# CoPheScan: phenome-wide association studies accounting for linkage disequilibrium

Ichcha Manipur[1,2], Guillermo Reales[1,2], Jae Hoon Sul[3], Myung Kyun Shin[3], Simonne Longerich[3], Adrian Cortes[4], Chris Wallace[1,2,5]

[1]Cambridge Institute of Therapeutic Immunology & Infectious Disease (CITIID), Jeffrey Cheah Biomedical Centre, Cambridge Biomedical Campus, University of Cambridge, Cambridge CB2 0AW, UK
[2]Department of Medicine, University of Cambridge School of Clinical Medicine, Cambridge Biomedical Campus, Cambridge CB2 2QQ, UK
[3]Merck & Co., Inc., Rahway, NJ, USA
[4]Human Genetics and Genomics, GSK, Heidelberg, 69117, Germany
[5]MRC Biostatistics Unit, University of Cambridge, Cambridge, United Kingdom

## Abstract

Phenome-wide association studies (PheWAS) facilitate the discovery of associations between a single genetic variant with multiple phenotypes. For variants which impact a specific protein, this can help identify additional therapeutic indications or on-target side effects of intervening on that protein. However, PheWAS is restricted by an inability to distinguish confounding due to linkage disequilibrium (LD) from true pleiotropy. Here we describe CoPheScan (Coloc adapted Phenome-wide Scan), a Bayesian approach that enables an intuitive and systematic exploration of causal associations while simultaneously addressing LD confounding. We demonstrate its performance through simulation, showing considerably better control of false positive rates than a conventional approach not accounting for LD. We used CoPheScan to perform PheWAS of protein-truncating variants and fine-mapped variants from disease and pQTL studies, in 2275 disease phenotypes from the UK Biobank. Our results identify the complexity of known pleiotropic genes such as *APOE*, and suggest a new causal role for *TGM3* in skin cancer.

## Main

Phenome-wide association studies (PheWAS) are an inversion of the GWAS (Genome-Wide Association Studies) paradigm, where a single genetic variant is tested against a broad range of phenotypes. Phenome scale studies are facilitated by the availability of a broad array of phenotypes linked to genomic data in large-scale biobanks. PheWAS are a promising tool in the field of pharmacogenomics as they facilitate drug repurposing efforts and identification of potential adverse effects due to their ability to detect pleiotropy [1–3]. Often, PheWAS has been paired with other approaches such as Mendelian Randomisation to identify causal effects of exposures on outcomes and network analysis to identify interactions between phenotypes [4–6].

Prevailing methods for phenome-wide testing are built upon single variant tests and do not inherently tackle the spurious associations that can arise when traits are causally associated not with the index variant, but with another variant in LD with the index variant. For instance, a PheWAS of UK Biobank phenotypes with protein-truncating variants by DeBoever et al.[7] first revealed an association between an *ANKDD1B* variant, and high cholesterol, which was found

46  to reflect an indirect association, through LD with an intronic variant in *HMGCR* which is known
47  to be associated with cholesterol levels. Thus LD confounding necessitates the use of
48  additional follow-up tests such as colocalisation analyses, where pairs of traits are tested for
49  shared causal variants within a genomic region, to isolate associations that are truly causal
50  [3,8].

51  PheWAS hits are colocalised with molecular QTLs or disease traits on which the identified
52  variants have a prior known effect. However, this two-step approach is not feasible for variants
53  with known biological effects for which summary statistics are unavailable, such as those
54  involved in protein truncation.

55

56  In this work, we introduce a Bayesian approach to PheWAS, Coloc adapted Phenome-wide
57  Scan, (CoPheScan), that tests phenome-scale causal associations with a set of index variants
58  while handling confounding due to LD at the same time. CoPheScan can exploit external
59  covariate data, such as the genetic correlation between phenotypes, and can be run in
60  different ways depending on whether accurate LD information is available and whether the
61  analyst is prepared to make assumptions about the number of causal variants in the tested
62  genomic region. We demonstrate the utility and robustness of these different approaches on
63  simulated datasets. We also analysed causal variants selected from three real-world sources
64  and tested for causal associations against 2275 phenotypes from the UK Biobank using
65  CoPheScan.

66

## 67  Results

### 68  Overview of CoPheScan

69  CoPheScan is an adaptation of the coloc [9–11] approach, for the case where a variant known
70  to be causal either through fine-mapping or functional studies, is subjected to a phenome-wide
71  scan to test for causal associations with other phenotypes/traits. Coloc considers the genetic
72  association patterns for two traits in a genomic region and assesses whether it is likely they
73  share a causal variant in that region. It is a Bayesian approach and assumes prior probabilities
74  for each of the five possible hypotheses (no association with either trait, association with just
75  one trait or the other, association with both traits and different causal SNPs, or association
76  with both traits at the same causal SNP) are fixed and known.

77

78  We consider the case where a SNP of interest is known to be causal for a phenotype which is
79  often the case in PheWAS, and we are interested in determining if it is also causally associated
80  with another phenotype (Figure 1a). We will hereafter refer to the variant of interest as the
81  query variant, the phenotype for which the query variant is known to be causally associated
82  as the primary trait and the phenotype to be tested as the query trait. In a genomic region with
83  Q SNPs, and under the initial assumption of a single causal variant (which we will relax later),
84  there are $Q+1$ possible ways or "configurations", (Supplementary Figure 1) to describe where
85  the single causal variant may lie, each corresponding to exactly one of three hypotheses:
86  $H_n$: No association of any variant with the query trait (one configuration)
87  $H_a$: Causal association of a variant other than the query variant with the query trait ($Q-1$
88  configurations)
89  $H_c$: Causal association of the query variant with the query trait (one configuration)
90

91     The posterior odds for each hypothesis ($H$) given the data ($D$) for the query trait with respect
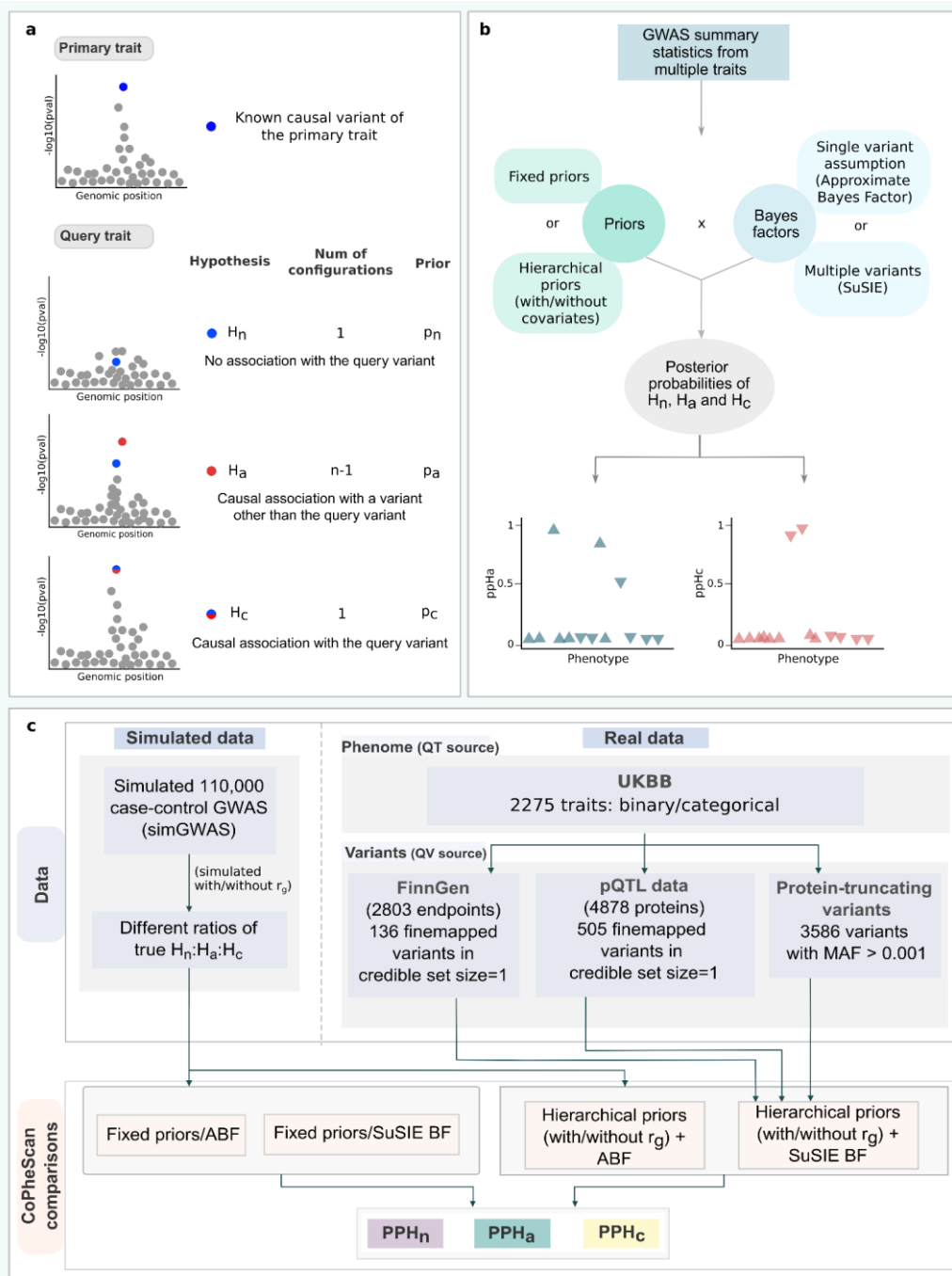92     to the null hypothesis ($H_n$) is given by,
93

$$\frac{P(H|D)}{P(H_n|D)} = \frac{P(H)}{P(H_n)} \times \frac{P(D|H)}{P(D|H_n)} \tag{1}$$

94

95     In equation (1), the first ratio in the right-hand side is the prior odds and the second ratio is the
96     Bayes Factor (BF). Thus, the three prior probabilities that have to be specified are: $p_n = P(H_n)$,
97     $p_a = P(H_a)/(Q-1)$, and $p_c = P(H_c)$, subject to the constraint that $p_n + (Q-1)p_a + p_c = 1$.
98     Beyond the difference in the hypothesis space described above, CoPheScan differs from coloc
99     in two further ways. First, because we have reduced the hypothesis space, we can examine
100     many variants simultaneously, allowing us to learn the priors from the data in a hierarchical
101     Bayesian manner with Markov Chain Monte Carlo (MCMC) sampling (Supplementary
102     methods). In contrast, coloc assumes priors are fixed and known, which is a weakness
103     because inference must rely on the investigators' judgement on prior probabilities of
104     colocalisation. Second, because we are using this hierarchical approach, we can exploit
105     additional external information about the variants and/or the traits in the form of covariates
106     which can be included when learning the priors. This allows the priors to vary depending on
107     the query trait/query variant pairs being considered. Here, we include the genetic correlation
108     ($r_g$) between the primary trait and each query trait tested (see Supplementary Methods).
109

110     The restriction to a single causal variant allows us to count the possible configurations ($Q+1$),
111     and if the assumption is deemed valid, CoPheScan can be run directly on summary GWAS
112     data using Wakefield's method[12], to compute approximate Bayes factors summarising the
113     relative support for a model where the SNP is associated with a trait compared to the null
114     model of no association. However, this assumption is not broadly valid, and an alternative is
115     to use the Sum of Single Effects (SuSiE) Bayesian fine mapping regression framework[13,14] to
116     partition the evidence into configurations corresponding to each of multiple possible causal
117     variants and use these in a similar manner to allowing for multiple causal variants in coloc[10].
118     The SuSiE approach works best with either raw genotype data or summary GWAS data when
119     in-sample LD information is available[15].
120

121     Hence, CoPheScan has the flexibility to be run in several ways (Figure 1b) depending on: (i)
122     the assumption about the number of causal variants, (ii) the specification of either fixed or
123     hierarchical priors, and (iii) the inclusion/exclusion of covariates if the hierarchical model is
124     used to infer priors. A detailed description of the CoPheScan method is available in the
125     Supplementary methods. A summary of the simulated data, variant and phenotype sources
126     used for the analysis with the real data can be found in (Figure 1c), while a detailed description
127     is provided in the Methods.

**Figure 1: Introduction and evaluation of the CoPheScan method.**



(**a**) CoPheScan methodology: Hypotheses with illustrations of the configurations of genetic variants within the genomic region and corresponding priors. (**b**) Schematic of the CoPheScan workflow. The inputs are GWAS summary statistics from multiple traits and the position of the query variant. Computation of the posterior probabilities of the three hypotheses is performed with priors and Bayes factors computed using different CoPheScan approaches. (**c**) Study design for evaluation: Simulated data - Generated using SimGWAS and all CoPheScan approaches were run on this set. Real data - Phenotypes tested were obtained from UK Biobank and variants from fine-mapping FinnGen and a proteome dataset[16]. Hierarchical priors and SuSIE BF were used on the real data to identify SNP-disease associations. (QV - query variant, QT - query trait)

128 **Simulations show CoPheScan is more accurate than a standard method which**
129 **does not account for LD confounding**

130 We simulated regional GWAS summary data for traits with either zero, one or two causal
131 variants (Methods) such that they corresponded to the three CoPheScan hypotheses. We also
132 allowed the probability of Hc to vary according to a simulated genetic covariance between
133 primary and query traits and considered whether including this information in the analysis
134 increased inferential accuracy. We analysed the same data in parallel using a conventional
135 PheWAS approach of testing each of the set of query SNPs for association, controlling either
136 the FDR or the family-wise error rate via Bonferroni correction. We compared these to the
137 results from CoPheScan, using a hierarchical model (with and without the covariate data) or
138 fixed priors chosen as described in the Supplementary methods which broadly matched the
139 proportion of Hn, Ha, and Hc in the sample.

140

141 First, we considered the appropriate threshold on the posterior probability of Hc, ppHc, to call
142 an association. We estimated the FDR internally, as 1-mean(ppHc) | ppHc > t for different
143 values of threshold t (Supplementary Figure 4). We found that ppHc > 0.6 maintained an FDR
144 < 0.05 across all analyses of simulated data. Using this threshold, CoPheScan appeared less
145 sensitive to the presence of a single causal variant (true Hc) than the conventional BH
146 approach but more sensitive than the Bonferroni approach (Figure 2). CoPheScan
147 demonstrated control of the FDR (0.026-0.039) estimated as the proportion of significant calls
148 that were truly Hn or Ha , traits where the query variant was not causal, for the different
149 CoPheScan approaches compared to   0.219 and 0.308 for the conventional BH and
150 Bonferroni approaches respectively, (Supplementary Table 1). The majority of the false
151 positives obtained from these conventional approaches were true Ha but called as associated
152 due to LD confounding. All CoPheScan approaches performed well in the case of a single
153 causal variant, but when there were two causal variants (True Hc2), using SuSIE resulted in
154 approximately 30% higher sensitivity to correct Hc predictions than the ABF approach (Figure
155 2). This was balanced against marginally lower (<0.5%) sensitivity to Hc with SuSiE when
156 traits truly had only a single causal variant (True Hc) when compared to the CoPheScan
157 approaches that assumed a single causal variant.

158
159
160
161
162
163
164
165
166
167
168
169
170
171

**Figure 2: Results for hypotheses discrimination in simulated data.**

We called a single result for each simulated trait as described in Methods. The x axis shows the percentage of hypothesis calls using the different approaches shown on the y axis. For CoPheScan (top 6 rows), the three labelled columns on the y-axis, from right to left, indicate the type of priors used, the method used to calculate Bayes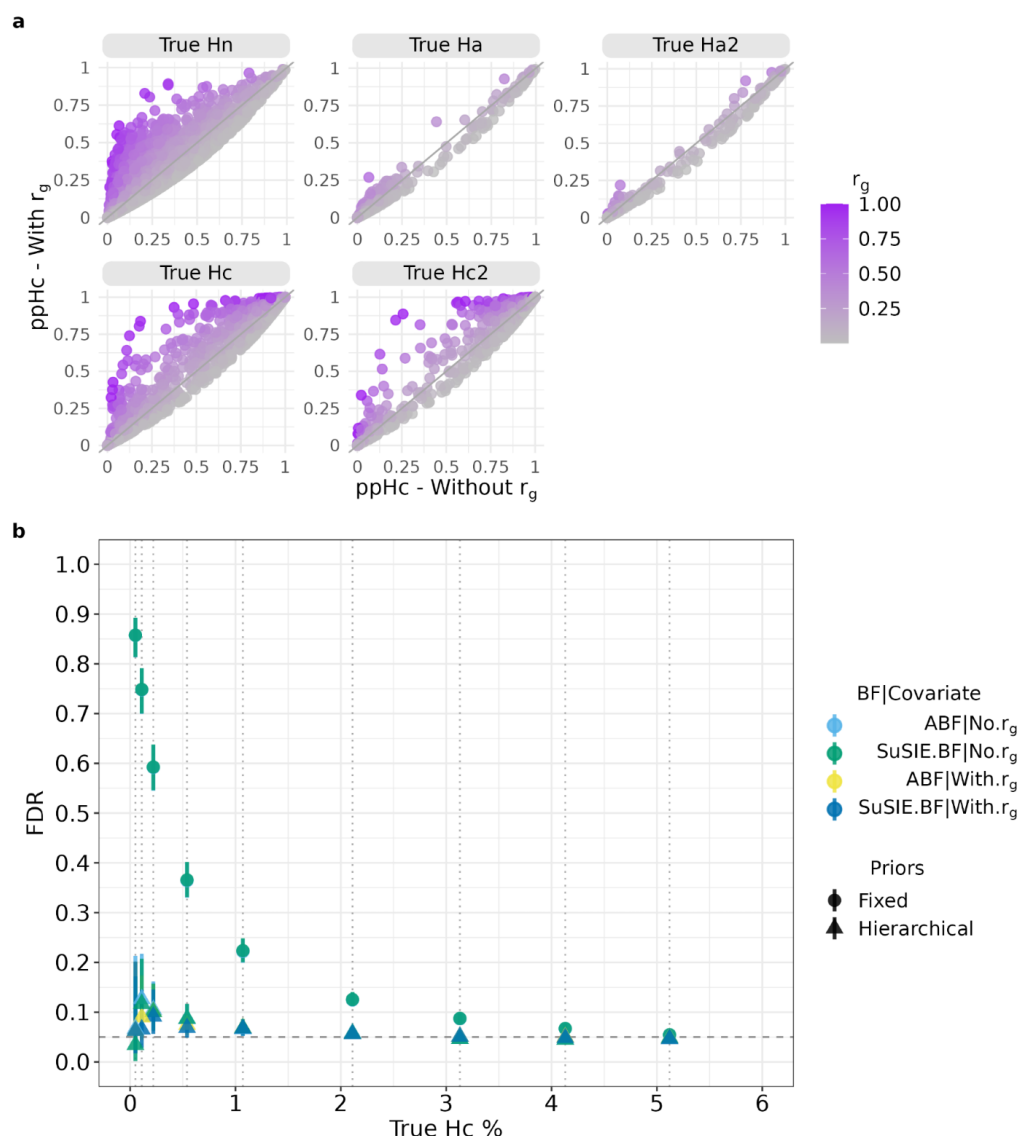 factors, and whether or not genetic correlation (rWe called a single result for each simulated trait as described in Methods. The x axis shows the percentage of hypothesis calls using the different approaches shown on the y axis. For CoPheScan (top 6 rows), the three labelled columns on the y-axis, from right to left, indicate the type of priors used, the method used to calculate Bayes factors, and whether or not genetic correlation ($r_g$) was used. The last two rows show conventional approaches controlling the FDR (BH - Benjamini-Hochberg) or the FWER (Bonf - Bonferroni) at 0.05. The top bar shows an illustration of the configuration of SNPs in the genomic region corresponding to the different simulated traits (Methods), with the queried variant at position 1 and causally associated (non-associated) variants indicated by filled (open) circles. [True Hn: no causal variant, True Ha/Ha2: one/two causal non-query variants, True Hc: causal query variant, True Hc2: causal query variant and one causal non-query variant].

172

173     Although the effect of including covariate information was minor overall, Figure 3a shows that it
174     had a substantial effect in a minority of cases, bringing ppHc from below to above 0.6 in 2.79%
175     of true Hc and Hc2 cases (80/2867), although also in 0.088% of true Hn and 0.011% of true Ha
176     and Ha2.
177     Finally, these initial simulations showed that the hierarchical model recovered very similar
178     results to the fixed prior model, where we chose our fixed prior values to broadly match the
179     simulation scenarios, i.e., an optimal scenario. This offers reassurance that the hierarchical
180     model can perform just as well as a method that "knows" the correct prior values. However, in
181     real data, we will not know the true proportion of Hn, Hc, or Ha in our data, so we explored the
182     robustness of both approaches to variations in these proportions. We found that using over-

6

183 optimistic fixed priors, i.e. when the prior probability for Hc (P(Hc)=0.091) exceeded the
184 proportion of Hc in our data, led to dramatically high FDR, whilst the hierarchical model
185 correctly adapted to the different datasets so that the FDR was controlled except at the very
186 lowest true proportions of Hc (Figure 3b).
187



**Figure 3: Effects of covariate inclusion and varying proportions of simulated hypotheses.**
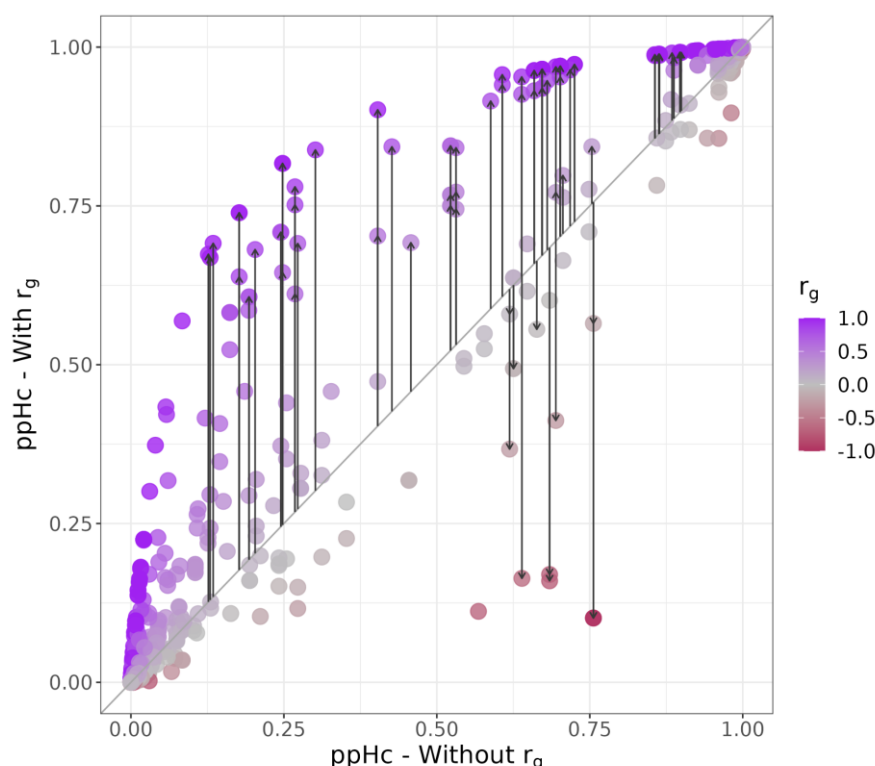
(**a**) Comparison of the posterior probability of Hc (ppHc) obtained with (y axis) and without (x axis) the inclusion of genetic correlation ($r_g$) in the hierarchical model (using ABF). The panels represent the traits of different simulated hypotheses (Methods). (**b**) The proportion of Hc traits was varied as shown in the x axis (dotted vertical lines), to compare Hc predictions using the fixed and hierarchical priors with different BF, both with and without the inclusion of the genetic correlation ($r_g$) covariate. The y axis represents the estimated FDR - the proportion of traits assigned as Hc in each dataset which were simulated as Hn or Ha with 95% confidence intervals (dashed line - 0.05 FDR).

**Using genetic correlation as a covariate increases detection of associations with disease-causal variants**

We explored the performance of CoPheScan (Supplementary Figure 5) using a variety of causal variants sets to perform PheWAS in three sets of query variants in up to 2275 query traits (Supplementary Figures 6 and 7) from the UK Biobank summary data provided by the Neale Lab (http://www.nealelab.is/uk-biobank/). First, 136 disease-causal variants were identified as single variant credible sets in fine mapping data from FinnGen disease endpoints (primary traits, https://www.finngen.fi/en/access_results). We identified causal associations in UKBB at 43 (31.62%) of these, predominantly amongst query traits identical or related to the primary trait. Out of 101 unique query-variant-primary trait pairs with exact query variant-query trait matched pairs in UKBB, 32 were found to be Hc (Supplementary Figure 7), and 65 Hn due to a lack of power in UKBB (p-value> $10^{-5}$). Four cases were called Ha, and in these the UKBB p value was small, but the fine mapping produced different results in UKBB and FinnGen (Supplementary Figure 8).

Genetic correlation information ($r_g$) for only 1582 out of the 2275 traits used in analysis without $r_g$ was available. $r_g$ values between the 1582 query traits and 69 UKBB traits which were matched with the FinnGen primary traits were used as a covariate (130697 query trait-query variant pairs tested). Including $r_g$ in the hierarchical model made a larger difference here than in the simulated data, perhaps reflecting a stronger effect than we anticipated in our simulations. Overall, ppHc values for traits with higher $r_g$ with the primary traits increased and, conversely, decreased for traits with lower (negative) $r_g$ (Figure 4). Incorporating the $r_g$ resulted in the identification of 19 additional associations (Supplementary Table 8). For example, the variants rs3217893_C>T and rs2476601_A>G, fine-mapped for type 2 diabetes and rheumatoid arthritis (RA) in FinnGen respectively, were found to have associations with medications gliclazide, which is a sulfonylurea used in the treatment of Type 2 diabetes, and steroid prednisolone which can be used to treat RA, only when the genetic correlation information was included.
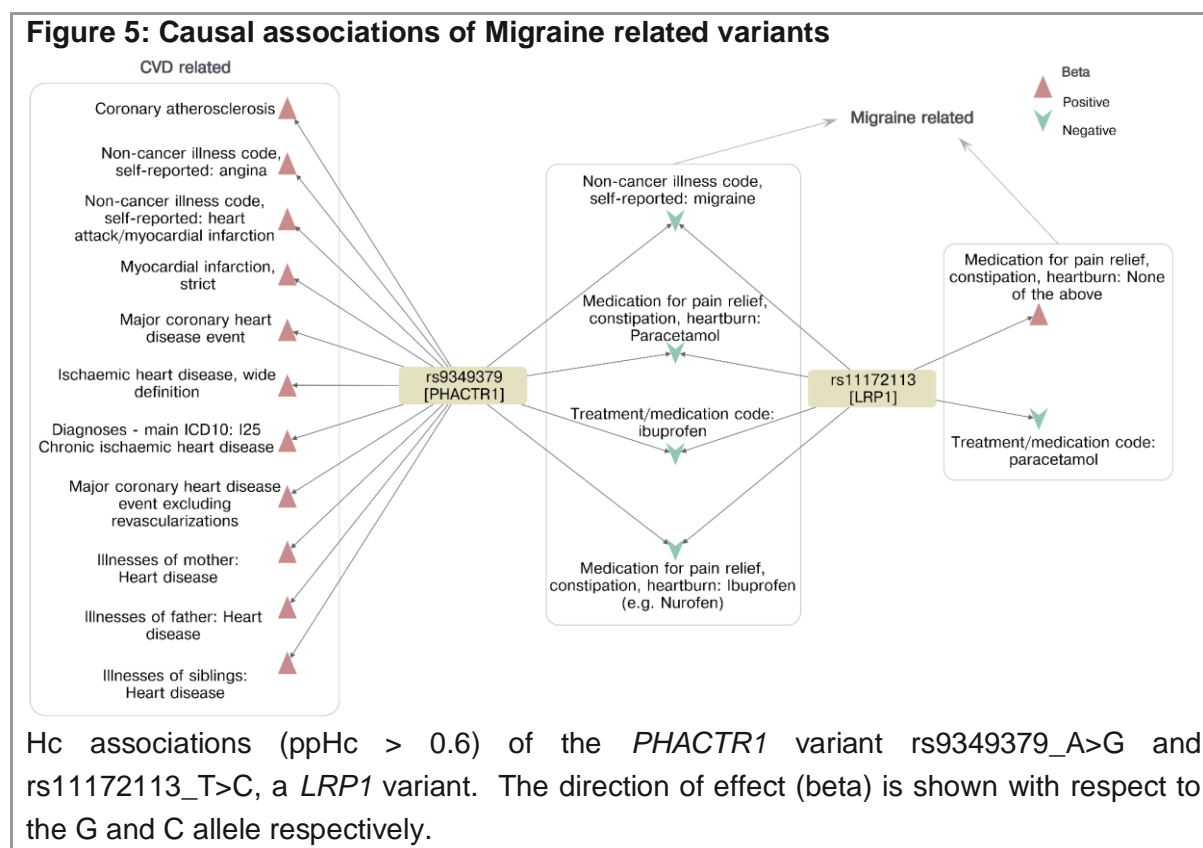
**Figure 4: Genetic correlation detects additional phenotypes**

Hierarchical models of the FinnGen/UKBB dataset with/without genetic correlation ($r_g$). The posterior probability of Hc (ppHc) of traits with and without the inclusion of genetic correlation ($r_g$) are shown on the y and x axes respectively. The arrows represent the traits which show a difference of > 0.1 ppHc after inclusion of $r_g$ (compared to the model without) and also have a ppHc > 0.6. The traits are coloured to represent their $r_g$ with the primary trait.

Query variants were often associated with multiple UKBB traits (median 5) that reflected related diseases and medications (Supplementary Table 8). For instance, rs11591147_G>T, a missense variant of *PCSK9*, identified as a disease-causal variant in FinnGen for statin medication was found associated with the UKBB traits related to different statin medications along with several cardiovascular traits. Less commonly, we found evidence for causal association of variants to seemingly unrelated traits. For example, rs9349379_A>G, an intron variant and eQTL for *PHACTR1*, identified by fine-mapping the FinnGen primary trait - triptan, which is a medication used to manage migraine, was found to be associated with several UKBB traits related to migraine such as the phenotype itself, migraine medications such as sumatriptan, ibuprofen and paracetamol and also the presence of family history. However, we also found associations with angina, myocardial infarction and ischaemic heart disease, with the migraine-protective allele acting as a risk factor for cardiovascular traits. This matches results from a Mendelian randomisation study of migraine and cardiovascular disease[17] but is in contrast to observational studies where migraine is considered positively associated with cardiovascular traits[18]. Such discrepancies between genetic and observational studies in other traits have often been resolved in favour of the genetic result, through the identification of some confounding factor which led the observational studies to report inverse relationships, and it has been suggested that certain non-triptan migraine therapies might act to increase cardiovascular risk[17]. However, this pleiotropy did not appear at another migraine-identified

250  variant, rs11172113_T>C, an intronic variant of *LRP1*, which was fine-mapped for the same
251  FinnGen primary trait of migraine, and found to be independently associated with several
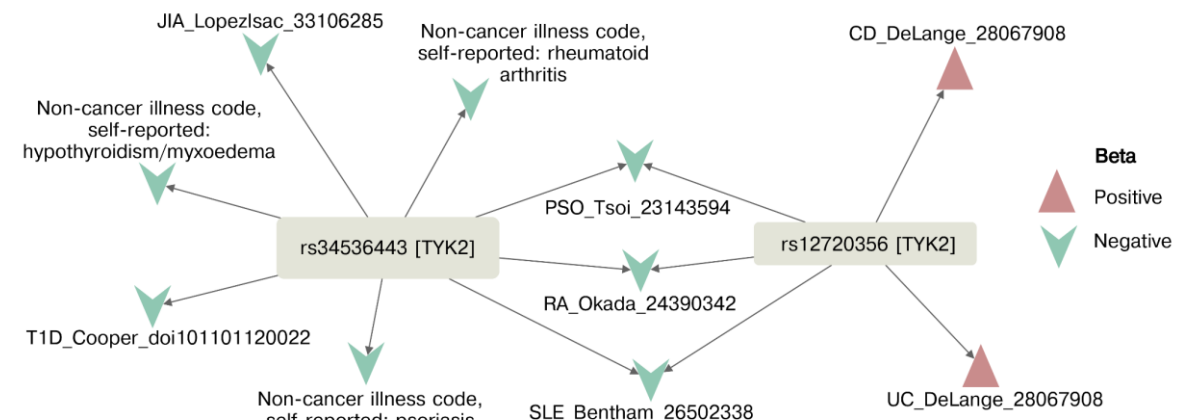252  migraine-related phenotypes in UKBB but not with any of the cardiovascular traits (Figure 5).
253



**Figure 5: Causal associations of Migraine related variants**

Hc associations (ppHc > 0.6) of the *PHACTR1* variant rs9349379_A>G and rs11172113_T>C, a *LRP1* variant. The direction of effect (beta) is shown with respect to the G and C allele respectively.

254

255  Other examples of pleiotropic variants include rs2476601, a non-synonymous variant in
256  *PTPN22* which we found to be causally associated with multiple autoimmune diseases and
257  their treatments as well as skin cancer, with the autoimmune-protective allele increasing risk
258  of cancer (Supplementary Figure 9). We also found a complex set of associations with two
259  variants in *APOE, rs*429358 and rs7412 that jointly define the three major structural isoforms
260  of APOE[19], $\varepsilon4$, $\varepsilon3$ and $\varepsilon2$ (Supplementary Figure 10). $\varepsilon2$ represents the TT haplotype
261  corresponding to the rs429358 and rs7412 variants, $\varepsilon3$ is represented by TC and $\varepsilon4$ by the
262  CC haplotype[20]. We found associations with increased risk of Alzheimer's disease, statin
263  medication, angina and ischemic heart disease with the $\varepsilon4$ allele with reference to the $\varepsilon2/\varepsilon3$
264  genotype. We also found a protective effect of $\varepsilon4$ compared to $\varepsilon2/\varepsilon3$ on traits related to a
265  family history of diabetes and blood pressure which correspond to similar traits found in
266  FinnGen as well as a protective effect of $\varepsilon3/\varepsilon4$ compared to $\varepsilon2$ for deep venous thrombosis
267  might be related to the ε3/ε4 genotype with reference to ε2 and might indicate the ε2 allele.
268  These findings align with previous studies on disease associations with different APOE
269  genotypes[21] and highlight the ability of SuSiE to map traits to distinct alleles in LD.

**Individual variant analyses**

270

271   CoPheScan can also be used to study single variants if sensible prior values can be supplied.
272   We considered exemplar non-synonymous variants in two genes, *TYK2* with established
273   allelic heterogeneity and associations to multiple immune-mediated diseases, and *SLC39A8,*
274   with established pleiotropic function. We ran CoPheScan with SuSiE BF and priors inferred
275   from the disease-causal variant analysis above ($p_a \approx$ 3.82e-5 and $p_c \approx$ 1.82e-3),
276   considering as query traits 2275 UKBB and 56 additional traits potentially related to either
277   gene from the GWAS catalog (Supplementary Table 3).

278

279   *TYK2* which encodes the tyrosine kinase 2 enzyme has multiple missense variants that have
280   been associated with a range of immune-mediated diseases (Supplementary Table 11). We
281   considered four: rs35018800_G>A (MAF: 0.0082), rs34536443_G>C (MAF: 0.0465),
282   rs12720356_A>C (MAF: 0.0979), and rs55882956_G>A (MAF: 0.0017). rs35018800_G>A
283   and rs55882956_G>A with the lowest MAF showed no association with any trait.
284   rs34536443_G>C was associated with 3 UKBB and 5 GWAS catalog traits, all immune-related
285   and previously established associations, including psoriasis, RA, JIA (Juvenile Idiopathic
286   Arthritis), Type 1 DM, and hypothyroidism. The variant rs12720356_A>C was associated with
287   ulcerative colitis, psoriasis, Crohn's disease, SLE (Systemic Lupus Erythematosus) and RA
288   traits from the GWAS catalog, but not with any of the UKBB traits (Figure 6).

289



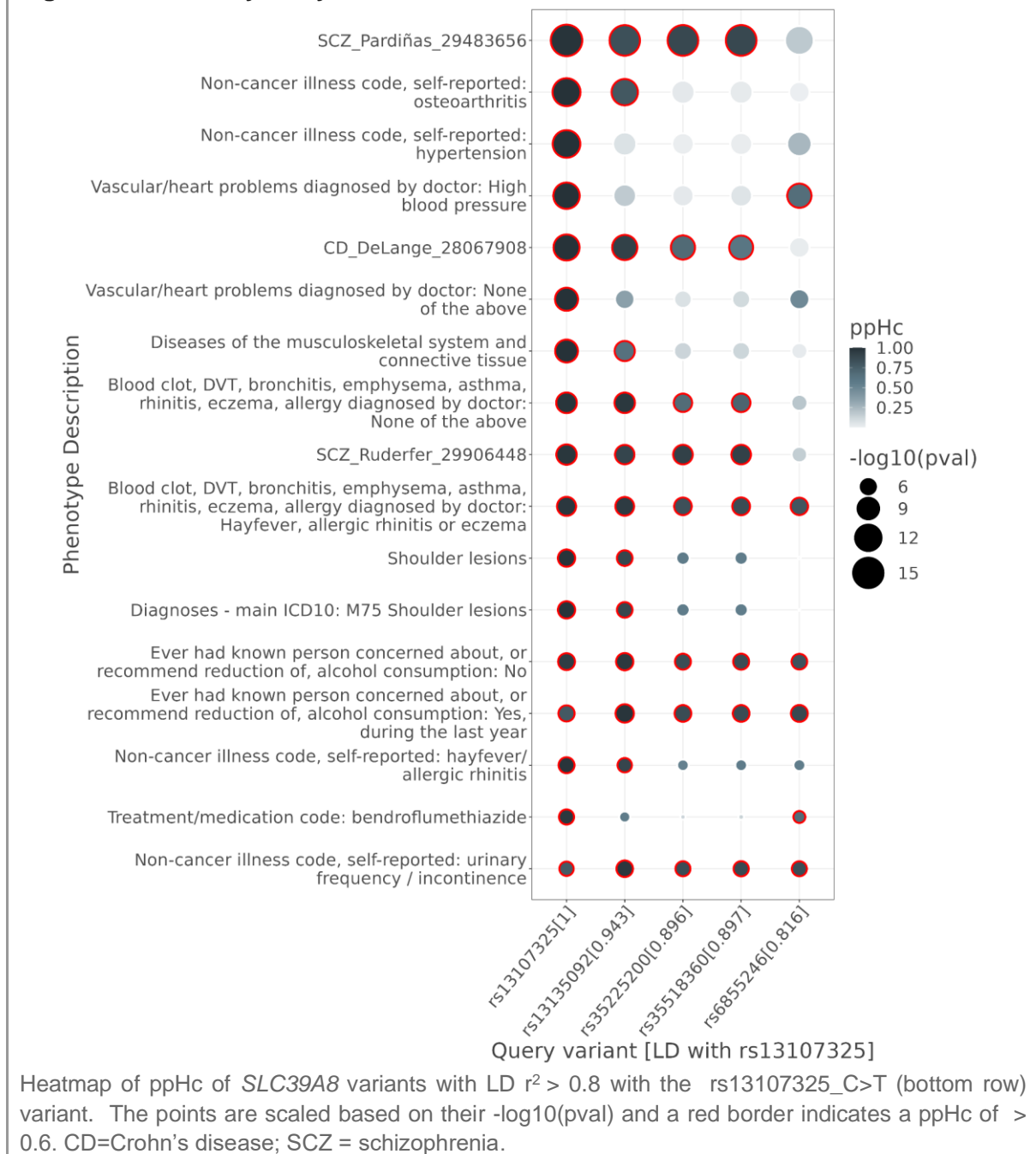**Figure 6: CoPheScan analysis of a gene with allelic heterogeneity: *TYK2***

Plots showing Hc associations (ppHc > 0.6) of *TYK2* variants rs34536443_G>C and rs12720356_A>C. The direction of beta is shown with respect to the ALT allele (C in both cases). T1D = Type 1 Diabetes Mellitus, JIA = Juvenile Idiopathic Arthritis, PSO = Psoriasis, RA = Rheumatoid Arthritis, SLE = Systemic Lupus Erythematosus, CD = Crohn's Disease, UC = Ulcerative Colitis

290

291   The highly pleiotropic variant, rs13107325_C>T, of *SLC39A8* (solute-carrier family gene which
292   encodes the ZIP8 protein), was associated with 14 UKBB and 3 GWAS catalog phenotypes,
293   replicating several known associations[22] with hypertension, schizophrenia, Crohn's disease,
294   urinary incontinence, musculoskeletal system-related traits such as osteoarthritis and traits
295   related to alcohol dependence.

296

11

297  We used this region to perform a sensitivity analysis, selecting four variants - rs6855246,
298  rs35225200, rs35518360, rs13135092, in LD with rs13107325_C>T (r2=0.816 - 0.943) and
299  running CoPheScan as if each had been selected as the causal variant.  This allows us to
300  explore two related questions: either, to what extent can two causal variants in LD cause false
301  positive findings, or, to what extent CoPheScan might still detect an association if the "causal"
302  variants supplied to CoPheScan are not really causal, but in LD with the causal variant. We
303  found that CoPheScan was indeed sensitive to this misspecification, where out of the 17 traits
304  identified as causally associated with rs13107325, 4 had ppHc < 0.6 with rs13135092
305  (r2=0.943) and 11 with rs6855246 (r2=0.816). The results were increasingly discrepant as the
306  r2 with rs13107325_C>T decreased (Figure 7 and Supplementary Figure 11). The group of
307  traits with high ppHc across multiple variants tended to have larger minimum p values in the
308  region compared to those for which ppHc was low across multiple variants, suggesting that
309  CoPheScan will be best at discriminating between potential causal variants in LD when the
310  association signal in the query data is strong.

311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340

**Figure 7: Sensitivity analysis with SLC39A8 variants**



Heatmap of ppHc of *SLC39A8* variants with LD $r^2 > 0.8$ with the rs13107325_C>T (bottom row) variant. The points are scaled based on their -log10(pval) and a red border indicates a ppHc of > 0.6. CD=Crohn's disease; SCZ = schizophrenia.

341

342  Finally, we sought to verify previously proposed causal associations between the *HMGCR*
343  variant rs12916_T>C and metabolic traits. *HMGCR* encodes HMG-CoA reductase which is
344  targeted by statins to lower LDL cholesterol. Previously, *HMGCR* variants have been used as
345  a proxy for statin effect to show a higher risk of type 2 diabetes and body mass index (BMI) in
346  MR studies[23]. But the validity of this has been challenged with evidence that there may be
347  distinct causal variants underlying type 2 diabetes, BMI and HMGCR levels[24]. We performed
348  CoPheScan analysis on the UKBB traits: LDL, BMI, type 2 diabetes, waist circumference and
349  weight. We identified a known causal association with LDL (ppHc = 1). Despite significant
350  observed p-values at rs12916 at BMI, weight and waist circumference, however, CoPheScan
351  consistently concluded that while the region contained a causal variant for each trait, it was

13

352   not rs12916 (ppHa > 0.99). In fact, no credible sets were identified in the *HMGCR* gene region
353   and the SuSIE signals from these traits indicate the presence of an alternative causal *POC5*
354   variant (Supplementary Figure 12). This implies that genetic studies that demonstrated a
355   relationship between statin therapy and BMI/T2DM through *HMGCR* variants as a proxy might
356   be incorrect[24] as they studied the SNPs in isolation while ignoring their regional context[25].
357   CoPheScan is thus valuable in verifying assumptions in instrumental variable analyses.

### PheWAS of protein-associated variants

359   One challenge of GWAS has been to link disease associations to their causal genes. PheWAS
360   allows us to start with variants with known causal function on a protein and ask which diseases
361   are also causally associated, exploiting the low false positive rate of CoPheScan.  We began
362   with 505 plasma protein QTLs[16] identified as single variant credible sets in fine-mapping of
363   527 plasma proteins. Nine variants were identified to be associated with UKBB traits (Table 1
364   and Supplementary Table 9). Among the established associations, we found an association
365   between a pQTL for APOC1 and high cholesterol, as well as reported treatment with the
366   cholesterol-lowering simvastatin. Both associations make sense given the known biology of
367   APOC1, but only the first would have been detected in scanning for significant p values, as
368   the p-value for high cholesterol at this SNP (p=6.19 x $10^{-19}$) is much lower than for simvastatin
369   (p=9.59 x $10^{-4}$), emphasising the value of exploiting the additional information that we believe
370   the variant to have a causal effect on a measurable phenotype (Supplementary Figure 13).

**Table 1: Hc associations detected with pQTL variants**

| Query variant | Protein | Direction | UKBB traits detected as Hc |
|---|---|---|---|
| rs11591147 | PCSK9 | risk | High cholesterol, cholesterol-lowering medication, ischaemic heart disease |
| rs5743618 | TLR1 | protect | Asthma |
| rs3775291 | TLR3 | risk | Hypothyroidism/myxoedema |
| rs3136516 | F2_Prothrombin | risk | Venous thromboembolism |
| rs34324219 | TCN1 | protect | Pernicious anaemia |
| rs964184 | APOC3 | risk | Cholesterol-lowering medication |
| rs116843064 | ANGPTL4 | risk | High cholesterol |
| rs5112 | APOC1 | risk | High cholesterol, cholesterol-lowering medication |
| rs214830 | TGM3 | risk | Skin cancer |

374   We also found a novel association, rs214830_G>C, a pQTL for TGM3, was associated with
375   skin cancer (ppHc=0.75). TGM3 is required for skin development and is normally expressed
376   in the spinous/granular layers of the epidermis. Its expression was found to be absent in
377   melanoma and squamous cell carcinoma of the skin but strongly expressed in basal cell
378   carcinoma (BCC), suggesting it could be a specific marker for BCC diagnosis[26] . Association
379   of variants in *TGM3* with BCC have also been reported[27–29] but rs214830_G>C was not
380   always the top variant and GWAS associations can mark causal effects in neighbouring

381  genes[30]. Our analysis suggests this association could be directly causal, with TGM3 involved
382  in the development of BCC as well as acting as a biomarker.
383
384  Finally, we considered 3586 variants labelled as protein-truncating (PTV) in the UKBB
385  summary data with MAF > 0.001, consisting of those predicted by VEP to be stop_gained,
386  frameshift, splice_acceptor and splice_donor. The fraction of query variants that were found
387  to be causally associated with at least one trait in UKBB was much lower for PTV (~0.31%)
388  than for disease-causal variants identified in FinnGen (~40%) and pQTL (~1.8%) (Table 2,
389  Supplementary Figures 5 and 6).
390
391  Examination of the Markov chain Monte Carlo (MCMC) chains showed issues with mixing for
392  the PTV example which were not seen with the other datasets (Supplementary Figure 5).
393  When we examined the inferred priors (Supplementary Table 12) obtained from this model,
394  we observed that the $p_c/p_a$ ratio was ~1.02, indicating that the inferred $p_a$ and $p_c$ priors were
395  almost the same. Our PTV consisted of four VEP classes, but while the MAF distribution of
396  the stop-gained PTV was similar to missense variants, those of the other PTV (frameshift,
397  splice donor and splice acceptor) were similar to synonymous variants (Supplementary Figure
398  14a). As selection can constrain MAF, we hypothesised that the VEP stop_gained class might
399  be more enriched for functional variation than the set of four classes we had used. We
400  considered two ways to enrich the PTV set for functional variation: either using just this subset
401  of the stop-gained PTVs or using the PTVs which were also defined as high confidence
402  homozygous predicted loss-of-function (pLoF) variants in gnomAD.[31] pLoF were
403  predominantly rare, such that the pLoF subset of PTV variants had a higher number of rare
404  variants compared to the stop-gained subset (Supplementary Figure 14b).
405
406  We ran the hierarchical models for these two subsets of PTVs (Supplementary Figure 15).
407  Comparing the priors (Supplementary Table 12) across the different datasets tested we
408  observed that the ratio of prior probabilities for the query variant or a non-query variant to be
409  causal, $p_c/p_a$ (Table 2) obtained using the pLoF variants (2.59) was second only to the ones
410  obtained using the FinnGen disease-related variants. The ratio from the stop-gained variant
411  model (1.39) was similar to the pQTL variant model (1.28). This shows that sets of query
412  variants which have a higher functional enrichment are expected to have a high $p_c/p_a$ ratio.
413
414
415
416
417
418
419
420
421
422
423
424

**Table 2: Summary of tested variants and phenotypes from real data**

| Query variant set | N QV | N QT | N QV-QT pairs* | N QV-QT detected as $H_c$ | N (%) unique variants detected as $H_c$ | pc/pa |
|---|---|---|---|---|---|---|
| FinnGen (with $r_g$) | 75 | 1582 | 130697 | 184 | 30 (40%) | 95.6 |
| FinnGen (no $r_g$) | 136 | 2275 | 193706 | 328 | 43 (31.62%) | 47.5 |
| pQTL (no $r_g$) | 505 | 2275 | 954616 | 29 | 9 (1.78%) | 1.28 |
| PTV all (no $r_g$) | 3586 | 2275 | 4359271 | 26 | 11 (0.31%) | 1.02 |
| PTV gnomAD (no $r_g$) | 366 | 2275 | 292787 | 7 | 2 (0.54%) | 2.59 |
| PTV stop gained (no $r_g$) | 911 | 2275 | 837060 | 15 | 6 (0.66%) | 1.39 |

425   QV - query variant, QT - query trait, N QV-QT number of trait-variant associations. The number of QT
426   were lower for the FinnGen/UKBB dataset for the 'with $r_g$' case as only traits having $r_g$ data with the
427   primary traits available were retained (Methods).
428
429   26 associations were identified using all the PTV variants. All 15 associations detected with
430   the stop-gained PTVs and 7 from the pLOF overlapped with those from the whole set. Of the
431   combined 26 PTV-trait associations (Supplementary Table 10), many corresponded to known
432   effects. One of them is, rs2066847_G>GC, a *NOD2* frameshift mutation, which is reported as
433   a pathogenic variant for inflammatory bowel disease in ClinVar and was associated with
434   several phenotypes related to Crohn's disease and mouth ulcers in our analysis
435   (Supplementary Figure 16). However, as seen with migraine and cardiovascular disease
436   above, the association with mouth ulcers occurs in the opposite direction to the established
437   comorbidity of Crohn's disease and mouth ulcers in the population, with the Crohn's disease
438   risk allele appearing protective for mouth ulcers. Note that in the mouth ulcer trait, the effect
439   sizes were opposite in two other SNPs identified as a credible set in SuSiE analyses of both
440   traits (Supplementary Figure 16).

## Discussion

442   Detection of pleiotropic effects of genetic variants is an essential component of target
443   discovery and drug repositioning. PheWAS typically takes information from marginal statistics
444   at query variants in isolation of their neighbours, which can lead to false positives when
445   multiple causal variants exist in some LD. CoPheScan considers not only how small a p-value
446   is at a given variant, but how small it is in comparison to its neighbours, and estimates how
447   much upweighting should be applied due to the information that the variant is in a query variant
448   set. In our simulations, CoPheScan showed considerably better control of false positive calls
449   compared to a standard PheWAS approach, at the cost of lower sensitivity where multiple
450   causal variants exist in a region. Whilst the higher false positive rate for standard PheWAS
451   testing can be mitigated by the use of a second-stage analysis testing for colocalisation, that
452   is not possible in the case of query SNPs selected for their known effects on a protein, such
453   as the PTV considered here.
454

455  CoPheScan learns how much to upweight query variants through the prior parameter $p_c$ and
456  the ratio of average $p_c$ to average $p_a$ is a useful measure of enrichment of causal variants for
457  the set of query traits amongst the set of query variants. This measure can be used to assess
458  the quality of any choice of variant set, with values close to 1 indicating a weak choice. It may
459  vary considerably across query variant sets for the same set of query traits, as seen in the
460  PTV analyses. However, while restriction to a smaller set of query variants with greater
461  enrichment is likely to find a higher proportion of causal associations with the smaller set, this
462  will not necessarily enhance discovery: whilst the majority of the discoveries found using the
463  smaller, more enriched sets of PTV were also found in the larger unfiltered set, this restriction
464  also meant losing plausible discoveries that didn't fall into either of the more restricted classes.
465
466  We allow $p_c$ to vary between variants by exploiting additional external information in a
467  regression framework. In our disease-variant focused analysis, we used the genetic
468  correlation between index and query traits, but this could also be a categorical variable, such
469  as the predicted deleteriousness of a missense variant, or the level of evidence for the
470  functional effect of a PTV. Our model can exploit covariate information that relates to query
471  trait-query variant pairs, but would need to be extended to accommodate other information.
472  For example, we might see modest evidence for causal association of a medication trait with
473  a given query variant, but intuitively trust the result is true because of stronger evidence at the
474  same variant with the disease itself. The difference in inference in such a case might be
475  explained by the smaller numbers of individuals reporting use of a specific medication. We
476  could consider exploiting genetic correlation between query traits by using a multivariate prior,
477  with covariance linked to the genetic correlations. While this is beyond the scope of the current
478  study, we hope our use of covariate-informed priors illustrates the potential for external
479  information to be exploited when conducting PheWAS and other genetic studies.
480
481  While the simulations emphasised the importance of learning $p_c$ in a hierarchical model for
482  accurate inference, point estimates can be substituted if required. This borrowing of priors
483  from a larger dataset is beneficial in scenarios where we might want to use CoPheScan to test
484  associations between a small set of variants and phenotypes, as running a hierarchical model
485  on limited data will not result in optimal prior estimates. However, we strongly advise that
486  careful consideration is needed to ensure the larger dataset in which the priors are learnt is a
487  good match for the limited dataset under consideration.
488
489  One of the advantages of incorporating SuSIE in CoPheScan is the ability to detect allelic
490  heterogeneity at a locus. We demonstrated this with two well-known distinct variants in the
491  *TYK2* gene which were associated with overlapping sets of immune-mediated disorders. This
492  analysis also highlighted the importance of surveying disease-specific GWAS studies and not
493  relying solely on biobanks which may hold relatively low numbers of cases of any individual
494  disease. For example, only three UKBB traits showed any association compared to seven of
495  our curated immune-mediated disease GWAS, and while psoriasis in UKBB (4192 cases) was
496  identified with one variant, psoriasis in Tsoi's GWAS study (10558 cases) was identified with
497  two. While biobanks remain incredibly useful for common traits such as cardiovascular and
498  metabolic diseases, carefully curated bespoke GWAS of less common traits should be

499  included in any PheWAS to complement the biobank resources and reveal the full spectrum
500  of pleiotropy. This is particularly important because predicted beneficial effects of targeting a
501  protein may be countered by on-target side effects on other traits, as we saw where the
502  autoimmune-protective variant in *PTPN22* was associated with an increased risk of skin
503  cancer.
504
505  Our CoPheScan approach has some specific limitations. The signals obtained from the
506  multiple variant assumption rely on the available LD information. Zou et al. demonstrated that
507  the performance of SuSIE degrades when presented with out-of-sample LD matrices[14]. In
508  cases where in-sample LD matrices are not accessible, it is recommended to utilise out-of-
509  sample LD from large reference panels. Our experience with SuSIE is that as the number of
510  causal variants increase, the performance (ability to detect all causal variants and/or ability to
511  detect the correct causal variants) may decrease. In the case that a true causal variant signal
512  is missed for our query variant, as occurred in around 35% of our simulations with two causal
513  variants, CoPheScan concludes Ha - ie a false negative. SuSiE is likely to miss a higher
514  fraction of secondary causal variants when the true number of causal variants increase, which
515  would be expected to lead to greater false negatives for CoPheScan. Importantly, we do not
516  see any increase in false positives when simulating two compared to one causal variant
517  Analysis of rare events in large samples with standard methods can cause bias in regression
518  summary statistics. Here, we used careful QC, thresholding on the number of events / MAF,
519  but a better approach would be to use methods specifically developed to deal with this such
520  as REGENIE[32] to generate input to CoPheScan. The current form of CoPheScan only allows
521  single-ancestry studies which will be addressed in future iterations and allow an increase in
522  power to detect rare variants.
523
524  GWAS causal variants, even when identified with confidence, remain challenging to interpret
525  partly because it can be hard to link them with confidence to their causal genes. Protein-
526  altering variants have thus become increasingly important because their function on a gene is
527  presumed known. The different relative enrichments in different sets of PTV we ran suggests
528  that incorporating external evidence on the plausibility of a putative PTV having a functional
529  effect will increase accuracy in PheWAS of these variants. However, as highlighted here, they
530  often have very low minor allele frequencies. Thus, larger biobanks are still needed both for
531  analysis of less common traits with common variants and for analysis of rare functional
532  variants. It is thus encouraging that UKBiobank and FinnGen studied here are complemented
533  by the Japan Biobank[33], the Million Veteran Program[34] and the Uganda Genome Resource[35],
534  which should allow CoPheScan, together with efforts at multiple ancestry fine mapping[36], to
535  reveal more completely the pleiotropic spectrum of protein-altering genetic variation.


## Methods

### Simulated data

538  We simulated case-control summary statistics using the EUR samples in the 1000 Genomes
539  phase 3 reference data[37]. LD-independent blocks were identified using lddetect[38] and
540  haplotypes containing 1000 SNPs with MAF > 0.01 were extracted from the reference data[10,39].

18

541 We used simGWAS[39] to simulate summary statistics with either one or two causal variants for
542 the corresponding LD blocks with 10000 cases and 10000 controls.

543 We simulated GWAS summary statistics for 110000 traits to evaluate hypothesis
544 discrimination, with all genomic regions containing 1000 SNPs. We sampled query causal
545 variants at random from the 1000 SNPs and simulated each trait to correspond to one of the
546 three hypotheses. The traits within the simulated dataset were divided based on the number
547 and position of the causal variants within their genomic region:

    1. **True Hn**: No causal variants within the genomic region.
549     2. **True Ha**: A trait with a single causal variant that is not the query variant.
550     3. **True Ha2**: A genomic region simulated with two distinct causal variants, none of which
551        are the query variant.
552     4. **True Hc**: A trait with a single causal variant that is the same as the query variant.
553     5. **True Hc2**: Two distinct causal variants, where one of them is the same as the query
554        variant.

555 The 110000 simulated traits were comprised of 88048 true Hn, 6276 each of true Hc and Hc2,
556 and 4700 each of Ha and Ha2 traits. We also simulated genetic correlation values for each of
557 these traits where the Hc traits were assigned a higher proportion of high $r_g$ values when
558 compared to the Hn and Ha traits (Supplementary Figure 2).

559 We used conventional PheWAS approaches, based on selecting associations that cross a
560 threshold p-value after accounting for multiple testing. We used the Benjamini & Hochberg to
561 control the FDR < 0.05 which corresponded to a p-value < 7.5e-3 and Bonferroni correction
562 with a p-value < 4.55e-7 to control the family-wise error rate (FWER) at 0.05. In parallel, we
563 used different approaches of CoPheScan to analyse this dataset:

    1. Fixed priors and Approximate Bayes Factors (ABF).
565        Fixed priors used with the simulated data were adapted from coloc (described in the
566        Supplementary Methods) where $p_a \approx$ 0.81 and $p_c \approx$ 0.091.
567     2. Fixed priors and SuSIE Bayes factors
568     3. Hierarchical priors, ABF, with and without genetic correlation ($r_g$)
569     4. Hierarchical priors, SuSIE Bayes factors, with and without $r_g$

570

571 The hierarchical model of CoPheScan was run for 3e5 iterations for the models without $r_g$ and
572 1e6 iterations for the ones with $r_g$ and the chain was thinned by retaining every 30th and 100th
573 observation respectively (Supplementary Figure 3). The first 50% of the remaining 1e4
574 observations were discarded and the average prior ($p_n$, $p_a$, $p_c$) and posterior probabilities
575 ($ppH_n$, $ppH_a$, $ppH_c$) were calculated.

576

577 Therefore, the output of each analysis for each trait was summarised by the Hc hypothesis
578 when $ppH_c$ > 0.6, and Hn when $ppH_n$ > 0.2 and Ha for the remaining traits. When
579 multiple signals were detected by SuSIE in the same genomic region, CoPheScan was run on
580 each of them. Here, we assigned the hypothesis to each signal as the thresholds specified
581 above. The first hypothesis to occur in the ranking order of Hc, Hn, and Ha was assigned to
582 the trait, i.e., when there was at least one signal which was assigned as Hc, the trait was taken
583 to be Hc. Next, any signal with Hn but no Hc was assigned as Hn, because we wanted to be
584 conservative in calling Ha which might rule out a pleiotropic effect. In the absence of a Hc and

585  Hn signal, we checked for the presence of Ha and where there were multiple Ha signals we
586  report the minimum Ha of all the signals.
587
588  We also simulated datasets where we varied the number of true Hc traits {50, 100, 200, 500,
589  1000, 2000, 4000, 5000} while maintaining the true Hn and Ha traits the same at 88048 and
590  4700 respectively. The percentage of Hc traits in the datasets corresponded to {0.05%, 0.11%,
591  0.22%, 0.54%, 1.07%, 2.11%, 3.13%, 4.13%, 5.12%}.

## Disease causal query variants from FinnGen

593  We downloaded SuSIE[13,14] fine-mapping results from FinnGen[40] release R5, which has a
594  sample size of 218,792 with 2,803 endpoints (https://www.finngen.fi/en/access_results). For
595  each endpoint, we filtered variants that belonged to a credible set of size 1 and were also
596  present in the UKBB dataset. We retained 136 variants, fine-mapped from 141 FinnGen traits,
597  for further analysis. 69 out of these 141 FinnGen primary traits, had closely matching UKBB
598  traits with pre-computed genetic correlation data. Thus, 75 variants from these matching traits
599  were used for the hierarchical model using $r_g$ (Supplementary Tables 4 and 5).

## pQTL variants

601  We downloaded summary data from a GWAS of plasma protein levels measured with 4,907
602  aptamers (corresponding to 4719 proteins) in 35,559 Icelanders from Ferkingstad et al.[16]. We
603  fine-mapped the region around each signal under a single variant assumption. This is
604  equivalent to taking only the first signal in a stepwise fine-mapping procedure. We made this
605  conservative choice to address the lack of access to an LD matrix for the Icelandic population,
606  making it difficult to trust secondary signals found by stepwise regression or other multiple
607  causal variant methods such as SuSIE. We obtained 505 SNPs associated with 527 proteins
608  for testing associations with the UKBB phenotypes  (Supplementary Table 6).

## Protein truncating variants

610  We selected 3586 protein truncating variants (PTV) with MAF > 0.001 (Supplementary Table
611  7) from the UKBB variants (https://broad-ukb-sumstats-us-east-
612  1.s3.amazonaws.com/round2/annotations/variants.tsv.bgz), which were annotated as
613  frameshift (883), stop gained (911), splice acceptor (682) and splice donor (1110) variants,
614  using VEP[41] (The Ensembl Variant Effect Predictor, 85).
615  We also downloaded homozygous pLoF from gnomAD (v2.1.1). Out of these, we selected 366
616  variants, which were classified as either 'lof' or 'likely_lof' and were in common with the 3586
617  UKBB PTV variants[31] .

## Individual variant analyses

619  We chose four *TYK2* variants: rs35018800, rs34536443, rs12720356 and rs55882956 a
620  *SLC39A8* variant, rs13107325, and a *HMGCR* variant,  rs12916, to examine the performance
621  of CoPheScan for region-specific analyses[42–44].

622 **Query phenotypes**

623 We used 2275 phenotypes from the UK Biobank (http://www.nealelab.is/uk-biobank). We
624 obtained in-sample linkage disequilibrium matrices from https://broad-alkesgroup-ukbb-
625 ld.s3.amazonaws.com/UKBB_LD[15]. We included all the 2275 traits in the CoPheScan analysis
626 of the FinnGen, pQTL and PTV variants.

627 We downloaded genetic correlation[45] data between UK Biobank traits and disorders estimated
628 using LD score regression[46] from https://ukbb-rg.hail.is/. In the FinnGen/UKBB dataset, only
629 1582 out of the 2275 traits had genetic correlation ($r_g$) estimates with the UKBB traits mapped
630 to the FinnGen primary traits. So, only these traits were used for the hierarchical model that
631 included $r_g$. Additionally, we checked the allele counts (AC) of the query variant in the
632 phenotype files and only retained the QV-QT pairs for association testing when the AC in the
633 cases > 25. After reviewing the results, to increase stringency, we further removed QV-QT
634 pairs, identified as being causally associated with AC < 30 to reduce false positive detection.
635 Individual results presented in tables have been trimmed to reflect this more stringent criterion
636 (removing five Hc results), but estimates directly from models (eg pa/pc) include observations
637 with AC between 26 and 30.

638

639 For the phenome-wide scan of the *TYK2* and *SLC39A8* variants, we downloaded 56 additional
640 publicly available GWAS summary statistics of phenotypes related to immune-mediated and
641 psychiatric diseases (Supplementary Table 3)[47–49]. In the case of the *HMGCR* variant we used
642 additional quantitative UKBB phenotypes: LDL direct, Body mass index (BMI) and Weight,
643 We used UKBB LD matrices for data from European populations, and for other populations,
644 we extracted LD from the 1000 genomes phase 3 reference data[37].
645 The lists of phenotypes used with the different variants are provided in Supplementary Tables
646 2 and 3. We used Phase II HapMap[50] obtained from
647 (https://ftp.ncbi.nlm.nih.gov/hapmap/recombination/2011-01_phaseII_B37/) to subset regions
648 from the summary statistics data around the query variants. We excluded variants in the HLA
649 region (20MB - 40MB) from the analysis.

650

651 Visualisations of the causal trait-variant associations identified with CoPheScan were done
652 using Cytoscape[51] 3.9.0.

653 **Functional annotation**

654 Previously reported variant/gene associations with diseases were obtained from the Open
655 Target Platform and Open Target Genetics[48,52]. We used the DrugBank online resource for
656 indications of medications that were associated with the variants[53].

657 # Data availability

658 The simulated summary statistics and processed files are available on figshare. Source data
659 are provided with this paper.

**Code availability**

661 The CoPheScan R package is available on CRAN at https://cran.r-
662 project.org/package=cophescan.
663 A shiny app to browse the results is available here:https://ichcha-m.shinyapps.io/cophescan-
664 app/, the code for which can be found at https://github.com/ichcha-m/cophescan-app.
665
666 The code to reproduce the simulated summary statistics and processed datasets are
667 available here: https://github.com/chr1swallace/cophescan-manuscript-sim-summary-data
668 and https://github.com/ichcha-m/cophescan-paper.


**URLs**

670 UKBB summary statistics, http://www.nealelab.is/uk-biobank); Phase II HapMap,
671 https://ftp.ncbi.nlm.nih.gov/hapmap/recombination/2011-01_phaseII_B37/; UKBB in-sample
672 LD matrices, https://broad-alkesgroup-ukbb-ld.s3.amazonaws.com/UKBB_LD ; FinnGen
673 Freeze 5 cohort, https://www.finngen.fi/en/access_results; GWAS catalog,
674 https://www.ebi.ac.uk/gwas/; GWAS of plasma protein levels used for pQTL fine-mapping:
675 https://www.decode.com/summarydata/, ClinVar variant annotation:
676 https://platform.opentargets.org/downloads; gnomAD:
677 https://gnomad.broadinstitute.org/downloads; DrugBank, https://go.drugbank.com/drugs/


**References**

679 1. Denny, J. C. *et al.* PheWAS: demonstrating the feasibility of a phenome-wide scan to

680    discover gene–disease associations. *Bioinformatics* **26**, 1205–1210 (2010).

681 2. Rastegar-Mojarad, M., Ye, Z., Kolesar, J. M., Hebbring, S. J. & Lin, S. M. Opportunities

682    for drug repositioning from phenome-wide association studies. *Nat Biotechnol* **33**, 342–

683    345 (2015).

684 3. Diogo, D. *et al.* Phenome-wide association studies across large population cohorts

685    support drug target validation. *Nat. Commun.* **9**, 4285 (2018).

686 4. Millard, L. A. C. *et al.* MR-PheWAS: hypothesis prioritization among potential causal

687    effects of body mass index on many outcomes, using Mendelian randomization. *Sci Rep*

688    **5**, 16645 (2015).

689 5. Zuber, V. *et al.* Combining evidence from Mendelian randomization and colocalization:

690    Review and comparison of approaches. *Am J Hum Genet* **109**, 767–782 (2022).

691 6. Verma, A. *et al.* Human-Disease Phenotype Map Derived from PheWAS across 38,682

692    Individuals. *Am J Hum Genet* **104**, 55–64 (2019).

693    7. DeBoever, C. *et al.* Medical relevance of protein-truncating variants across 337,205

694        individuals in the UK Biobank study. *Nat. Commun.* **9**, 1612 (2018).

695    8. Veturi, Y. *et al.* A unified framework identifies new links between plasma lipids and

696        diseases from electronic medical records across large-scale cohorts. *Nat Genet* **53**, 972–

697        981 (2021).

698    9. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic

699        association studies using summary statistics. *PLoS Genet* **10**, e1004383 (2014).

700    10.     Wallace, C. A more accurate method for colocalisation analysis allowing for multiple

701        causal variants. *PLoS Genet* **17**, e1009440 (2021).

702    11.     Wallace, C. Eliciting priors and relaxing the single causal variant assumption in

703        colocalisation analyses. *PLoS Genet* **16**, e1008720 (2020).

704    12.     Wakefield, J. Bayes factors for genome-wide association studies: comparison with P-

705        values. *Genet Epidemiol* **33**, 79–86 (2009).

706    13.     Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to

707        variable selection in regression, with application to genetic fine mapping. *J. R. Stat. Soc.*

708        *Ser. B Stat. Methodol.* **82**, 1273–1300 (2020).

709    14.     Zou, Y., Carbonetto, P., Wang, G. & Stephens, M. Fine-mapping from summary data

710        with the "Sum of Single Effects" model. *PLOS Genet.* **18**, e1010299 (2022).

711    15.     Weissbrod, O. *et al.* Functionally informed fine-mapping and polygenic localization of

712        complex trait heritability. *Nat Genet* **52**, 1355–1363 (2020).

713    16.     Ferkingstad, E. *et al.* Large-scale integration of the plasma proteome with genetics

714        and disease. *Nat Genet* **53**, 1712–1721 (2021).

715    17.     Daghlas, I., Guo, Y. & Chasman, D. I. Effect of genetic liability to migraine on

716        coronary artery disease and atrial fibrillation: a Mendelian randomization study. *Eur J*

717        *Neurol* **27**, 550–556 (2020).

718    18.     Kurth, T. *et al.* Migraine and risk of cardiovascular disease in women: prospective

719        cohort study. *BMJ* **353**, (2016).

720    19.     Huebbe, P. & Rimbach, G. Evolution of human apolipoprotein E (APOE) isoforms:

721        Gene structure, protein function and interaction with dietary factors. *Ageing Res. Rev.* **37**,

722    146–161 (2017).

723    20.    Babenko, V. N. *et al.* Haplotype analysis of APOE intragenic SNPs. *BMC Neurosci.*

724    **19**, 16 (2018).

725    21.    Lumsden, A. L., Mulugeta, A., Zhou, A. & Hyppönen, E. Apolipoprotein E (APOE)

726    genotype-associated disease risks: a phenome-wide, registry-based, case-control study

727    utilising the UK Biobank. *eBioMedicine* **59**, (2020).

728    22.    Nebert, D. W. & Liu, Z. SLC39A8 gene encoding a metal ion transporter: discovery

729    and bench to bedside. *Hum Genomics* **13**, 51 (2019).

730    23.    Swerdlow, D. I. *et al.* HMG-coenzyme A reductase inhibition, type 2 diabetes, and

731    bodyweight: evidence from genetic analysis and randomised trials. *The Lancet* **385**, 351–

732    361 (2015).

733    24.    Holm, H. *et al.* Abstract 18777: The Low-Density Lipoprotein Cholesterol and Body

734    Mass Index/Type-2 Diabetes Signals in the HMGCR Region Are Not Explained by a

735    Single Variant. *Circulation* **134**, A18777–A18777 (2016).

736    25.    VanderWeele, T. J., Tchetgen, E. J. T., Cornelis, M. & Kraft, P. Methodological

737    challenges in Mendelian randomization. *Epidemiol. Camb. Mass* **25**, 427 (2014).

738    26.    Smirnov, A. *et al.* Transglutaminase 3 is expressed in basal cell carcinoma of the

739    skin. *Eur. J. Dermatol.* **29**, 477–483 (2019).

740    27.    Stacey, S. N. *et al.* Germline sequence variants in TGM3 and RGS22 confer risk of

741    basal cell carcinoma. *Hum. Mol. Genet.* **23**, 3045–3053 (2014).

742    28.    Ue, L. *et al.* Combined analysis of keratinocyte cancers identifies novel genome-wide

743    loci. *Hum. Mol. Genet.* **28**, 3148–3160 (2019).

744    29.    Adolphe, C. *et al.* Genetic and functional interaction network analysis reveals global

745    enrichment of regulatory T cell genes influencing basal cell carcinoma susceptibility.

746    *Genome Med.* **13**, 19 (2021).

747    30.    Smemo, S. *et al.* Obesity-associated variants within FTO form long-range functional

748    connections with IRX3. *Nature* **507**, 371–375 (2014).

749    31.    Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation

750    in 141,456 humans. *Nature* **581**, 434–443 (2020).

24

751    32.    Mbatchou, J. *et al.* Computationally efficient whole-genome regression for

752    quantitative and binary traits. *Nat. Genet.* **53**, 1097–1103 (2021).

753    33.    Nagai, A. *et al.* Overview of the BioBank Japan Project: Study design and profile. *J.*

754    *Epidemiol.* **27**, S2–S8 (2017).

755    34.    Gaziano, J. M. *et al.* Million Veteran Program: A mega-biobank to study genetic

756    influences on health and disease. *J. Clin. Epidemiol.* **70**, 214–223 (2016).

757    35.    Fatumo, S. *et al.* Uganda Genome Resource: A rich research database for genomic

758    studies of communicable and non-communicable diseases in Africa. *Cell Genomics* **2**,

759    100209 (2022).

760    36.    Yuan, K. *et al.* Fine-mapping across diverse ancestries drives the discovery of

761    putative causal variants underlying human complex traits and diseases.

762    2023.01.07.23284293 Preprint at https://doi.org/10.1101/2023.01.07.23284293 (2023).

763    37.    1000 Genomes Project Consortium *et al.* A global reference for human genetic

764    variation. *Nature* **526**, 68–74 (2015).

765    38.    Berisa, T. & Pickrell, J. K. Approximately independent linkage disequilibrium blocks in

766    human populations. *Bioinformatics* **32**, 283–285 (2016).

767    39.    Fortune, M. D. & Wallace, C. simGWAS: a fast method for simulation of large scale

768    case–control GWAS summary statistics. *Bioinformatics* **35**, 1901–1906 (2018).

769    40.    Kurki, M. I. *et al.* FinnGen provides genetic insights from a well-phenotyped isolated

770    population. *Nature* **613**, 508–518 (2023).

771    41.    McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol* **17**, 122

772    (2016).

773    42.    Diogo, D. *et al.* TYK2 protein-coding variants protect against rheumatoid arthritis and

774    autoimmunity, with no evidence of major pleiotropic effects on non-autoimmune complex

775    traits. *PLoS One* **10**, e0122271 (2015).

776    43.    Dendrou, C. A. *et al.* Resolving TYK2 locus genotype-to-phenotype differences in

777    autoimmunity. *Sci Transl Med* **8**, 363ra149 (2016).

778    44.    Motegi, T. *et al.* Identification of rare coding variants in TYK2 protective for

779    rheumatoid arthritis in the Japanese population and their effects on cytokine signalling. *Ann.*

25

780  *Rheum. Dis.* **78**, 1062–1069 (2019).

781  45.  Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and

782  traits. *Nat Genet* **47**, 1236–1241 (2015).

783  46.  Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from

784  polygenicity in genome-wide association studies. *Nat Genet* **47**, 291–295 (2015).

785  47.  Morales, J. *et al.* A standardized framework for representation of ancestry data in

786  genomics studies, with application to the NHGRI-EBI GWAS Catalog. *Genome Biol.* **19**, 21

787  (2018).

788  48.  Ochoa, D. *et al.* Open Targets Platform: supporting systematic drug–target

789  identification and prioritisation. *Nucleic Acids Res* **49**, D1302–D1310 (2020).

790  49.  Mease, P. J. *et al.* Efficacy and safety of selective TYK2 inhibitor, deucravacitinib, in

791  a phase II trial in psoriatic arthritis. *Ann Rheum Dis* **81**, 815–822 (2022).

792  50.  Frazer, K. A. *et al.* A second generation human haplotype map of over 3.1 million

793  SNPs. *Nature* **449**, 851–861 (2007).

794  51.  Shannon, P. *et al.* Cytoscape: A Software Environment for Integrated Models of

795  Biomolecular Interaction Networks. *Genome Res.* **13**, 2498–2504 (2003).

796  52.  Ghoussaini, M. *et al.* Open Targets Genetics: systematic identification of trait-

797  associated genes using large-scale genetics and functional genomics. *Nucleic Acids Res* **49**,

798  D1311–D1320 (2021).

799  53.  Wishart, D. S. *et al.* DrugBank 5.0: a major update to the DrugBank database for

800  2018. *Nucleic Acids Res* **46**, D1074–D1082 (2018).

## Acknowledgements

## Ethics declarations

Competing interests

CW and IM receive funding from GSK and MSD. CW is a part-time employee of GSK, but this research was conducted within her academic role.

## Supplementary information

1. **Supplementary Information**
   Supplementary Methods, and Figures 1-16
2. **Supplementary Tables**
   Supplementary Tables 1-12