used to compare different models of association, such as for identifying heterogeneity of effect sizes between populations¹⁵; an M-by-M matrix \mathbf{R}_{var} being the correlations of genetic effects among genetic variants, which may reflect the assumption that all the PTVs in a gene may have the same biological consequence^{9,10,16} or prior information obtained through integration of additional data sources, such as functional assay data^{5,17}, otherwise zero correlation of genetic effects may be assumed, which is used in dispersion tests like C-alpha^{18,19} and SKAT¹⁴; and a K-by-K matrix \mathbf{R}_{phen} being the correlations of genetic effects among phenotypes, which may be obtained from common variant data^{20–22}. The variance-covariance matrix of the effect sizes may be obtained from readily available summary statistic data such as in-study LD matrices, effect size estimate (or log odds ratio), and the standard error of the effect size estimate (Online Methods).

For gene discovery, the scale may be used to denote the prior spread of effect sizes. For instance, emerging empirical genetic studies have shown that within a gene protein truncating variants may have stronger effects than missense variants²³.

Thus far we have assumed that the prior mean, or location, of genetic effects is zero. This makes it feasible to analyze a large number of phenotypes without enumerating the prior mean

across all phenotypes. However, in practice, we may want to proactively identify genetic variants that have effects that are consistent with a protective profile for a disease, which may exploit information from mendelian randomization studies of common variants, such as recent rare variant findings in *PCSK9* where truncating loss-of-function variants are found to decrease LDL and triglyceride levels, and decrease CAD risk^{11,24–26}. This is straightforward to consider by including a non-zero vector as a prior mean of genetic effect (Online Methods).

Although we see advantages in adopting a Bayesian perspective, our approach could be used in a frequentist context by calculating a Bayes factor and using it as a test statistic to compute p-values (Methods, Fig. 2).

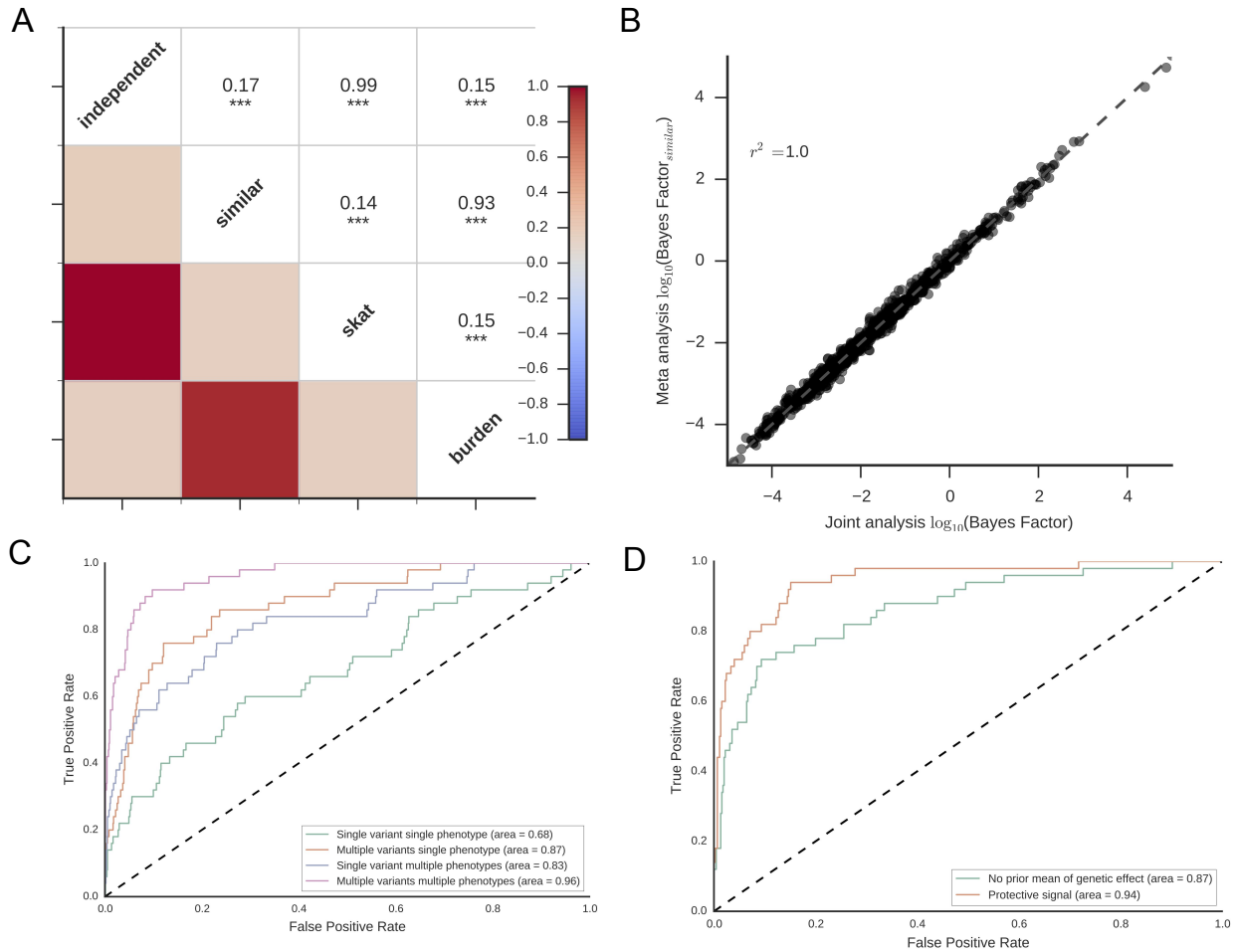


Figure 2 Simulation studies. A. Comparison of $-\log_{10}(\text{p-values})$ from frequentist BF_{MRP} approximation for an independent effects and a similar effects model to commonly used gene-based statistical tests (skat and burden). B. Comparison of $\log_{10}(\text{Bayes Factors})$ obtained when raw genotype and phenotype data is available to a scenario where summary statistics only was available and similar effects across studies is assumed. C. From single variant and single phenotype to multiple variants and multiple phenotypes gene discovery: ROC curves for detecting gene association to any of the phenotypes using single variant/single phenotype association (blue) to multiple variants and multiple phenotypes

association (purple). D. ROC curves for detecting gene association when incorporating prior mean of genetic effects (orange) to identify protective alleles.

2.2 Simulation studies

We first verified the analytical derivations and examined the properties of the approach under a simulation framework.

2.2.1 Comparison to frequentist gene tests

For the analysis of multiple rare variants and a single phenotype we compared it to the burden test and the SKAT test, commonly used statistical tests in rare variant association studies of a single phenotype. We observe concordance between the frequentist methods and the Bayesian models. To compare the Bayesian models we compute p-values by using the Bayes Factor as the test statistic and approximating it using distribution properties of quadratic forms (Appendix). As expected, an independent effects model has high correlation with the gene-based test SKAT ($r^2 = 0.99$), whereas the similar effects model has high correlation with the burden test ($r^2 = 0.93$, Fig. 2A).

2.2.2 Summary statistic data

To study the behavior of MRP using summary statistics we simulate two scenarios: first, the scenario where analysts have access to all the raw genotype and phenotype data; and second, the scenario where analysts only have access to summary statistics data⁷. We conducted 1000 simulation experiments where we let K (the number of phenotypes) = 3, M (the number of variants) = 10, S (the number of studies) = 2, N_0 (number of individuals in study with access to all the data) = 10000, N_1 (meta-analysis study 1) = 5000, N_2 (meta-analysis study 2) = 5000. We find that, under the scenario where similar effects are assumed across studies, the Bayes Factors obtained using summary statistics alone are strongly correlated ($r^2 = 1$) to Bayes Factors obtained by the full genotype and phenotype data (Fig. 2B).

2.2.3 From single variant and single phenotype analysis to multiple variants and multiple phenotypes

To validate the flexibility of the approach we conducted a simulation experiment where we assumed an allelic architecture consistent to that discovered for *APOC3* in relation to coronary artery disease (CAD), triglycerides (TG), low-density lipoprotein cholesterol (LDL-C), and high-density lipoprotein cholesterol (HDL-C)^{26–29}. We simulated three studies and applied the model comparison unit jointly to summary statistic data obtained for each study (Supplementary Note). Overall, we observed that considering the joint effects across multiple studies in a group of variants

and phenotypes may improve ability to detect gene-based signals (Fig. 2C), and that considering prior mean of genetic effects should aid in efforts to identify protective modifiers of disease risk (Fig. 2D).

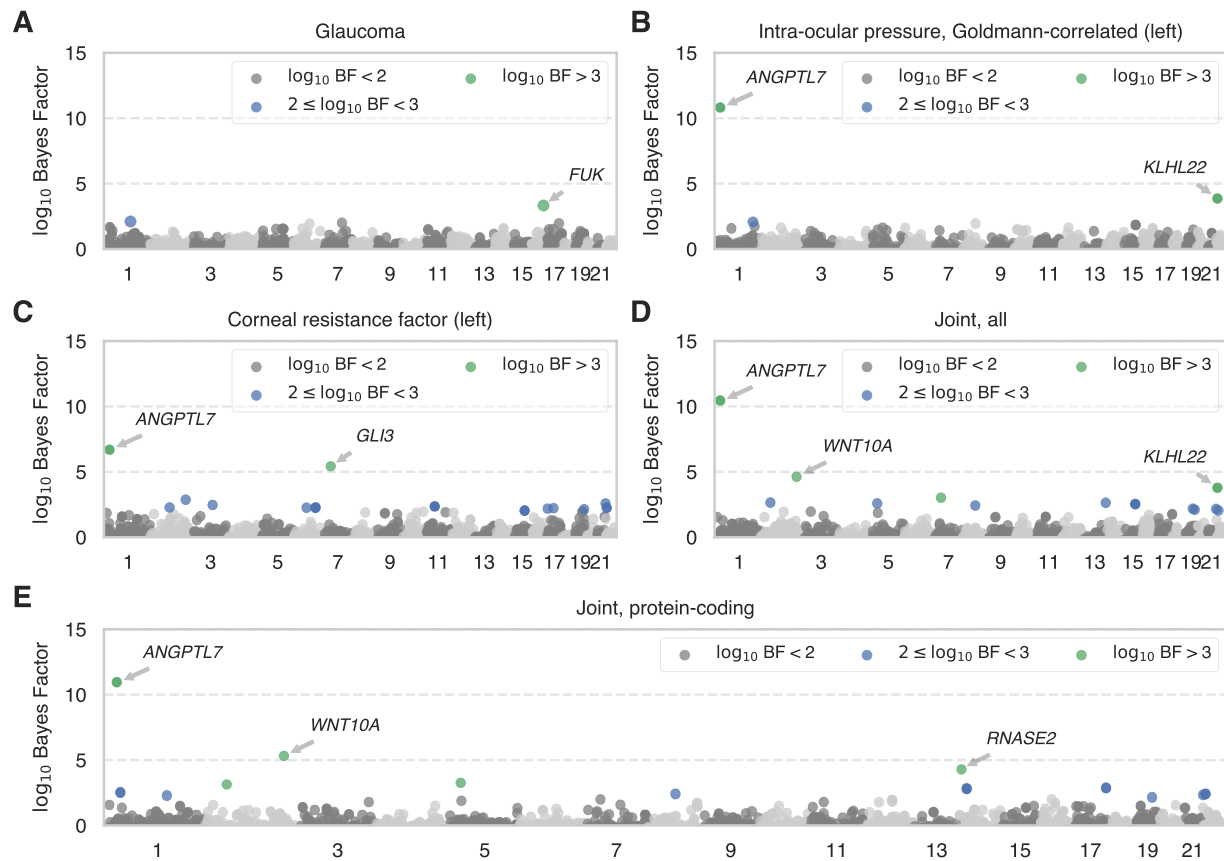
2.3 Applications

We demonstrate its use in rare variant analysis in the UK Biobank.

2.3.1 Model comparison applied to asthma, eosinophil count, forced expiratory volume in 1-second, and forced vital capacity summary statistic data

We applied the MRP model comparison approach to summary statistic data for coding variants on the UK Biobank array. The summary statistic data was generated from single variant logistic regression and linear regression analysis applied to asthma, eosinophil count, forced expiratory volume in 1-second, and forced vital capacity phenotypes. We applied the model to missense variants and PTVs for each phenotype separately (Fig. ??A-D) as well as to all phenotypes jointly (Fig. ??E) and obtained \log_{10} BF for each gene. We observed evidence that missense variants and PTVs in *IL33* affect eosinophil counts and offer protection from asthma from the single-phenotype analyses^{30,31}, though the evidence of association was strongest for the joint analysis (Table 1). We also found evidence for associations between missense variants and PTVs in several other genes and the four phenotypes in the joint model that were not identified when considering asthma alone.

For instance, *GPR20*, *CSF2RB*, *PACS2*, and *CCR3* all had \log_{10} BF greater than 9 in the eosinophil count-only analysis but \log_{10} BF less than 0 for the asthma-only analysis. In the joint model, the \log_{10} BFs for these genes were all greater than six indicating possible associations with these traits. Besides *IL33*, none of the genes with \log_{10} BF greater than six in the joint model were reported in a large GWAS for allergic disease including asthma³² though *CCR3* is near a locus associated with atopy in a previous meta-analysis³³. These results illustrate the ability of applying model comparison to multiple traits simultaneously to identify new genes involved in disease.



208

Figure 3 \log_{10} Bayes Factors from applying MRP to summary statistics for missense and protein-truncating variants from (A) eosinophil counts (INI30150), (B) forced vital capacity

210

(FVC, INI3062), (C) forced expiratory volume in 1-second (FEV1, INI3063), (D) asthma (HC382), and (E) all four traits jointly. The four genes outside of chromosome 6 with the largest Bayes Factors are labeled in each plot. Only \log_{10} Bayes Factors greater than -5 are plotted.

3 Discussion

In this study, we developed a Bayesian model comparison approach MRP that shares information across both variants and phenotypes to identify rare variant associations. We used simulations to compare MRP to the widely used burden and SKAT tests for identifying rare variant associations and found that jointly considering both variants and phenotypes can improve the ability to detect associations. We also applied the MRP model comparison framework to summary statistic data from four traits from the UK Biobank: asthma, eosinophil counts, forced expiratory volume in 1-second (FEV1), and forced vital capacity (FVC). We identified strong evidence for the previously described association between *IL33* and asthma^{30,31} as well as evidence of association for several additional genes that have not been previously reported (Table 1).

As genetic data linked to high-dimensional phenotype data continues to be 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

229 on summary statistics from marginal association analyses which can be shared with less privacy
230 concerns compared to raw genotype and phenotype data. Combining joint analysis of variants and
231 phenotypes with meta-analysis across studies offers new opportunities to identify gene-disease
232 associations.

4 Methods

4.1 MRP model comparison for association testing

We consider the multivariate linear regression model

$$\underset{(N \times K)}{\mathbf{Y}} = \underset{(N \times K)}{\Psi} + \underset{(N \times M)}{\mathbf{X}} \underset{(M \times K)}{\mathbf{B}} + \underset{(N \times K)}{\mathbf{E}},$$

where the matrices $\mathbf{Y} = [y_{ik}]$, $\mathbf{X} = [x_{im}]$, $\mathbf{B} = [\beta_{mk}]$ and $\mathbf{E} = [e_{ik}]$ describe the phenotype values (y_{ik}), copies of minor allele (x_{im}), variant-phenotype effects (β_{mk}), and residual errors (e_{ik}), for individual i , phenotype k , and variant m . We assume that each phenotype has been transformed to a standard normal distribution and that the columns of \mathbf{X} have been centered, which means that the estimate for the intercept term Ψ is 0 and independent of the estimate of \mathbf{B} . We use vectorized notation where the rows of \mathbf{B} form vector $\beta = (\beta_1, \dots, \beta_M)^\top$ of length MK .

We define the MRP model comparison as a Bayes factor (BF) between the alternative model, where at least one variant affects at least one phenotype, and the null model, where all variant-phenotype effects are zero. BF is the ratio of the marginal likelihoods for these two models:

$$\text{BF} = \frac{\int_{\beta} p(\text{Data}|\beta) p(\beta|\text{ALT}) d\beta}{\int_{\beta} p(\text{Data}|\beta) p(\beta|\text{NULL}) d\beta},$$

where Data can correspond either to the effect size estimates $\hat{\beta}$ and the estimated variance-covariance

242 matrix of $\hat{\beta}$, \hat{V}_β , or to the original phenotypes and genotypes, $\underset{(N \times K)}{Y}$ and $\underset{(N \times M)}{X}$, and any other co-
 243 variates that we want to regress out from the phenotypes.

244 The prior distribution for the null model, $p(\beta|\text{NULL})$, is simply the point mass at $\beta = 0$.
 245 In subsection 4.1.1 we show how we approximate the likelihood function for β , $p(\text{Data}|\beta)$, in
 246 subsection 4.1.2 we define the prior distribution $p(\beta|\text{ALT})$ for the alternative model, and finally,
 247 in subsection 4.1.3, we compute the BF.

248 **4.1.1 Likelihood function**

A maximum likelihood estimator of B is given by the ordinary least-squares method

$$\hat{B} = (X^T X)^{-1} X^T Y,$$

that in vectorized form is denoted by $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_M)^T$. An estimator of the variance-covariance
 of $\hat{\beta}$ is given by

$$\hat{V}_\beta = (X^T X)^{-1} \otimes \hat{V}_Y,$$

249 where \hat{V}_Y is the estimated residual variance-covariance matrix of Y given X .

250 Following Band et al.¹⁵, we approximate the likelihood function of β by a multivariate normal
 251 distribution with mean $\hat{\beta}$ and variance-covariance matrix \hat{V}_β . Note that by approximating \hat{V}_Y by
 252 the trait correlation matrix, this likelihood approximation does not require access to the individual

level data \mathbf{X} and \mathbf{Y} but only to the summary data of effect sizes $\hat{\beta}$, LD-matrix $\mathbf{X}^T\mathbf{X}$ and a trait correlation estimate.

4.1.2 Prior of β in the alternative model

We construct the prior distribution $p(\beta|\text{ALT})$ for the alternative model in three steps allowing user to specify correlations between effects of different variants on different traits across different studies.

In a single study, the prior density for β incorporates the expected correlation of genetic effects among a group of variants (\mathbf{R}_{var}) and among a group of phenotypes (\mathbf{R}_{phen}). In addition, we incorporate an expected spread of the effect size of each variant by scaling \mathbf{R}_{var} as

$$\mathbf{S}_{\text{var}} = \Delta(\sigma_m) \mathbf{R}_{\text{var}} \Delta(\sigma_m),$$

where $\Delta(\sigma_m)$ is a diagonal matrix with entries σ_m determining the spread of the effect size distribution for each variant $m \leq M$. Thus, we can model settings where, e.g., protein-truncating variants have larger effect sizes ($\sigma = 1$) than missense variants ($\sigma = 0.1$). Note that when $\sigma_m = 1$ for all m then $\mathbf{S}_{\text{var}} = \mathbf{R}_{\text{var}}$. All in all, our prior density for β under alternative model is

$$\beta|\text{ALT} \sim \mathcal{N}(\mathbf{0}, \mathbf{U}), \text{ where } \mathbf{U} = \mathbf{S}_{\text{var}} \otimes \mathbf{R}_{\text{phen}}.$$

When we have data from multiple studies we allow for possible differences in genetic effects across ethnicities or populations extending the Approximate Bayes Factors of Band et al.¹⁵ and the summary statistics approach of RAREMETAL⁷ from univariate to multivariate phenotypes. Let $\hat{\beta} = (\hat{\beta}_{s,m,k}) = (\hat{\beta}_{1,1,1}, \hat{\beta}_{1,1,2}, \dots, \hat{\beta}_{1,1,K}, \hat{\beta}_{1,2,1}, \dots, \hat{\beta}_{1,2,K}, \dots, \hat{\beta}_{1,M,K}, \hat{\beta}_{2,1,1}, \dots, \hat{\beta}_{S,M,K})$, where S is the number of studies, M is the number of variants, and K is the number of phenotypes. As with a single study, we incorporate the expected correlation of genetic effects between a pair of variants and a single phenotype using the matrix \mathbf{S}_{var} , between a variant and a pair of phenotypes using the matrix \mathbf{R}_{phen} , and we introduce the matrix $\mathbf{R}_{\text{study}}$ to specify prior on the similarity in effect sizes across the studies. Thus, the prior is

$$\beta \sim \mathcal{N}(\mathbf{0}, \mathbf{U}), \text{ where } \mathbf{U} = \mathbf{R}_{\text{study}} \otimes (\mathbf{S}_{\text{var}} \otimes \mathbf{R}_{\text{phen}}).$$

It is straightforward to include a non-zero vector μ as a prior mean of genetic effects, in which case the prior is

$$\beta \sim \mathcal{N}(\mu, \mathbf{U}).$$

259 We use this, for example, when screening for protective rare variants that have a pre-specified
 260 beneficial profile on a set of risk factors.

The Bayes Factor is the ratio of the marginal likelihoods between the alternative and the null model.

The marginal likelihood for the alternative model is

$$\int_{\beta} p(\text{Data}|\beta) p(\beta|\text{ALT}) d\beta = c \times \mathcal{N}(\hat{\beta}; \mu, \hat{V}_{\beta} + \mathbf{U})$$

and the marginal likelihood for the null model is

$$\int_{\beta} p(\text{Data}|\beta) p(\beta|\text{NULL}) d\beta = c \times \mathcal{N}(\hat{\beta}; 0, \hat{V}_{\beta}).$$

In Appendix (5.1) we show how we compute the Bayes Factor

$$\text{BF}_{\text{MRP}} = \frac{\det(\hat{V}_{\beta} + \mathbf{U})^{-\frac{1}{2}} \exp\left[-\frac{1}{2}(\hat{\beta} - \mu)^{\top}(\hat{V}_{\beta} + \mathbf{U})^{-1}(\hat{\beta} - \mu)\right]}{\det(\hat{V}_{\beta})^{-\frac{1}{2}} \exp\left[-\frac{1}{2}\hat{\beta}^{\top}\hat{V}_{\beta}^{-1}\hat{\beta}\right]}.$$

262 When $\mu = 0$, BF_{MRP} is an increasing function of the following quadratic form

$$Q(\hat{\beta}; \hat{V}_{\beta}, \mathbf{U}) = \hat{\beta}^{\top} \left(\hat{V}_{\beta}^{-1} - (\hat{V}_{\beta} + \mathbf{U})^{-1} \right) \hat{\beta}. \quad (1)$$

263 Furthermore, this quadratic form is the only part of the BF_{MRP} that depends on $\hat{\beta}$. Thus, by deriving
 264 a distribution of $Q(\hat{\beta}; \hat{V}_{\beta}, \mathbf{U})$ under the null model we can compute a p-value when BF_{MRP} is
 265 used as a test statistic. According to basic properties of quadratic forms of Gaussian variables,

266 $Q(\hat{\beta}; \hat{V}_\beta, \mathbf{U}) \sim \sum_{i=1}^n d_i \chi_i^2$, where χ_i^2 are an independent sample from χ_1^2 distribution (chi-square
267 with one degree of freedom), and d_i are the eigenvalues of matrix $I - (\hat{V}_\beta + \mathbf{U})^{-1} \hat{V}_\beta$. The
268 distribution function for a mixture of chi-squares can be numerically evaluated by the R-package
269 ‘CompQuadForm’³⁴.

270 4.2 UK Biobank Data

271 4.2.1 GWAS Summary Statistics

272 We performed genome-wide association analysis using PLINK v2.00a(17 July 2017) as previ-
273 ously described³⁰. For asthma, we used the Firth-fallback in PLINK, a hybrid algorithm which
274 normally uses the logistic regression code described in Hill (2017), but switches to a port of lo-
275 gistf() (<https://cran.r-project.org/web/packages/logistf/index.html>) in two cases: (1) one of the cells
276 in the 2x2 allele count by case/control status contingency table is empty (2) logistic regression was
277 attempted since all the contingency table cells were nonzero, but it failed to converge within the
278 usual number of steps. We used the following covariates in our analysis: age, sex, array type, and
279 the first four principal components, where array type is a binary variable that represents whether an
280 individual was genotyped with UK Biobank Axiom Array or UK BiLEVE Axiom Array. For vari-
281 ants that were specific to one array, we did not use array as a covariate. Asthma cases were defined
282 using both Hospital Episode Statistics and verbal questionnaire responses. We used the provided
283 values from the UK Biobank for eosinophil counts, forced vital capacity (FVC), and forced ex-

piratory volume in 1-second (FEV1). The phenotype codes used throughout (asthma=HC382, eosinophil count=INI30150, Forced Expiratory Volume in 1-second (FEV1)=INI3063, and Forced Vital Capacity (FVC)=INI3062) correspond to the phenotype codes used in on the Global Biobank Engine [<https://biobankengine.stanford.edu>].

4.2.2 Genetic Correlations

We calculated the genetic between asthma, eosinophil counts, forced vital capacity (FVC), and forced expiratory volume in 1-second (FEV1) using the MultiVariate Polygenic Mixture Model (MVPMM)³⁵. Briefly, MVPMM estimates genetic correlation given GWAS summary statistics (effect size and standard error of effect size estimate) by modeling GWAS summary statistics as generated from one of two mixture components. Summary statistics from variants in the null component are modeled as being drawn from a multivariate normal distribution with zero mean and covariance matrix that captures correlation in the summary statistics due to the use of shared subjects or other sources of correlation. Summary statistics from variants in the non-null component are modeled as being drawn from a multivariate normal distribution with zero mean, but the covariance matrix for the non-null component combines the covariance matrix from the null component with another covariance matrix that captures the genetic correlation between the phenotypes being considered.

4.2.3 UK Biobank Asthma Application

For the Manhattan plots 3 and table 1, we removed any genes with non-unique gene symbols. We also removed the sense overlapping gene *CTD-3064M3.3* and the antisense gene *ZFAT-ASI* since these genes overlap *GPR20* and *ZFAT* respectively and therefore have the same \log_{10} Bayes Factors as those genes.

4.3 Author Contributions

M.A.R. and M.P. designed the method and derived all analytical calculations. M.A.R., M.P., and C.D. wrote the manuscript. Y.T., M.A., and C.D. provided analysis and designed figures. T.P. designed hDF5 tables and implementation of loaders. M.J.D. provided critical feedback on methodology.

4.4 Acknowledgments

This research has been conducted using the UK Biobank Resource under Application Number 24983. We thank all the participants in the UK Biobank study. M.A.R. and C.D. are supported by Stanford University and a National Institute of Health center for Multi- and Trans-ethnic Mapping of Mendelian and Complex Diseases grant (5U01 HG009080). C.D. is supported by a postdoctoral fellowship from the Stanford Center for Computational, Evolutionary, and Human Genomics

and the Stanford ChEM-H Institute. Y .T. is supported by Funai Overseas Scholarship from Funai Foundation for Information Technology and the Stanford University Biomedical Informatics Training Program. The primary and processed data used to generate the analyses presented here are available in the UK Biobank access management system (<https://amsportal.ukbiobank.ac.uk/>) for application 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>), and the results are displayed in the Global Biobank Engine (<https://biobankengine.stanford.edu>). We would like to thank the Customer Solutions Team from Paradigm4 who helped us implement efficient databases for queries and application of inference methods to the data, and also implemented optimized versions of truncated singular value decomposition. M.A.R. is a paid consultant in Genomics PLC.

5 Appendix

5.1 MRP Bayes Factor computation

To compute the Bayes Factor

$$\text{BF}_{\text{MRP}} = \frac{\det(\hat{\mathbf{V}}_{\beta} + \mathbf{U})^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (\hat{\beta} - \mu)^{\top} (\hat{\mathbf{V}}_{\beta} + \mathbf{U})^{-1} (\hat{\beta} - \mu) \right]}{\det(\hat{\mathbf{V}}_{\beta})^{-\frac{1}{2}} \exp \left[-\frac{1}{2} \hat{\beta}^{\top} \hat{\mathbf{V}}_{\beta}^{-1} \hat{\beta} \right]},$$

we first consider the term inside the exponential function:

$$\mathcal{E}(\hat{\beta}, \mu, \hat{V}_\beta, \mathbf{U}) = \frac{1}{2} \hat{\beta}^\top \hat{V}_\beta^{-1} \hat{\beta} - \frac{1}{2} (\hat{\beta} - \mu)^\top (\hat{V}_\beta + \mathbf{U})^{-1} (\hat{\beta} - \mu).$$

Since \hat{V}_β and \mathbf{U} are typically defined through Kronecker products of smaller matrices, their inverses are easier to compute than the inverse of their sum. Hence we use Woodbury matrix identity to write

$$\mathcal{E}(\hat{\beta}, \mu, \hat{V}_\beta, \mathbf{U}) = \frac{1}{2} \hat{\beta}^\top \hat{V}_\beta^{-1} \hat{\beta} - \frac{1}{2} (\hat{\beta} - \mu)^\top \left(\hat{V}_\beta^{-1} - \hat{V}_\beta^{-1} (\mathbf{U}^{-1} + \hat{V}_\beta^{-1})^{-1} \hat{V}_\beta^{-1} \right) (\hat{\beta} - \mu).$$

To simplify the determinant calculation we write

$$\det(\hat{V}_\beta + \mathbf{U}) = \det(\hat{V}_\beta) \det(\mathbf{I} + \hat{V}_\beta^{-1} \mathbf{U}).$$

The logarithm of the Bayes Factor is then

$$\log(\text{BF}_{\text{MRP}}) = -\frac{1}{2} \log(\det(\mathbf{I} + \hat{V}_\beta^{-1} \mathbf{U})) + \mathcal{E}(\hat{\beta}, \mu, \hat{V}_\beta, \mathbf{U}).$$

If studies do not share individuals, $\hat{\mathbf{V}}_{\beta}$ is a block-diagonal matrix

$$\hat{\mathbf{V}}_{\beta} = \begin{bmatrix} \hat{\mathbf{V}}_{\beta}^1 & 0 & \cdots & 0 \\ 0 & \hat{\mathbf{V}}_{\beta}^2 & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \hat{\mathbf{V}}_{\beta}^S \end{bmatrix}.$$

If studies share individuals, e.g., controls, we take the approach of Cichonska et al.⁸ to use summary level data to estimate the correlation structure of the non-diagonal blocks caused by overlapping individuals.

5.2 Implementation details: HDF5 tables

Although summary statistics are quicker to read and process than raw data, the number of studies meta-analyzed in this work is expected to be sufficiently large to require optimizations in data representation and processing (Figure S1). Our solution was the use of the HDF5 (Hierarchical Data Format 5) data representation to enable rapid processing of effect size, uncertainty, and cross-trait estimate data. HDF5 is a fast and lightweight file format designed for scientific data. It has bindings for R, Python, C/C++, Java, and nearly every other population programming language. Reading data from a table within a HDF5 file can be an order of magnitude faster than reading text files from a Unix file, and it makes it easier to organize data within an internal structure.

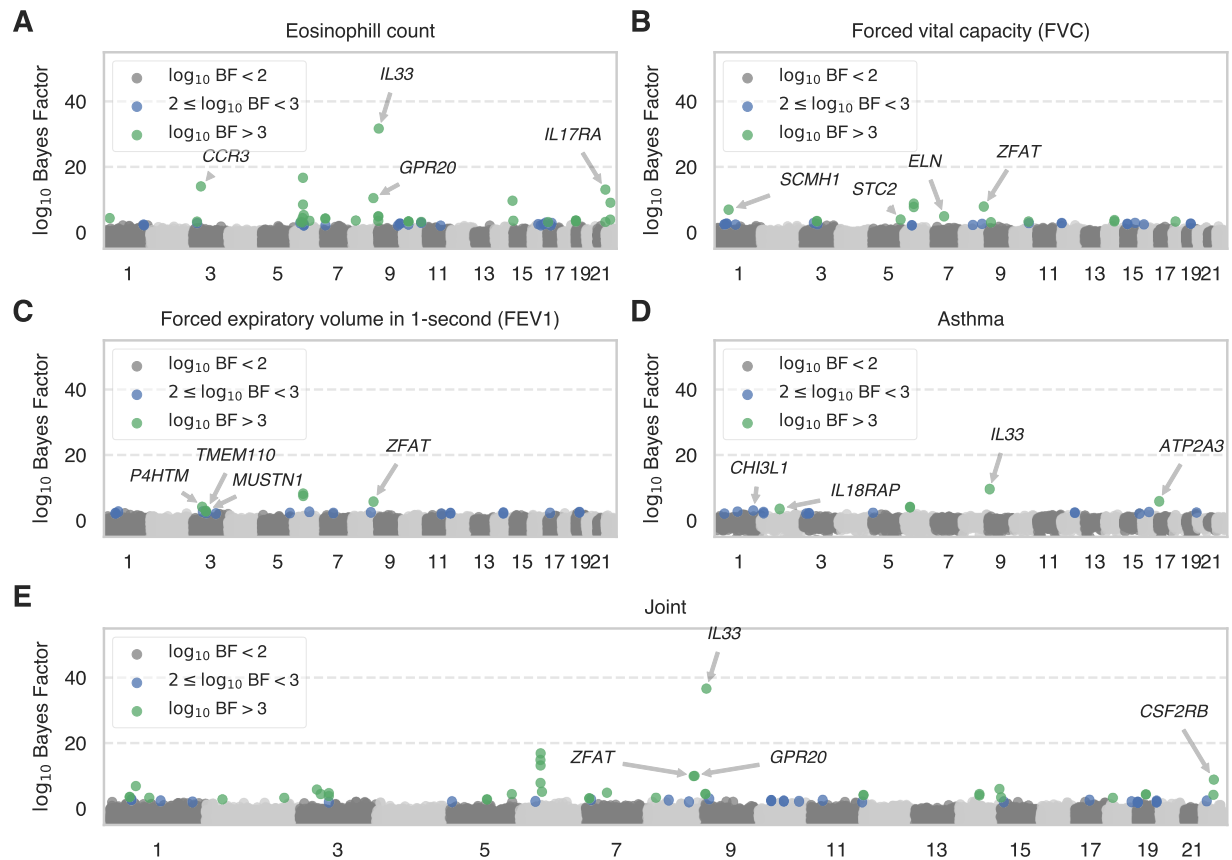


Figure 4 \log_{10} Bayes Factors from applying MRP to summary statistics for missense and protein-truncating variants from (A) eosinophil counts (INI30150), (B) forced vital capacity (FVC, INI3062), (C) forced expiratory volume in 1-second (FEV1, INI3063), (D) asthma (HC382), and (E) all four traits jointly. The four genes outside of chromosome 6 with the largest Bayes Factors are labeled in each plot. Only \log_{10} Bayes Factors greater than -5 are plotted.

	Joint	Eosinophil count	FVC	FEV1	Asthma
<i>IL33</i>	36.6	31.5	-2.3	-1.9	9.4
<i>GPR20</i>	10.2	10.4	-0.1	-0.2	
<i>ZFAT</i>	10.1	-3.1	7.8	5.9	-0.5
<i>CSF2RB</i>	9.0	9.1	0.2	-0.1	-0.2
<i>SCMH1</i>	7.1	-1.5	6.8	2.6	-0.9
<i>PACS2</i>	6.2	9.6	-1.3	-1.6	-0.5
<i>CCR3</i>	6.0	14.0	-3.3	-3.7	-1.0
<i>ELN</i>	5.0	-1.5	4.9	2.3	-0.6

Table 1: \log_{10} Bayes Factors for genes outside of chromosome 6 with \log_{10} Bayes Factors greater than 5 in the joint analysis.

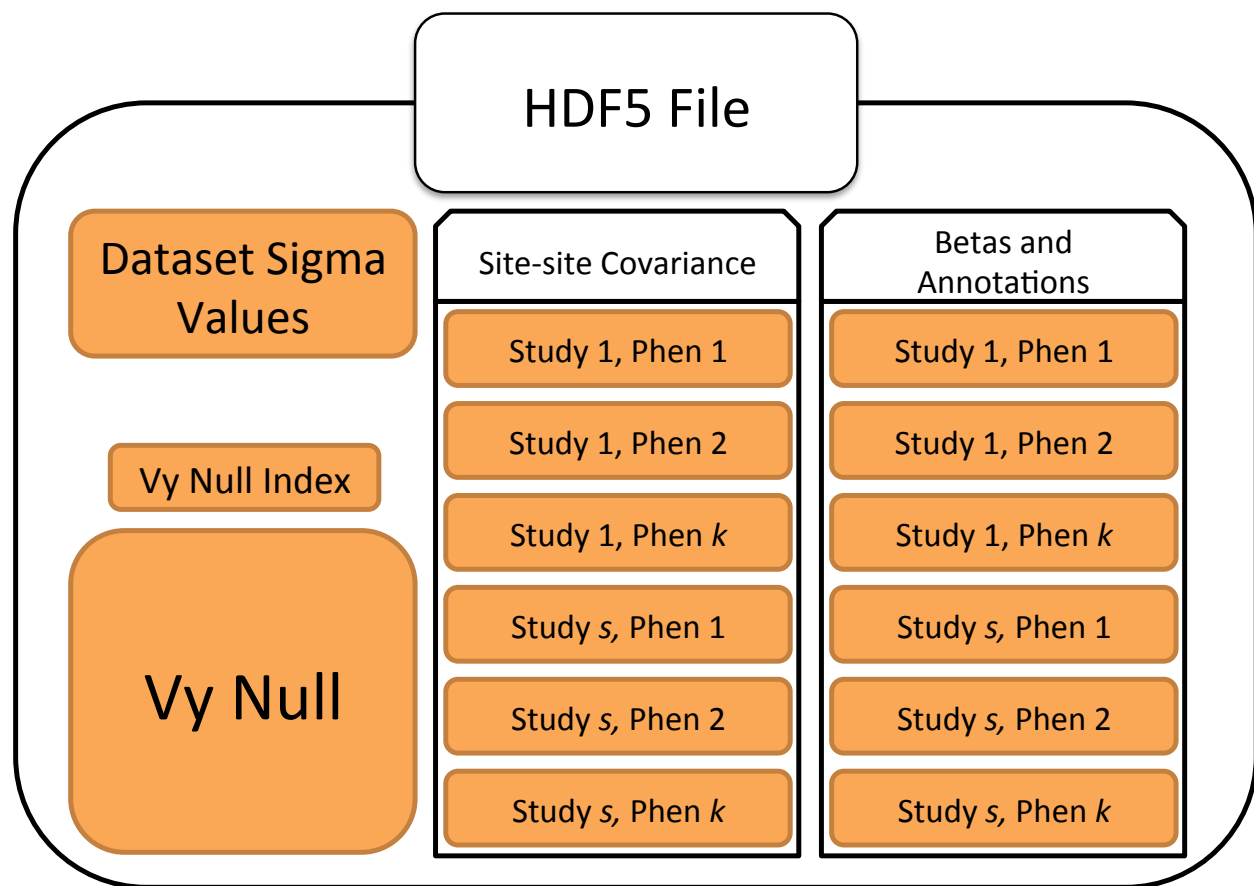


Figure 5 Our HDF5 implementation contained the following components: first, a group with one table per annotation file. All effect size (beta) values and study-specific annota-

352 tions were contained here, and the number of tables is limited by S (the number of studies) \times
 353 K (the number of traits). Second, a group with site-site covariance data. While these
 354 covariance matrices may have dimension M (the number of variants) $\times M$, we store the
 355 data as tables, each row specifying the covariance between two variants. The number
 356 of tables should be the same as the previous set, capped by S (the number of studies) \times
 357 K (the number of traits). Third, we store one table with sigma values for each study/phe-
 358 notype combination. In the event that the traits were rank-normal transformation was
 359 performed these sigma values are equal to 1. These are used to compute correlation
 360 between two datasets. Finally, we store a matrix/table pair for V_y null and its index. The
 361 V_y null matrix has dimensions $(S \times K) \times (S \times K)$ each entry specifying the estimated
 362 correlation of effect sizes between two datasets. The index table encodes row/column
 363 position of each dataset.

1. Nejentsev, S., Walker, N., Riches, D., Egholm, M. & Todd, J. A. Rare variants of *IFIH1*, a gene implicated in antiviral responses, protect against type 1 diabetes. *Science* **324**, 387–389 (2009).
2. 1000 Genomes Project Consortium. A map of human genome variation from population scale sequencing. *Nature* **467**, 1061–1073 (2010).
3. Rivas, M. A. *et al.* Deep resequencing of GWAS loci identifies independent rare variants associated with inflammatory bowel disease. *Nature Genetics* **43**, 1066–1073 (2011).
4. The 1000 Genomes Consortium. An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56–65 (2012).
5. Majithia, A. R. *et al.* Rare variants in *PPARG* with decreased activity in adipocyte differentiation are associated with increased risk of type 2 diabetes. *Proceedings of the National Academy of Sciences* **111**, 13127–13132 (2014).
6. Lee, S., Abecasis, G. R., Boehnke, M. & Lin, X. Rare-variant association analysis: study designs and statistical tests. *The American Journal of Human Genetics* **95**, 5–23 (2014).
7. Liu, D. J. *et al.* Meta-analysis of gene-level tests for rare variant association. *Nature Genetics* **46**, 200–204 (2014).
8. Cichonska, A. *et al.* metacca: summary statistics-based multivariate meta-analysis of genome-wide association studies using canonical correlation analysis. *Bioinformatics* **32**, 1981–1989 (2016).

- 384 9. Rivas, M. A. *et al.* Assessing association between protein truncating variants and quantitative
385 traits. *Bioinformatics* **29**, 2419–2426 (2013).
- 386 10. Rivas, M. A. *et al.* Effect of predicted protein-truncating genetic variants on the human tran-
387 scriptome. *Science* **348**, 666–669 (2015).
- 388 11. Cohen, J. *et al.* Low LDL cholesterol in individuals of African descent resulting from frequent
389 nonsense mutations in *PCSK9*. *Nature Genetics* **37**, 161–5 (2005).
- 390 12. Cohen, J. C., Boerwinkle, E., Mosley Jr, T. H. & Hobbs, H. H. Sequence variations in *PCSK9*,
391 low LDL, and protection against coronary heart disease. *New England Journal of Medicine*
392 **354**, 1264–1272 (2006).
- 393 13. Sullivan, D. *et al.* Effect of a monoclonal antibody to PCSK9 on low-density lipoprotein
394 cholesterol levels in statin-intolerant patients: the GAUSS randomized trial. *JAMA* **308**, 2497–
395 2506 (2012).
- 396 14. Wu, M. C. *et al.* Rare-variant association testing for sequencing data with the sequence kernel
397 association test. *American Journal of Human Genetics* **89**, 82–93 (2011).
- 398 15. Band, G. *et al.* Imputation-based meta-analysis of severe malaria in three african populations.
399 *PLoS genetics* **9**, e1003509 (2013).
- 400 16. MacArthur, D. G. *et al.* A systematic survey of loss-of-function variants in human protein-
401 coding genes. *Science* **335**, 823–8 (2012).

- 402 17. Findlay, G. M., Boyle, E. A., Hause, R. J., Klein, J. C. & Shendure, J. Saturation editing of
403 genomic regions by multiplex homology-directed repair. *Nature* **513**, 120–123 (2014).
- 404 18. Neale, B. M. *et al.* Testing for an unusual distribution of rare variants. *PLoS Genetics* **7**,
405 e1001322 (2011).
- 406 19. Clarke, G. M., Rivas, M. A. & Morris, A. P. A flexible approach for the analysis of rare
407 variants allowing for a mixture of effects on binary or quantitative traits. *PLoS Genetics* **9**,
408 e1003694 (2013).
- 409 20. Cotsapas, C. *et al.* Pervasive sharing of genetic effects in autoimmune disease. *PLoS genetics*
410 **7**, e1002254 (2011).
- 411 21. Solovieff, N., Cotsapas, C., Lee, P. H., Purcell, S. M. & Smoller, J. W. Pleiotropy in complex
412 traits: challenges and strategies. *Nature Reviews Genetics* **14**, 483–495 (2013).
- 413 22. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and traits.
414 *Nature genetics* (2015).
- 415 23. Do, R. *et al.* Exome sequencing identifies rare *ldlr* and *apoa5* alleles conferring risk for my-
416 ocardial infarction. *Nature* **518**, 102–106 (2015).
- 417 24. Cohen, J. *et al.* Low LDL cholesterol in individuals of African descent resulting from frequent
418 nonsense mutations in *pcsk9*. *Nature Genetics* **37**, 161–165 (2005).
- 419 25. Do, R. *et al.* Common variants associated with plasma triglycerides and risk for coronary
420 artery disease. *Nature Genetics* **45**, 1345–1352 (2013).

26. The TG and HDL Working Group of the Exome Sequencing Project, National Heart, Lung, and Blood Institute. Loss-of-function mutations in *APOC3*, triglycerides, and coronary disease. *New England Journal of Medicine* **371**, 22–31 (2014).
27. Pollin, T. I. *et al.* A null mutation in human *APOC3* confers a favorable plasma lipid profile and apparent cardioprotection. *Science* **322**, 1702–5 (2008).
28. Hofker, M. H. *APOC3* null mutation affects lipoprotein profile *APOC3* deficiency: from mice to man. *European Journal of Human Genetics* **18**, 1–2 (2010).
29. Jørgensen, A. B., Frikke-Schmidt, R., Nordestgaard, B. G. & Tybjaerg-Hansen, A. Loss-of-function mutations in *APOC3* and risk of ischemic vascular disease. *New England Journal of Medicine* **371**, 32–41 (2014).
30. DeBoever, C. *et al.* Medical relevance of protein-truncating variants across 337,208 individuals in the uk biobank study. *bioRxiv* (2017). URL <https://www.biorxiv.org/content/early/2017/09/02/179762>.
31. Smith, D. *et al.* A rare *il33* loss-of-function mutation reduces blood eosinophil counts and protects from asthma. *PLOS Genetics* **13**, 1–24 (2017). URL <https://doi.org/10.1371/journal.pgen.1006659>.
32. Ferreira, M. A. *et al.* Shared genetic origin of asthma, hay fever and eczema elucidates allergic disease biology. *Nature genetics* **49**, 1752–1757 (2017).
33. Ober, C. & Yao, T.-C. The genetics of asthma and allergic disease: a 21st century perspective. *Immunological reviews* **242**, 10–30 (2011).

- 441 34. Duchesne, P. & de Micheaux, P. L. Computing the distribution of quadratic forms: Further
442 comparisons between the liu-tang-zhang approximation and exact methods. *Computational*
443 *Statistics and Data Analysis* **54**, 858–862 (2010).
- 444 35. DeBoever, C. & Rivas, M. A. Harnessing digital phenotyping to enhance genetic studies of
445 human diseases. *Submitted* (2017).