1    **The impact of non-additive genetic associations on age-related complex diseases.**

2

3    Marta Guindo-Martínez[1,*], Ramon Amela[1,*], Silvia Bonàs-Guarch[1,2,3], Montserrat Puiggròs[1],

4    Cecilia Salvoro[1], Irene Miguel-Escalada[1,2,3], Caitlin E Carey[4,5], Joanne B. Cole[6,7,8,9], Sina

5    Rüeger[10], Elizabeth Atkinson[4,5,11], Aaron Leong[8,12], Friman Sanchez[1], Cristian Ramon-

6    Cortes[1], Jorge Ejarque[1], Duncan S Palmer[4,5,13], Mitja Kurki[10], FinnGen Consortium[14], Krishna

7    Aragam[11,15,16], Jose C Florez[6,7,17], Rosa M. Badia[1], Josep M. Mercader[6,7,1,#], David

8    Torrents[1,18,#].

9

10    1 - Barcelona Supercomputing Center (BSC), Barcelona, Spain
11    2 - Regulatory Genomics and Diabetes, Centre for Genomic Regulation, The Barcelona Institute of
12    Science and Technology, Barcelona, Spain
13    3 - CIBER de Diabetes y Enfermedades Metabólicas Asociadas, Madrid, Spain
14    4 - Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA,
15    USA
16    5 - Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General
17    Hospital, Boston, MA, USA
18    6 - Programs in Metabolism and Medical and Population Genetics, Broad Institute of MIT and
19    Harvard, Cambridge, MA, USA
20    7 - Diabetes Unit and Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA,
21    USA
22    8 - Harvard Medical School, Boston, Massachusetts, USA
23    9 - Division of Endocrinology and Center for Basic and Translational Obesity, Research, Boston
24    Children's Hospital, Boston, MA, USA
25    10 - Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland
26    11 - Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge,
27    Massachusetts, USA
28    12 - Department of Medicine, Massachusetts General Hospital, Boston, Massachusetts, USA
29    13 - Current address: GENOMICS plc, Oxford, UK
30    14 - Members of the consortium are provided in Appendix S1
31    15 - Cardiology Division, Massachusetts General Hospital, Boston, Massachusetts, USA
32    16 - Cardiovascular Research Center, Massachusetts General Hospital, Boston, Massachusetts, USA
33    17 - Department of Medicine, Harvard Medical School, Boston, MA, USA
34    18 - Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain
35
36    * - Both authors contributed equally to this work.
37    # - These authors jointly directed this work.
38
39    Corresponding authors:
40
41    Josep M Mercader
42    Programs in Metabolism and Medical and Population Genetics
43    Broad Institute of Harvard and MIT
44    75 Ames St
45    02142, Cambridge, MA
46    United States of America
47    E-mail: mercader@broadinstitute.org
48
49    and
50    David Torrents
51    Life Science Department

52    Barcelona Supercomputing Center (BSC)

53    Institució Catalana de Recerca i Estudis Avançats (ICREA)

54    C/Jordi Girona 29, Edifici Nexus II

55    08003 Barcelona, Catalunya, Spain

56    Phone: 3493 413 40 74

57    E-mail: david.torrents@bsc.es

58

59    Keywords

60    Genotype Imputation, genome-wide association studies (GWAS), inheritance models,

61    recessive analysis, reference panels, phenome-wide association studies (PheWAS),

62    complex diseases, age-related diseases.

63

64    **Abstract**

65    Genome-wide association studies (GWAS) are not fully comprehensive as current strategies

66    typically test only the additive model, exclude the X chromosome, and use only one

67    reference panel for genotype imputation. We implemented an extensive GWAS strategy,

68    GUIDANCE, which improves genotype imputation by using multiple reference panels,

69    includes the analysis of the X chromosome and non-additive models to test for association.

70    We applied this methodology to 62,281 subjects across 22 age-related diseases and

71    identified 94 genome-wide associated loci, including 26 previously unreported. We observed

72    that 27.6% of the 94 loci would be missed if we only used standard imputation strategies and

73    only tested the additive model. Among the new findings, we identified three novel low-

74    frequency recessive variants with odds ratios larger than 4, which would need at least a

75    three-fold larger sample size to be detected under the additive model. This study highlights

76    the benefits of applying innovative strategies to better uncover the genetic architecture of

77    complex diseases.

78

**Introduction**

Genome-wide association studies (GWAS) have been successful in identifying thousands of associations between genetic variation and human complex diseases and traits [1]. Nevertheless, for most complex diseases, only a small fraction of their genetic architecture is known and a small amount of the estimated heritability is explained [2]. Variants that individually have small contributions to the risk of disease, and/or are rare in the population, are often missed by the majority of GWAS even though their role in the pathophysiology of complex diseases can be crucial. Some of the current limitations of GWAS could be overcome by increasing sample sizes and, as recently demonstrated, by applying more comprehensive analytical methods with improved imputation strategies [3]. Though the increase of sample size might allow the detection of more genetic signals, it also imposes major methodological and computational requirements. These can require scientists to restrict and simplify the analysis by limiting it to autosomal chromosomes, a single reference panel for imputation, and a single (additive) inheritance model for association testing, leaving a relevant fraction of the genetic architecture of the disease unexplored [4].

The genetic variants that modify the risk to develop a particular complex disease may contribute to the final phenotype through different functional mechanism defined by a particular model of inheritance, which is further reflected in a characteristic distribution of affected alleles across patients and healthy individuals in GWAS. For example, the additive inheritance model, which is often the only genetic model tested, assumes that the risk of the disease is proportional to the number of risk alleles in an individual, i. e., that the effect of the heterozygous genotype is halfway between the two possible homozygous genotypes. However, some variants follow non-additive inheritance models, which include dominant, recessive and heterodominant. The additive model is expected to capture a large fraction of the genetic risk for disease [5] and can identify some variants that follow non-additive inheritance patterns. However, the additive model is not sufficient to provide a comprehensive overview of the genetic architecture of diseases. In particular, most GWAS may have insufficient power to identify low-frequency variants that show recessive effects [6, 7]. The importance of evaluating non-additive inheritance models is well reported in the context of Mendelian diseases [8] and occasionally for complex traits as well, such as the recessive effects of the *FTO* locus in obesity [9], the *ITGA1* [10], *TBC1D4* [11] and *CDKAL1* [9, 12] genes in type 2 diabetes, as well as the known non-additive effects of HLA haplotypes in autoimmune diseases [13] and ulcerative colitis [14]. The increasing ability to capture low-frequency variants using modern imputation reference panels and the need to uncover the still missing heritability estimated for most complex diseases, call for comprehensive

114 association strategies that should include, among other improvements, the analysis of non-
115 additive inheritance models.

116 To fill this gap and to determine the prevalence and contribution of the different inheritance
117 patterns involved in the genetic architecture of complex diseases, we have designed and
118 implemented a comprehensive strategy for genetic association analysis that combines
119 dense imputation from multiple reference panels with association testing under five different
120 inheritance models across multiple phenotypes. We have applied this strategy to the Kaiser
121 Permanente Research Program on Genes, Environment and Health: A Genetic
122 Epidemiology Research on Adult Health and Aging (GERA) cohort [15], which includes 62,281
123 subjects from European ancestry and 22 diseases.

124 Finally, we release here both the summary statistics for all the models of inheritance as well
125 as the complete methodology, provided to the community as an easy-to-use and standalone
126 pipeline. This pipeline allows for analysis of existing and newly generated GWAS data with
127 better efficiency and more comprehensive testing, improving the chances of variant
128 discovery.

129

130 **Results**

131 In order to assess the potential benefits of applying more in-depth GWAS methodologies to
132 available genetic datasets, and to investigate the relative contribution of different inheritance
133 models to the risk to develop complex diseases, we have applied a global analysis strategy
134 to the GERA cohort, an age-related disease-based cohort with an average age of 63, well
135 powered to study a broad range of clinically defined age-related conditions. By using this
136 particular cohort, we expect to minimize a possible loss of power due to the misclassification
137 of controls, as often happens in datasets with younger individuals that can include cases at
138 pre-disease stages classified as controls.

139

140 **Genotype Imputation and association testing using multiple reference panels**

141 After applying strict genetic quality control to the GERA cohort (see Methods), we retained
142 56,637 individuals with European ancestry for further downstream analysis (Supplementary
143 Table 1). To cover the maximum number and type of genetic variants, we next applied an
144 extensive imputation strategy with four reference panels: the Genome of the Netherlands
145 (GoNL) [16, 17], the UK10K Project [18], the 1000 Genomes Project (1000G) phase 3 [19] and
146 Haplotype Reference Consortium (HRC) [20], and imputed 11.2 M, 11.4 M, 13.1 M, and 11.7
147 M high quality imputed variants (IMPUTE2 [21] info score ≥ 0.7 and minor allele frequency

148    [MAF] ≥ 0.001) with each panel, respectively. After combining the results of the four

149    reference panels by choosing, for each variant, the panel that provided the highest

150    imputation accuracy, we retained a total of 16,059,686 variants covering all the autosomes

151    and the X chromosome (Figure 1a). This strategy was particularly powerful to impute 2.6 M

152    and 5.5 M high quality, low-frequency (0.05 > MAF > 0.01) and rare variants (0.01 > MAF >

153    0.001), respectively, as well as 1.6 M indels. Note that as many as 684,393 common

154    variants (MAF ≥ 0.05), 255,106 low-frequency, 1.7 M rare, and all indels (1.6 M) would be

155    missed if only the HRC reference panel was used. This highlights the benefit of combining

156    different reference panels for comprehensive association testing (Figure 1b).

157    We next tested all the 16 M variants for association with the 22 conditions available in the

158    GERA cohort considering the entire genome and five different inheritance models

159    (Supplementary Figure 1-22). This analysis identified 94 independent loci associated with 17

160    phenotypes at a genome-wide significance level ($p$ < 5.0 × $10^{-8}$) of which 63 for 14

161    phenotypes were also experiment-wide significant ($p$ < 2.0 × $10^{-8}$) after considering

162    correction for the different models of inheritance (see methods) (Supplementary Table 2).

163    According to the GWAS catalog, 68 of the 94 genome-wide significant loci had been

164    previously reported to be associated with the same disease (Supplementary Table 3),

165    whereas 26 of them correspond to previously unreported loci with associations across 16

166    phenotypes (Table 1). Of these new loci, 16 correspond to common, 3 to low-frequency, and

167    7 to rare variants. Interestingly, only a fraction of the 26 new loci would have been genome-

168    wide significant by using individual imputation panels (Figure 1c), namely 20/26 using HRC,

169    14/26 using 1000G Phase 3, 14/26 using UK10K or 15/26 using GoNL. In addition, the lead

170    marker for three of the novel signals is an indel, further confirming the benefits of combining

171    multiple panels with our approach.

172

**Identification of recessive variants with large effects**

173

174    The implementation of refined GWAS strategies not only increases the number of

175    associated variants, but also allows the identification of loci with large impact on the disease.

176    Among the variants that were not detected under the additive model, and hence are

177    expected to be missed by the majority of current GWAS, we highlight three variants with

178    remarkably large recessive effects. First, an intronic indel in the *CACNB4* gene,

179    rs201654520, associated with a nearly twenty-fold increase in risk for cardiovascular

180    disease (MAF= 0.017, OR [CI 95%] = 19.0 [5.5 - 65.8], $p$ = 4.3 × $10^{-8}$). *CACNB4* encodes the

181    β4 subunit of the voltage-dependent calcium channel. This subunit contributes to the flux of

182    calcium ions into the cell by increasing peak calcium current and triggering muscle

183 contraction. Interestingly, an intronic single nucleotide polimorphism (SNP) within *CACNB4*,
184 rs150793926, was associated with idiopathic dilated cardiomyopathy in African Americans
185 [22], but this variant is not in linkage disequilibrium (LD) with rs201654520 (LD $r^2$ [23] = 0.0016).

186 A second recessive variant with large effect, rs77704739, near the *PELO* gene, is
187 associated with a four-fold risk for type 2 diabetes (MAF= 0.036, OR [CI 95%] = 4.3 [2.7 -
188 6.9], $p$ = 1.75 × $10^{-8}$). We also found this variant associated with type 2 diabetes (OR-
189 recessive [95% CI] = 1.9 [1.4 - 2.6], p = 4.95 × $10^{-4}$) and metformin use (OR-recessive [95%
190 CI] = 2.3 [1.6 - 3.4], p = 3.8 × $10^{-5}$) in the UK Biobank [24] (Supplementary Table 4), also only
191 under the recessive model. An independent signal that is about 112 K base pairs away
192 (rs870992, LD $r^2$ = 0.0009) was previously associated with type 2 diabetes in the
193 Greenlandic population, also with a recessive effect [10]. To provide insights into the
194 underlying molecular mechanisms in disease, we interrogated comprehensive catalogues of
195 genetic effects on gene expression; eQTLGen Consortium [25] and GTEx [26]. The rs77704739
196 variant was significantly associated with gene expression of *PELO* in multiple tissues,
197 including diabetes-relevant tissues such as adipose tissue, skeletal muscle, and pancreas.
198 Colocalization analyses showed a probability higher than 0.8 in several tissues, including
199 subcutaneous adipose tissue and skeletal muscle, suggesting this gene as the effector
200 transcript (Figure 3a, 3b, and Supplementary Table 5). In addition, we found that the lead
201 variants in the *PELO* locus overlap with active promoter annotations in human pancreatic
202 islets and open chromatin sites highly-bounded by islet specific transcription factors [27, 28]
203 (Figure 3c).

204 Third, a rare indel, rs557998486, located near the *THUMPD2* gene, is associated with age-
205 related macular degeneration (MAF= 0.009, OR = 10.5, $p$ = 2.75 × $10^{-8}$). Also under the
206 recessive model in UK Biobank, this indel was associated with age-related macular
207 degeneration (OR [CI 95%] = 7.6 [1.5-37.3], $p$ = 4.1 × $10^{-2}$), eye surgery (beta [CI 95%] = 1.6
208 [0.6-2.6], $p$ = 1.17 × $10^{-3}$) (Supplementary Table 4), and C-reactive protein, a known
209 biomarker for macular degeneration [29] (beta [CI 95%] =1.1 [0.7 - 1.5], $p$ = 1.15 × $10^{-4}$)
210 (Supplementary Table 6). Interestingly, the fact that we found no SNPs in LD with this lead
211 indel further confirms the benefits of multiple reference panel imputation strategies that
212 include alternative forms of variation. The lead indel rs557998486 overlaps DNAse I
213 hypersensitivity sites in retinal and iris cell lines [30], highlighting a candidate open chromatin
214 region that is also predicted to be an enhancer assigned to the *THUMPD2* gene according to
215 GeneHancer [31]. One of the variants with the highest LD with rs557998486 (rs116649730, LD
216 $r^2$= 0.32) is associated with reduced expression of its nearest gene, *THUMPD2* (Z-score = -
217 4.85, $p$ = 1.25 × $10^{-6}$), according to eQTLGen Consortium data.

218

7

**Replication using UK Biobank and FinnGen**

We sought replication of previously unreported loci using UK Biobank, a prospective cohort of ~500 K individuals aged between 40 to 69 [24]. Given the high heterogeneity in phenotype definitions in UK Biobank compared to GERA, we tested for replication with the same phenotype and related traits (Supplementary Table 4). Compared to GERA, some of the conditions may not be ascertained or have an age at onset later than the average age at ascertainment in UK Biobank (56.52 years [32]) which could affect the replication success. Despite these limitations, we tested the novel variants using the corresponding inheritance model, and replicated 4 new loci with the same phenotype (Table 2).

We further sought replication of the association within the *CACNB4* gene with cardiovascular disease in FinnGen, a cohort of ~218 K Finnish individuals with an average age of 63, as it includes individuals with a higher average age (63 vs 56 in UK Biobank) and the risk of developing a cardiovascular disease is well-known to increase with age [33]. In addition, FinnGen has a precise and richer classification of this particular phenotype than UK Biobank. In brief, we tested rs201654520 for association with 47 cardiovascular endpoints. Of all the conditions tested, four (hypertensive heart disease, hypertensive heart and/or renal disease, heart failure, and right bundle-branch block) were nominally associated ($p < 0.05$). All the associations had a direction of effect consistent with the effect observed in the GERA cohort (Supplementary Figure 23a). Although there is a high heterogeneity in the phenotype definitions between cohorts, we meta-analyzed the results from these endpoints from FinnGen with the result from "cardiovascular disease" phenotype from GERA, but none of them reach the genome-wide significance (see Methods) (Supplementary Figure 23). We did not include UK Biobank in this meta-analysis as the equivalent phenotypes were not available or had less than 350 cases in UK Biobank, therefore, underpowered for a recessive analysis. Notably, when analyzing the association of rs201654520 with related quantitative traits we found that those who were homozygous for the high-risk allele had lower systolic blood pressure ($p = 4.1 \times 10^{-3}$, beta = -0.23) (Supplementary Table 4). While lower systolic blood pressure has been associated with increased risk of myocardial infarction in particular circumstances, this is not the typical direction of association, and therefore merits additional study [34].

We also sought replication of the recessive association of rs557998486 near *THUMPD2* gene with macular degeneration in FinnGen. While rs557998486 was associated with increased risk of macular degeneration in UK Biobank under the recessive model (OR [CI 95%] = 7.6 [1.5-37.3], $p = 4.1 \times 10^{-2}$), it was not significantly associated in the FinnGen biobank although it showed the same direction of effect. However, the meta-analysis did not

254    reach the genome-wide significance (rs557998486 $p$ = 9.6 × $10^{-6}$) and had a high

255    heterogeneity (heterogeneity $I^2$ = 87.1, heterogeneity $p$ = 4.3 × $10^{-4}$).

256

**Detection ranges of the different inheritance models**

258    Our findings provide an empirical overview of the detection range of five different inheritance

259    models, and show how each of them captures a fraction of the genetic variants associated

260    with complex traits. As indicative of the power of current genetic studies that usually only

261    consider additive allelic effects, we found three different scenarios. Among all the 94

262    associated loci identified, 12 showed genome-wide significance only under the additive

263    model, 62 under both additive and non-additive models, and 20 showed genome-wide

264    significance only when non-additive tests were applied (Figure 2a). To further classify these

265    variants, we tested whether any of the 62 variants associated with both additive and non-

266    additive models deviate from additivity through a dominance deviation test [9]. Eleven of these

267    62 variants (17.7%) showed significant deviation from additivity (dominance deviation test $p$

268    < 0.05). Altogether, the dominance deviation test over the 93 autosomal loci identified 62

269    additive (66%) and 24 non-additive associations (25.5%) and 8 undetermined. Based on the

270    smallest GWAS $p$-value, we further classified non-additive associations into 9 recessive, 13

271    dominant, 8 heterodominat and 7 genotypic (Supplementary Table 2).

272    We also found that each of the available models for association testing has a different range

273    of detection. To identify the 94 genome-wide associated loci, the additive test, as expected,

274    was the most sensitive model (74 loci), followed by the genotypic (59 loci), the dominant (56

275    loci), the recessive (43 loci) and the heterodominant (32 loci). When considering known loci,

276    48 of the 68 previously reported loci were identified by more than one model in our analysis,

277    and almost half of these (22 loci) with all five models. In contrast, of the 26 newly discovered

278    variants, only 8 were identified with multiple models, whereas the majority of them (18 loci),

279    were detected only with the additive (6 loci), the genotypic (4 loci), the recessive (4 loci) and

280    the dominant (3 loci) model. Of note, 13 out of 26 (50%) novel loci were only identified by

281    non-additive models.

282    To further investigate to what extent the additive model captures non-additive signals, and

283    how much this depends on sample size, we carried out power calculations on loci that were

284    detected here only under a non-additive model, such as rs201654520 within *CACNB4* gene

285    and rs77704739 near the *PELO* gene. These power calculations showed that the additive

286    test would require a population sample size of at least 370,646 individuals to detect the

287    recessive association of rs201654520 in *CACNB4* (Figure 2b), and at least 188,637

288    individuals to capture the recessive signal of rs77704739 near the *PELO* gene (Figure 2c),

289 while the population sample size required for the recessive model was only 21,021 and
290 67,611, respectively. In this study, we were able to identify both associations with a modest
291 sample size by using the most well-suited disease model.

292

293 **The GUIDANCE framework**

294 We developed an integrated framework including our methodology used to analyze the
295 GERA cohort, called GUIDANCE. GUIDANCE allows the analysis of genome-wide
296 genotyped data in a single execution in distributed computing infrastructures without the
297 need for extensive computational expertise or constant user intervention. The GUIDANCE
298 workflow requires quality-controlled genotyped data as an input and provides association
299 results, graphical outputs and statistical summaries. Integrating state-of-the-art tools with in-
300 house code written in Java, bash and R [35], GUIDANCE efficiently performs large-scale
301 GWAS, including 1) the pre-phasing of haplotypes, 2) the imputation of genotypes using
302 multiple reference panels, 3) the association testing for different inheritance models and
303 integrating results from different panels, 4) a cross-phenotype analysis when more than one
304 phenotype is available in the cohort (Supplementary Table 7), and finally, 5) the generation
305 of summary statistics tables and graphic representations of the results (Supplementary
306 Figure 24), for both the autosomes and the X chromosome. While GUIDANCE can be
307 executed as a standalone compact program it can also be used in modules (Supplementary
308 Figure 25), which makes GUIDANCE adaptable to existing frameworks and provides an
309 even higher level of control to users.

310 GUIDANCE runs in distributed computing platforms, including the cloud, without requiring a
311 broad background in distributed environments. This is feasible since GUIDANCE was
312 implemented on top of the COMP Superscalar Programming Framework (COMPSs) [36]. With
313 COMPSs, the GUIDANCE workflow was implemented as a sequential Java program
314 containing the calls to the GWAS tools, encapsulated in Java methods, and selected as
315 tasks, while COMPSs controls the execution of those tasks on the underlying distributed
316 infrastructure. The source code, the pre-compiled binaries and documentation to use
317 GUIDANCE are available at http://cg.bsc.es/guidance.

318

319 **Discussion**

320 The increasingly large sample sizes in GWAS improve the statistical power to identify
321 genetic variants associated with complex diseases. At the same time, the emergence of
322 larger and denser reference panels allows genotype imputation at lower ranges of allele

323  frequencies previously unexplored. In this study, we demonstrate the value of applying a
324  comprehensive GWAS including denser imputation strategies, the X chromosome and non-
325  additive association tests to an existing large-scale genetic resource, the GERA cohort. We
326  show that by applying more powerful imputation protocols we increased the number and the
327  type of variants tested for association, including low-frequency and rare SNPs as well as
328  alternative forms of variation, such as indels. Our analysis in the GERA cohort shows that
329  between 13 and 20 of the genome-wide significant associations (14-21%) would not have
330  been identified when using a single reference panel. Likewise, our analysis in the GERA
331  cohort demonstrates that 21% of the associations would be missed by only testing the
332  additive model. Overall, 27.6% of associations would not have been identified by applying
333  the commonly used HRC and additive model association testing.

334  We here show the potential of identifying very large effect recessive associations by
335  maximizing the use of current reference panels and testing different inheritance models, as
336  exemplified by the associations with type 2 diabetes, cardiovascular disease and macular
337  degeneration with variants near *PELO*, *CACNB4*, and *THUMPD2*, respectively. This strategy
338  opens new avenues for future analyses in large scale biobanks, as demonstrated with our
339  power calculations, which show that even the largest available GWAS meta-analyses or
340  biobanks would not have enough power to identify these associations using only the additive
341  model. For example, the *CACNB4* gene, associated with cardiovascular disease, would
342  require a sample size equivalent to 370,000 individuals when using the additive test, 17
343  times larger than the required sample size under a recessive analysis. After considering all
344  the supporting evidence illustrated with many examples in this study, the results suggest that
345  this new associations deserve future validations and follow-up analysis, and demonstrate the
346  importance of a comprehensive analysis including non-additive models when performing
347  GWAS.

348  The inclusion of non-additive associations can also have an impact on the construction of
349  polygenic risk scores. Current polygenic scores (PRS) are calculated summing risk alleles
350  weighted by effect sizes from GWAS results, which have typically tested only the additive
351  model in the association test. Hence, large-scale genome-wide association data accounting
352  for different models of inheritance and including both SNPs and alternative forms of
353  variation, such as indels, will also be essential to develop more accurate genome-wide PRS,
354  which would weight each of the genotype carriers appropriately, rather than weighting the
355  heterozygous half-way between the homozygous of the effect and alternate alleles.

356  To easily apply this strategy to genetic studies we present GUIDANCE, a standalone and
357  easy-to-use application that allows an efficient and comprehensive GWAS analysis in
358  different computing platforms, such as cloud and high-performance computing architectures.

11

359 In a moment where the community is facing computational and methodological challenges
360 due to the growing complexity and size of genetic datasets, the availability of robust and
361 complete analysis platforms can improve the efficiency of genetic studies, standardizing
362 analysis strategies among large meta-analysis cohorts to ensure consistency.

363 Finally, to share our results with the community and to promote the analysis of non-additive
364 inheritance models in GWAS, a public searchable database including additive and non-
365 additive summary statistics for 16 M of variants and 22 phenotypes is available at the Type 2
366 Diabetes Knowledge Portal (http://www.type2diabetesgenetics.org and full summary
367 statistics at http://cg.bsc.es/guidance).

368

369 **Online Methods**

370

371 **GUIDANCE Workflow Description**

372 By combining and integrating state-of-the-art GWAS analysis tools into the COMP
373 Superscalar programming Framework (COMPSs), we developed GUIDANCE, a standalone
374 application that performs haplotype phasing, genome-wide imputation, association testing
375 and PheWAS analysis of large GWAS datasets (Supplementary Figure 24).

376 As shown in Supplementary Figure 24, GUIDANCE's workflow starts with quality-controlled
377 genotype data and ends with providing association results, graphical outputs and statistical
378 summaries.

379 Once everything is settled in the GUIDANCE configuration file, GUIDANCE performs an
380 efficent two-stage imputation procedure, by pre-phasing the genotypes into whole
381 haplotypes followed by genotype imputation itself [21]. SHAPEIT2 [37] or EAGLE2 [38] and
382 IMPUTE2 [39] or MINIMAC4 [40] can be used for pre-phasing and genotype imputation,
383 respectively. In addition, GUIDANCE accepts one or multiple reference panels, allowing the
384 integration of the results obtained from all panels by selecting for each variant the genotypes
385 from the reference panel that provides the highest imputation accuracy according to the
386 IMPUTE2 info score or MINIMAC2 $r^2$ (Supplementary Figure 26). GUIDANCE also performs
387 a post-imputation quality control to eliminate low-quality imputed variants under the basis of
388 the IMPUTE2 info score or MINIMAC2 $r^2$ and the MAF.

389 After genotype imputation and post-imputation quality control, GUIDANCE applies
390 SNPTEST for association testing, where additive, dominant, recessive, heterodominant and
391 genotype models can be analyzed. Here, the user can decide to include several covariates

12

392    for the association test, such as principal components to adjust for population stratification,
393    or any other confounders. GUIDANCE also allows testing for multiple phenotypes or for a
394    single phenotype with different covariates in the same execution. After association testing,
395    variants are filtered by the deviation from Hardy-Weinberg equilibrium (HWE) *p*-value.
396    Finally, GUIDANCE generates summary reports for each trait with all the inheritance models
397    tested in the association and the corresponding graphical representation, i.e., Manhattan
398    and Quantile-Quantile (Q-Q) plots (Supplementary Figure 1-22), also providing a matrix
399    identifying cross-phenotype associations (Supplementary Table 7).

400    GUIDANCE can be executed as a a standalone compact program or as independent
401    modules (see Supplementary Figure 25 for a list of independent modules) to facilitate the
402    use of GUIDANCE into existing frameworks.

403    Further details can be found in the configuration file from the GUIDANCE execution at
404    http://cg.bsc.es/guidance. Specific documentation to use this framework is available at
405    http://cg.bsc.es/guidance, as well as the source code and the pre-compiled binaries that are
406    available in the "download" section.

407

408    **The Analysis of GERA cohort**

409    **GERA cohort Description**

410    GERA cohort data was obtained through dbGaP under accession phs000674.v1.p1. For
411    further information about the specific phenotypes (ICD-9-CM codes) included in GERA,
412    please visit its website on dbGaP (https://www.ncbi.nlm.nih.gov/gap/). The Resource for
413    Genetic Epidemiology Research on Aging (GERA) Cohort was created by a RC2 "Grand
414    Opportunity" grant that was awarded to the Kaiser Permanente Research Program on
415    Genes, Environment, and Health (RPGEH) and the UCSF Institute for Human Genetics
416    (AG036607; Schaefer/Risch, PIs). The RC2 project enabled genome-wide SNP genotyping
417    (GWAS) to be conducted on a cohort of over 100 K adults who were members of the Kaiser
418    Permanente Medical Care Plan, Northern California Region (KPNC), and participating in its
419    RPGEH. The resulting GERA cohort is composed of 42% of males, 58% of females, and
420    ranges in age from 18 to over 100 years old with an average age of 63 years at the time of
421    the RPGEH survey (2007). 19% of the individuals are from non-European ancestry, while
422    81% are described as white non-Hispanic participants. After an explicit requirement of
423    consent by email, data from 78,486 participants was deposited in dbGaP, with similar
424    demographic characteristics to those of the initial genotyped cohort.

425

13

**Quality Control**

A subset of 62,281 subjects of European ancestry underwent quality control analyses. A 3-step quality control protocol was applied using PLINK [41, 42], and included 2 stages of SNP removal and an intermediate stage of sample exclusion.

The exclusion criteria for genetic markers consisted of: proportion of missingness ≥ 0.05, HWE $p \leq 1 \times 10^{-20}$ for all the cohort, and MAF < 0.001. This protocol for genetic markers was performed twice, before and after sample exclusion.

For the individuals, we considered the following exclusion criteria: gender discordance, subject relatedness (pairs with PI-HAT ≥ 0.125 from which we removed the individual with the highest proportion of missingness), sample call rates ≥ 0.02 and population structure showing more than 4 standard deviations within the distribution of the study population according to the first seven principal components (Supplementary Figure 27). After QC, 56,637 subjects remained for the analysis (Supplementary Table 1).

**Analyzing GERA cohort using GUIDANCE**

GUIDANCE pre-phased the genotypes to whole haplotypes with SHAPEIT2, and then performed genotype imputation with IMPUTE2 using 1000G phase 3, UK10K, GoNL, and HRC as reference panels. After filtering variants with an info score < 0.7 and a MAF < 0.001, we tested for association with additive, dominant, recessive, heterodominant and genotypic logistic regression using SNPTEST, and including seven derived principal components, sex and age as covariates. To maximize power and accuracy, we combined the association results from the four reference panels by choosing for each variant, the genotypes from the reference panel that provided the best IMPUTE2 info score. For chromosome X we restricted the analysis to non-pseudoautosomal (non-PAR) regions and stratified the association analysis by sex to account for hemizygosity for males, while for females, we followed an autosomal model. Finally, we excluded variants with HWE controls $p < 1 \times 10^{-6}$ in the final results.

**Identification of known and new associated loci**

After the association test, GUIDANCE provided a list of variants that passed the *p*-value threshold specified in the configuration file (i.e., $p \leq 5.0 \times 10^{-8}$). Using the "IRanges" R package [43], all the genome-wide significant variants were collapsed into ranges (500 kb) that define each associated locus.

459     To distinguish between known or new associated regions, for each top variant we looked for

460     any proxy variant with an LD $r^2$ > 0.35 in the GWAS catalog (accession 5 September 2019)

461     associated with the same phenotype or a related one (for example, bone mineral density,

462     cholesterol levels or diastolic/systolic blood pressure phenotypes for osteoporosis,

463     dyslipidemia or hypertension, respectively). HLA regions at chromosome 6 were excluded

464     since the particularities of these regions required further detailed studies on their LD pattern.

465     Proxies were selected using LDlink (https://ldlink.nci.nih.gov/) [44].

466     We defined an experiment-wide significant $p$-value cutoff of $p < 2.0 \times 10^{-8}$ by applying the

467     Bonferroni correction for 2.5 effective test ($5.0 \times 10^{-8}$ / 2.5 effective test). This factor of 2.5

468     was obtained from a simulation study when four genetic models (additive, dominant,

469     recessive and genotypic) are used [45] since the genetic models are not independent.

470     However, a new simulation study including the heterodominant model should be done for a

471     more accurate effective number of tests.

472

473     **Replication with UK Biobank**

474     **Phenotype Curation**

475     UK Biobank participants agreed to provide detailed information about their lifestyle,

476     environment and medical history, to donate biological samples (for genotyping and for

477     biochemical assays), to undergo measures and to have their health followed

478     (http://www.ukbiobank.ac.uk/).

479     When collecting and analyzing a wide range of phenotypes from the UK Biobank, a central

480     challenge was the curation and harmonization of the vast array of categorizations, variable

481     scalings, and follow-up responses. Fortunately, to this end, the PHEnome Scan ANalysis

482     Tool (or PHESANT: https://github.com/MRCIEU/PHESANT) [46] performs much of the

483     transformations and recodings required to generate meaningful, interpretable phenotypes.

484     We have made further adjustments based on user feedback, owing to the value of

485     transparency in generating our phenotype guidelines. Applying these changes to the

486     PHESANT source code, phenotypes were parsed using our modified version

487     (github.com/astheeggeggs/PHESANT) on a virtual machine on the Google Cloud Platform.

488     We first restricted to the subset of European individuals, before passing the resultant

489     phenotypic data to PHESANT. The 'variable list' file and 'data-coding' file, whose formats are

490     defined in the original version of PHESANT were updated as new phenotypes were added in

491     the latest UK Biobank release. Re-codings of variables, and inherent orderings of categorical

492    variables, are defined in the 'data-coding' file. The 'Excluded' column of the 'variable list' file

493    defines the collection of variables that we do not wish to interrogate.

494    A high level overview of the PHESANT pipeline, our defaults, and the associated short flags

495    for the phenomescan.r code are displayed in Supplementary Figure 28. In addition to the

496    inverse-rank normalization applied to the collection of continuous phenotypes, we also

497    consider the raw version of the continuous phenotype, with no transformation applied to the

498    data.

499    Curation of the ICD10 codes was carried out separately for computational efficiency. For the

500    ICD10 phenotype, individuals are assigned a vector of ICD10 diagnoses. We truncated

501    these codes to two digits, and assigned each individual to either case or control status for

502    that ICD10 code in turn by checking if their vector contains that code. Throughout, we

503    assumed the data contained no missingness, so the sum of cases and controls throughout

504    was the number of individuals in our 'European' subset of the UK Biobank data. As in the

505    PHESANT categorical (multiple) phenotypes, ICD10 code case/control phenotypes were

506    removed if less than 50 individuals had the diagnosis.

507

**Association testing and meta-analysis for UK Biobank phenotypes**

509    We performed the association testing for the curated phenotypes as implemented in

510    SNPTEST for additive, dominant, recessive, heterodominant and genotypic inheritance

511    models, as it has been described in the "Analyzing GERA cohort using GUIDANCE" section.

512    For all genotypic variants identified in the discovery stage, we assigned the recessive model

513    after we identified it as the underlying model.

514    After the association testing, we filtered and ordered all the phenotypes based on the *p*-

515    value for the best model of inheritance obtained from the GERA cohort analysis, with special

516    consideration to equivalent phenotypes or related traits.

517    With the association testing results of both GERA cohort and UK Biobank, we meta-

518    analyzed the results using METAL [47]. We use the inverse variance-weighted fixed effect

519    model for all the variants except for the rs557998486 variant associated with macular

520    degeneration, since its *beta*, calculated with the "em" method from SNPTEST, was inflated.

521    Therefore, we performed a sample size based meta-analysis, which converts the direction of

522    the effect and the *p*-value into a z-score.

523    For biomarkers, only the results from the first visit were taken into account since less than

524    10% of the cases where present in the second visit.

525

**Association testing and meta-analysis with FinnGen**

We used SAIGE [48] for recessive association testing using sex, age, PC1-10 and batch as covariates. We analyzed FinnGen release 5 that contains 218,792 individuals with a median age 62.6 and a mean age 59.8.

For the cardiovascular disease endpoints, we meta-analyzed the results using "rmeta" R package [49]. For macular degeneration, we meta-analyzed the results using METAL as described in the previous section.


**Dominance deviation test**

To detect genuine differences between additive and non-additive signals, we performed a dominance deviation test for all 93 autosomal genome-wide significant loci.

Dominance deviation was tested by a logistic regression analysis using PLINK (v1.90b6.9, www.cog-genomics.org/plink/1.9/). Sex, age and the first 7 PCs were included as covariates.


**Definition of 99% credible set of *PELO* locus**

For the *PELO* locus, the fraction of aggregated variants that have a 99% probability of containing the causal one was identified. The 99% credible set of variants for the region was defined with a Bayesian refinement approach [50], considering variants with an $r^2 > 0.1$ with the leading one.

For each variant within the *PELO* locus, the credible set provides a posterior probability of being the causal one [50]. The approximate Bayes factor (*ABF*) for each variant was estimated as

$$ABF = \sqrt{1 - r} \, e^{(rz^2/2)} \ ,$$

where

$$r = \frac{0.04}{(SE^2 + 0.04)},$$

$$z = \frac{\beta}{SE}.$$

The β and the SE result from a logistic regression model testing for association. The posterior probability for each variant was calculated as

$$Posterior \ Probability_i = \frac{ABF_i}{T},$$

17

551  where *ABFi* corresponds to the approximate Bayes' factor for the marker *i*, and *T* represents
552  the sum of all the *ABF* values enclosed in the interval. As commonly employed by
553  SNPTEST, this calculation assumes that the prior of the *β* is a Gaussian with mean 0 and
554  variance 0.04.

555  Finally, the cumulative posterior probability was calculated after ranking the variants
556  according to the *ABF* in decreasing order. Variants were included in the 99% credible set of
557  the region until the cumulative posterior probability of association got over 0.99.

558

### Gene expression and functional characterization

559

560  The eQTLGen Consortium (https://www.eqtlgen.org/cis-eqtls.html, last access on July 2019)
561  and GTEx portal (https://gtexportal.org/, last access on July 2019) were used to find
562  associations between our novel genetic associations and gene expression. When the variant
563  was not available in the resources, a proxy SNP was used instead.

564  To determine whether any identified overlap between GERA GWAS loci and eQTLGen or
565  GTEx eQTLs was due to a true shared association signal, we performed a colocalization
566  analysis. Colocalization was assessed by a Bayesian test using summary statistics from
567  both studies [51]; summary statistics from the *cis* eQTLGen and GTEx were downloaded from
568  the eQTLGen website and GTEx portal, respectively. The test was performed using the R
569  package coloc v3.2-1 [51, 52, 53]. The test provided a posterior probability for the GWAS locus
570  and the eQTL to share the same causal variant(s).

571  We integrated available epigenomic datasets to examine the role of human pancreatic islet
572  transcriptional regulation underlying rs77704739 association with type 2 diabetes. We used
573  the WashU EpiGenome Browser (http://epigenomegateway.wustl.edu/browser/, last access
574  on July 2019) and previously published RNA-seq, ATAC-seq and ChIP-seq assays of
575  H3K4me3, H3K27ac, Mediator, CTCF and islet transcription factors (FOXA2, MAFB,
576  NKX2.2, NKX6.1 and PDX1) in human pancreatic islets [27, 28] and islet regulome annotations
577  [28].

578

### Data Availability

579

580  The complete summary statistics are deposited at the Type 2 Diabetes Knowledge portal
581  (www.type2diabetesgenetics.org/) and can be also accessed from http://cg.bsc.es/guidance.
582  GUIDANCE is also available at http://cg.bsc.es/guidance.

583

19

619    The data used for the analyses described in this manuscript were obtained from the GTEx
620    Portal on 07/16/2019. We acknowledge PRACE for awarding us access
621    to both MareNostrum supercomputer from the Barcelona Supercomputing Center, based in
622    Spain at Barcelona, and the SuperMUC supercomputer of the Leibniz Supercomputing
623    Centre (LRZ), based in Garching at Germany (proposals numbers 2016143358 and
624    2016163985). The technical support group from the Barcelona Supercomputing Center is
625    gratefully acknowledged. Finally, we thank all the Computational Genomics group at the
626    BSC for their helpful discussions and valuable comments on the manuscript. We also
627    acknowledge Elias Rodriguez Fos for designing the GUIDANCE logo.

628

## Authors Contributions

630    M.G-M., R.A., J.M.M., and D.T. conceived, planned, and performed the main analyses. M.G-
631    M., J.M.M., and D.T. wrote the manuscript. M.G-M., R.A., M.P., C.R-C., F.S., J.E., C.D.,
632    E.T., and R.M.B. developed GUIDANCE. S.B-G. designed and performed the quality control.
633    S.B-G and I. M-E. performed the functional characterization. C.S performed the dominance
634    deviation test and the gene expression analysis. J.M.M., C.E.C., J.B.C, E.A., A.L., K.A.,
635    D.P., and J.C.F. contributed with UK Biobank data and analysis. S.R. and M.K. contributed
636    with FinnGen data and analysis. J.M.M. and D.T. designed and supervised the study. All
637    authors reviewed and approved the final manuscript.

638

## References

640    1.    Welter D*, et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait
641       associations. *Nucleic Acids Res* **42**, D1001-1006 (2014).
642
643    2.    Manolio TA*, et al.* Finding the missing heritability of complex diseases. *Nature* **461**,
644       747-753 (2009).
645
646    3.    Bonas-Guarch S*, et al.* Re-analysis of public genetic data reveals a rare X-
647       chromosomal variant associated with type 2 diabetes. *Nat Commun* **9**, 321 (2018).
648
649    4.    Tam V, Patel N, Turcotte M, Bosse Y, Pare G, Meyre D. Benefits and limitations of
650       genome-wide association studies. *Nat Rev Genet*, (2019).
651
652    5.    Zhu Z*, et al.* Dominance genetic variation contributes little to the missing heritability
653       for human complex traits. *Am J Hum Genet* **96**, 377-385 (2015).
654
655    6.    Salanti G*, et al.* Underlying genetic models of inheritance in established type 2
656       diabetes associations. *Am J Epidemiol* **170**, 537-545 (2009).
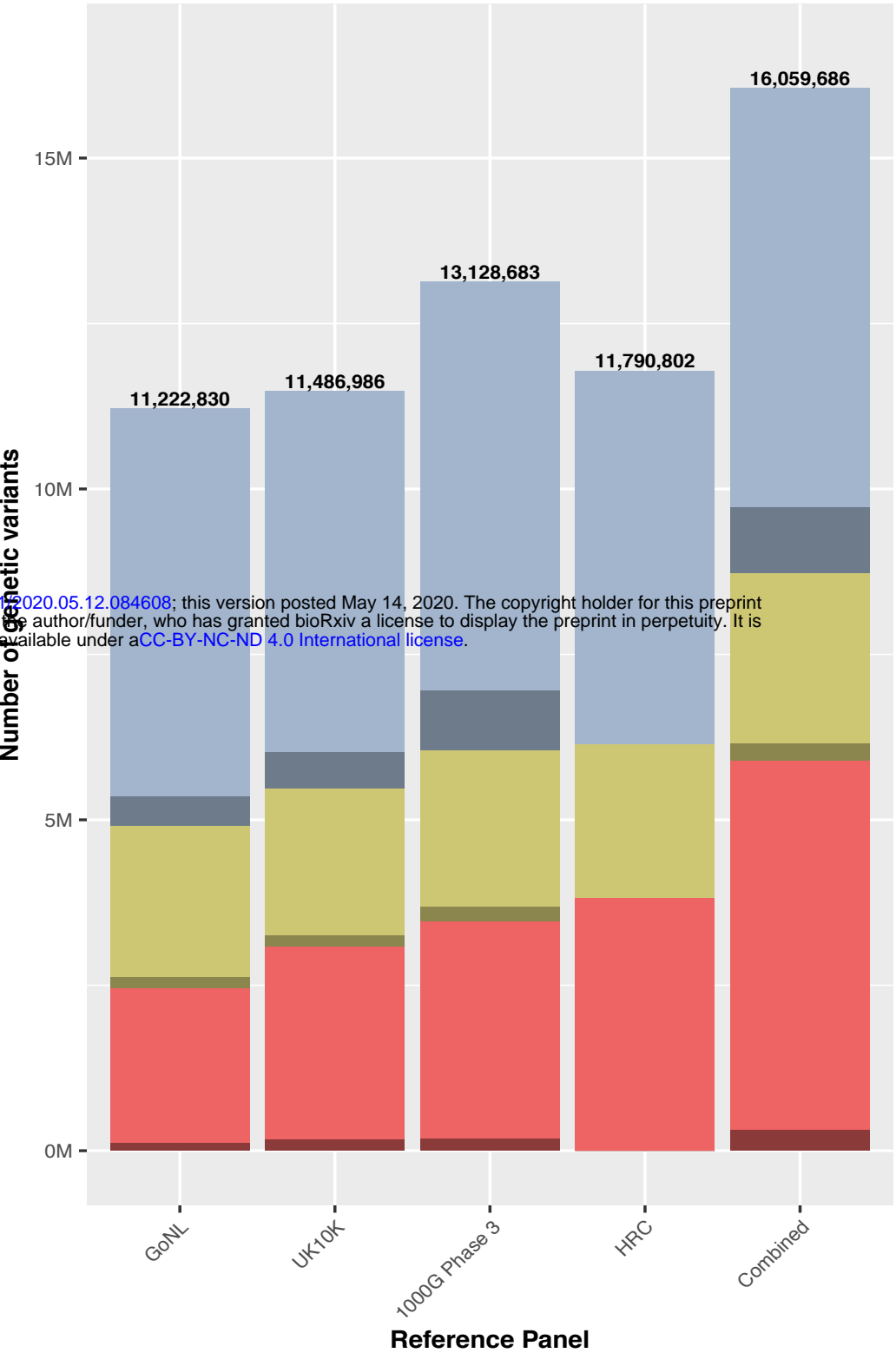657

7.   Lettre G, Lange C, Hirschhorn JN. Genetic model testing and statistical power in population-based association studies of quantitative traits. *Genet Epidemiol* **31**, 358-362 (2007).

8.   Antonarakis SE, Beckmann JS. Mendelian disorders deserve more attention. *Nat Rev Genet* **7**, 277-282 (2006).

9.   Wood AR*, et al.* Variants in the FTO and CDKAL1 loci have recessive effects on risk of obesity and type 2 diabetes, respectively. *Diabetologia* **59**, 1214-1221 (2016).

10.  Grarup N*, et al.* Identification of novel high-impact recessively inherited type 2 diabetes risk variants in the Greenlandic population. *Diabetologia* **61**, 2005-2015 (2018).

11.  Moltke I*, et al.* A common Greenlandic TBC1D4 variant confers muscle insulin resistance and type 2 diabetes. *Nature* **512**, 190-193 (2014).

12.  Steinthorsdottir V*, et al.* A variant in CDKAL1 influences insulin response and risk of type 2 diabetes. *Nat Genet* **39**, 770-775 (2007).

13.  Lenz TL*, et al.* Widespread non-additive and interaction effects within HLA loci modulate the risk of autoimmune diseases. *Nat Genet* **47**, 1085-1090 (2015).

14.  Goyette P*, et al.* High-density mapping of the MHC identifies a shared role for HLA-DRB1*01:03 in inflammatory bowel diseases and heterozygous advantage in ulcerative colitis. *Nat Genet* **47**, 172-179 (2015).

15.  Hoffmann TJ*, et al.* Imputation of the rare HOXB13 G84E mutation and cancer risk in a large population-based cohort. *PLoS Genet* **11**, e1004930 (2015).

16.  Boomsma DI*, et al.* The Genome of the Netherlands: design, and project goals. *Eur J Hum Genet* **22**, 221-227 (2014).

17.  Genome of the Netherlands Consortium. Whole-genome sequence variation, population structure and demographic history of the Dutch population. *Nat Genet* **46**, 818-825 (2014).

18.  UK10K Consortium*, et al.* The UK10K project identifies rare variants in health and disease. *Nature* **526**, 82-90 (2015).

19.  1000 Genomes Project Consortium*, et al.* An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56-65 (2012).

20.  McCarthy S*, et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet* **48**, 1279-1283 (2016).

21.  Howie B, Fuchsberger C, Stephens M, Marchini J, Abecasis GR. Fast and accurate genotype imputation in genome-wide association studies through pre-phasing. *Nat Genet* **44**, 955-959 (2012).

22.  Xu H*, et al.* A Genome-Wide Association Study of Idiopathic Dilated Cardiomyopathy in African Americans. *J Pers Med* **8**,  (2018).

23.  Hill WG, Robertson A. Linkage disequilibrium in finite populations. *Theor Appl Genet* **38**, 226-231 (1968).
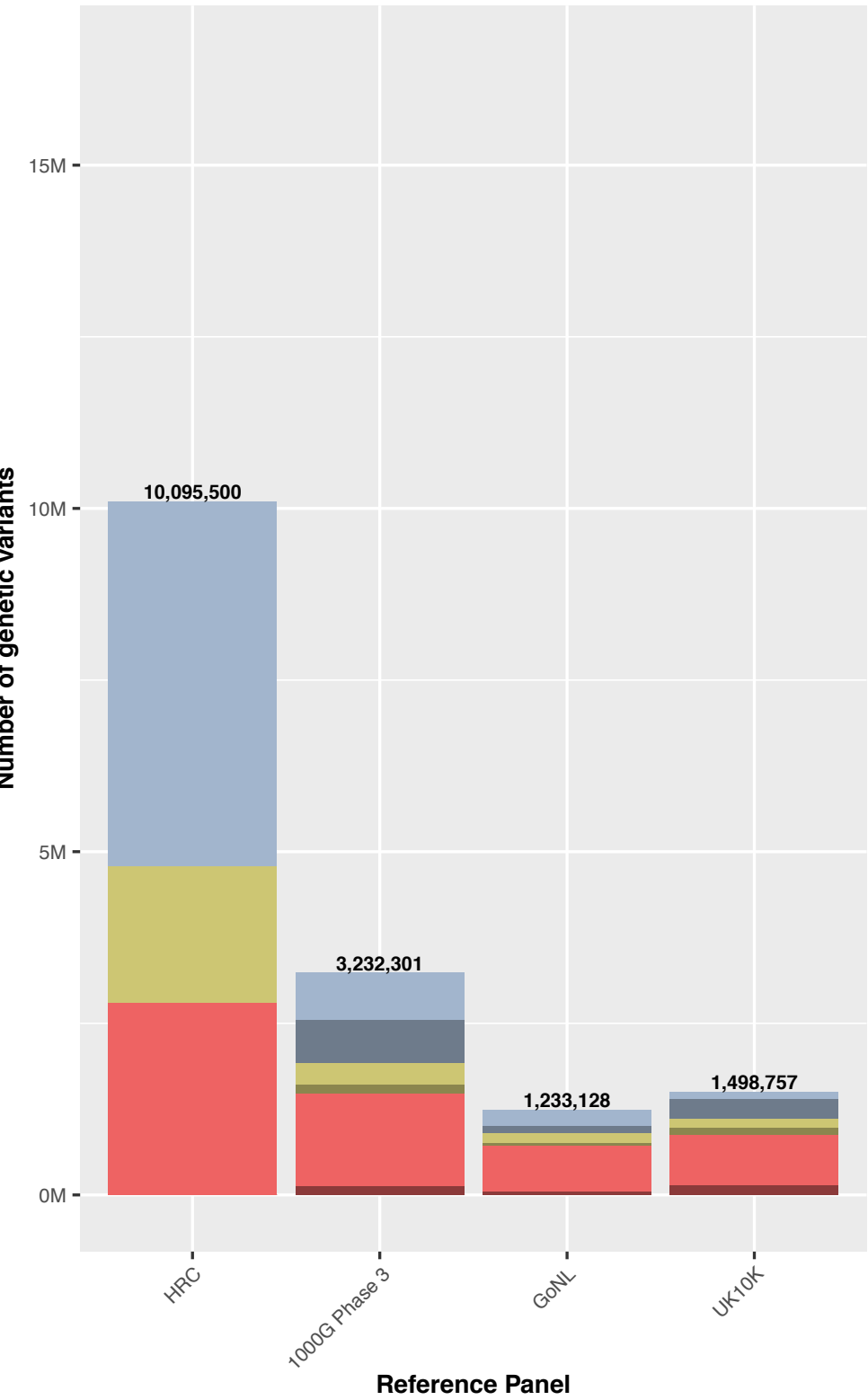
713

714   24.   Bycroft C*, et al.* The UK Biobank resource with deep phenotyping and genomic data.
715         *Nature* **562**, 203-209 (2018).
716

717   25.   Võsa U*, et al.* Unraveling the polygenic architecture of complex traits using blood
718         eQTL metaanalysis. *bioRxiv*, 447367 (2018).
719

720   26.   GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. *Nat Genet* **45**,
721         580-585 (2013).
722

723   27.   Pasquali L*, et al.* Pancreatic islet enhancer clusters enriched in type 2 diabetes risk-
724         associated variants. *Nat Genet* **46**, 136-143 (2014).
725

726   28.   Miguel-Escalada I*, et al.* Human pancreatic islet three-dimensional chromatin
727         architecture provides insights into the genetics of type 2 diabetes. *Nat Genet* **51**,
728         1137-1148 (2019).
729

730   29.   Molins B, Romero-Vazquez S, Fuentes-Prior P, Adan A, Dick AD. C-Reactive Protein
731         as a Therapeutic Target in Age-Related Macular Degeneration. *Front Immunol* **9**, 808
732         (2018).
733

734   30.   Consortium EP. The ENCODE (ENCyclopedia Of DNA Elements) Project. *Science*
735         **306**, 636-640 (2004).
736

737   31.   Fishilevich S*, et al.* GeneHancer: genome-wide integration of enhancers and target
738         genes in GeneCards. *Database (Oxford)* **2017**,  (2017).
739

740   32.   Hewitt J, Walters M, Padmanabhan S, Dawson J. Cohort profile of the UK Biobank:
741         diagnosis and characteristics of cerebrovascular disease. *BMJ Open* **6**, e009161
742         (2016).
743

744   33.   Yazdanyar A, Newman AB. The burden of cardiovascular disease in the elderly:
745         morbidity, mortality, and costs. *Clin Geriatr Med* **25**, 563-577, vii (2009).
746

747   34.   Vidal-Petiot E*, et al.* Cardiovascular event rates and mortality according to achieved
748         systolic and diastolic blood pressure in patients with stable coronary artery disease:
749         an international cohort study. *Lancet* **388**, 2142-2152 (2016).
750

751   35.   R Foundation for Statistical Computing. R: A Language and Environment for
752         Statistical Computing.) (2019).
753

754   36.   Lordan F*, et al.* ServiceSs: an interoperable programming framework for the Cloud.
755         *Journal of Grid Computing* **12**, 67-91 (2014).
756

757   37.   Delaneau O, Zagury JF, Marchini J. Improved whole-chromosome phasing for
758         disease and population genetic studies. *Nat Methods* **10**, 5-6 (2013).
759

760   38.   Loh PR*, et al.* Reference-based phasing using the Haplotype Reference Consortium
761         panel. *Nat Genet* **48**, 1443-1448 (2016).
762

763   39.   Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation
764         method for the next generation of genome-wide association studies. *PLoS Genet* **5**,
765         e1000529 (2009).
766

767  40.  Das S*, et al.* Next-generation genotype imputation service and methods. *Nat Genet*
768       **48**, 1284-1287 (2016).
769

770  41.  Purcell S*, et al.* PLINK: a tool set for whole-genome association and population-
771       based linkage analyses. *Am J Hum Genet* **81**, 559-575 (2007).
772

773  42.  Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation
774       PLINK: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 7 (2015).
775

776  43.  Lawrence M*, et al.* Software for computing and annotating genomic ranges. *PLoS*
777       *Comput Biol* **9**, e1003118 (2013).
778

779  44.  Machiela MJ, Chanock SJ. LDlink: a web-based application for exploring population-
780       specific haplotype structure and linking correlated alleles of possible functional
781       variants. *Bioinformatics* **31**, 3555-3557 (2015).
782

783  45.  Mercader JM*, et al.* Altered brain-derived neurotrophic factor blood levels and gene
784       variability are associated with anorexia and bulimia. *Genes Brain Behav* **6**, 706-716
785       (2007).
786

787  46.  Millard LAC, Davies NM, Gaunt TR, Davey Smith G, Tilling K. Software Application
788       Profile: PHESANT: a tool for performing automated phenome scans in UK Biobank.
789       *Int J Epidemiol* **47**, 29-35 (2018).
790

791  47.  Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of
792       genomewide association scans. *Bioinformatics* **26**, 2190-2191 (2010).
793

794  48.  Zhou W*, et al.* Efficiently controlling for case-control imbalance and sample
795       relatedness in large-scale genetic association studies. *Nat Genet* **50**, 1335-1341
796       (2018).
797

798  49.  Lumley T. rmeta: Meta-Analysis.). R package version 3.0 edn (2018).
799

800  50.  Wellcome Trust Case Control C*, et al.* Bayesian refinement of association signals for
801       14 loci in 3 common diseases. *Nat Genet* **44**, 1294-1301 (2012).
802

803  51.  Giambartolomei C*, et al.* A Bayesian framework for multiple trait colocalization from
804       summary association statistics. *Bioinformatics* **34**, 2538-2545 (2018).
805

806  52.  Wallace C. Statistical testing of shared genetic control for potentially related traits.
807       *Genet Epidemiol* **37**, 802-813 (2013).
808

809  53.  Wallace C*, et al.* Statistical colocalization of monocyte gene expression and genetic
810       risk variants for type 1 diabetes. *Hum Mol Genet* **21**, 2815-2824 (2012).
811

812  54.  Lim ET*, et al.* A novel test for recessive contributions to complex diseases implicates
813       Bardet-Biedl syndrome gene BBS10 in idiopathic type 2 diabetes and obesity. *Am J*
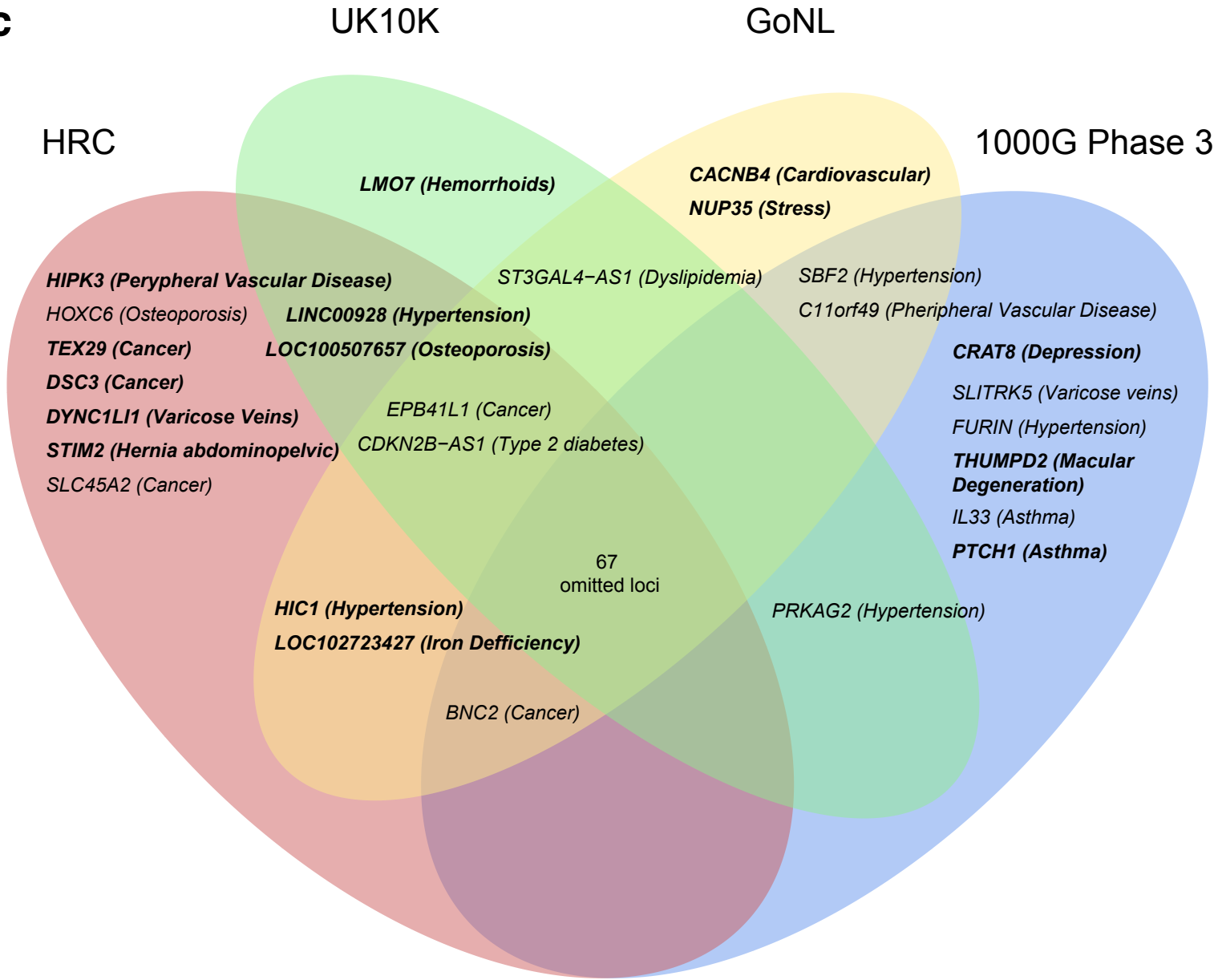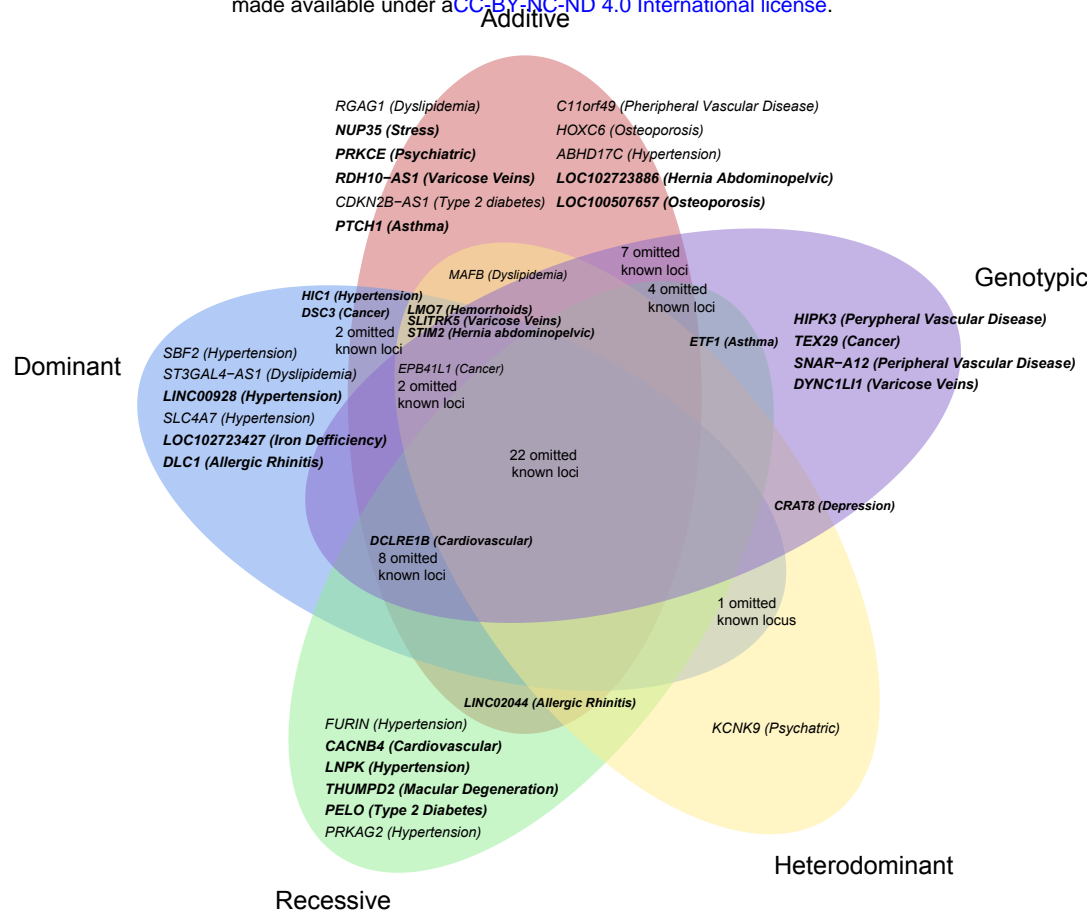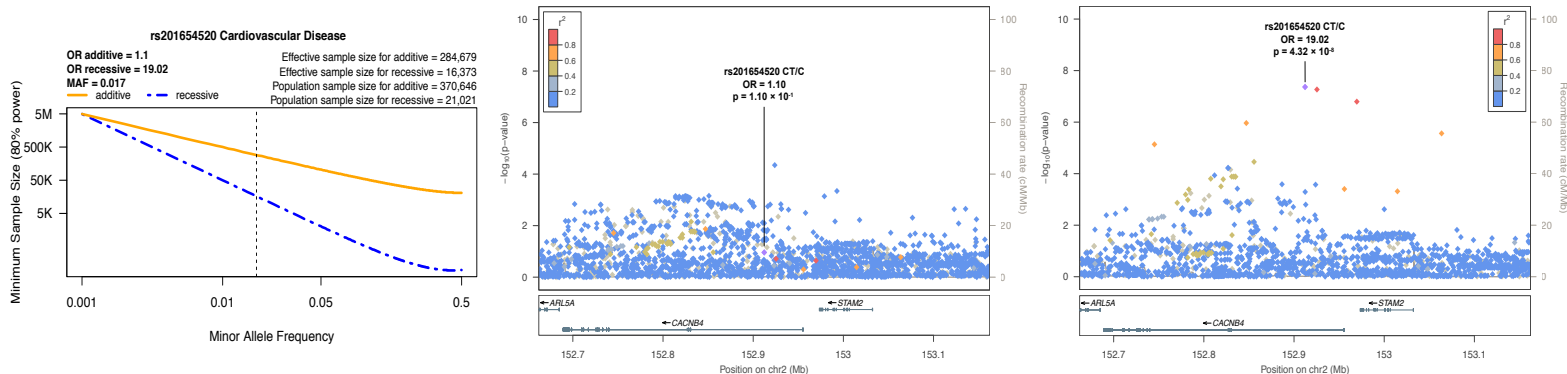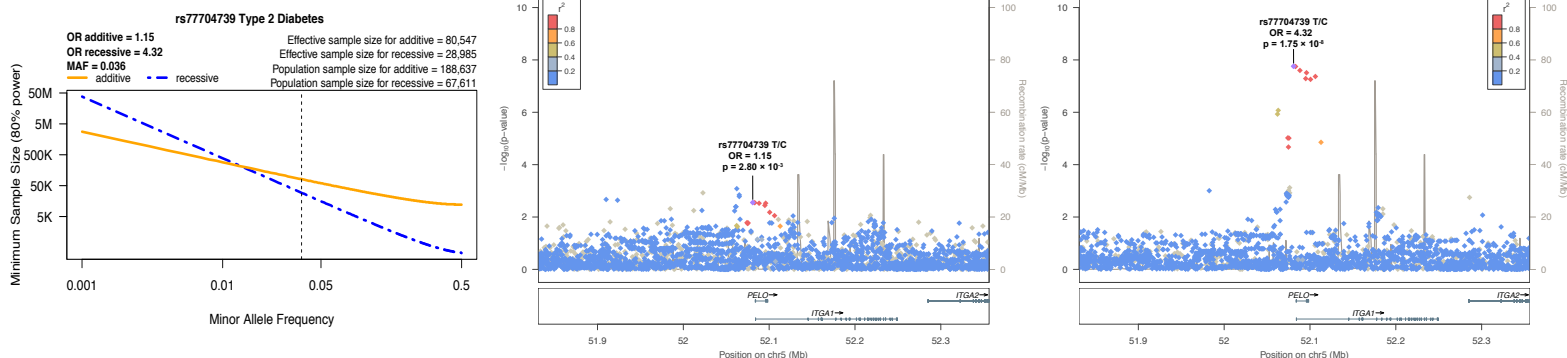814       *Hum Genet* **95**, 509-520 (2014).
815
816

**a**



**b**



**c**



**d**

**Figure legends.**

**Figure 1. Graphical representation illustrating the benefits of combining the results from different reference panels. a** Comparison of the number of variants after the imputati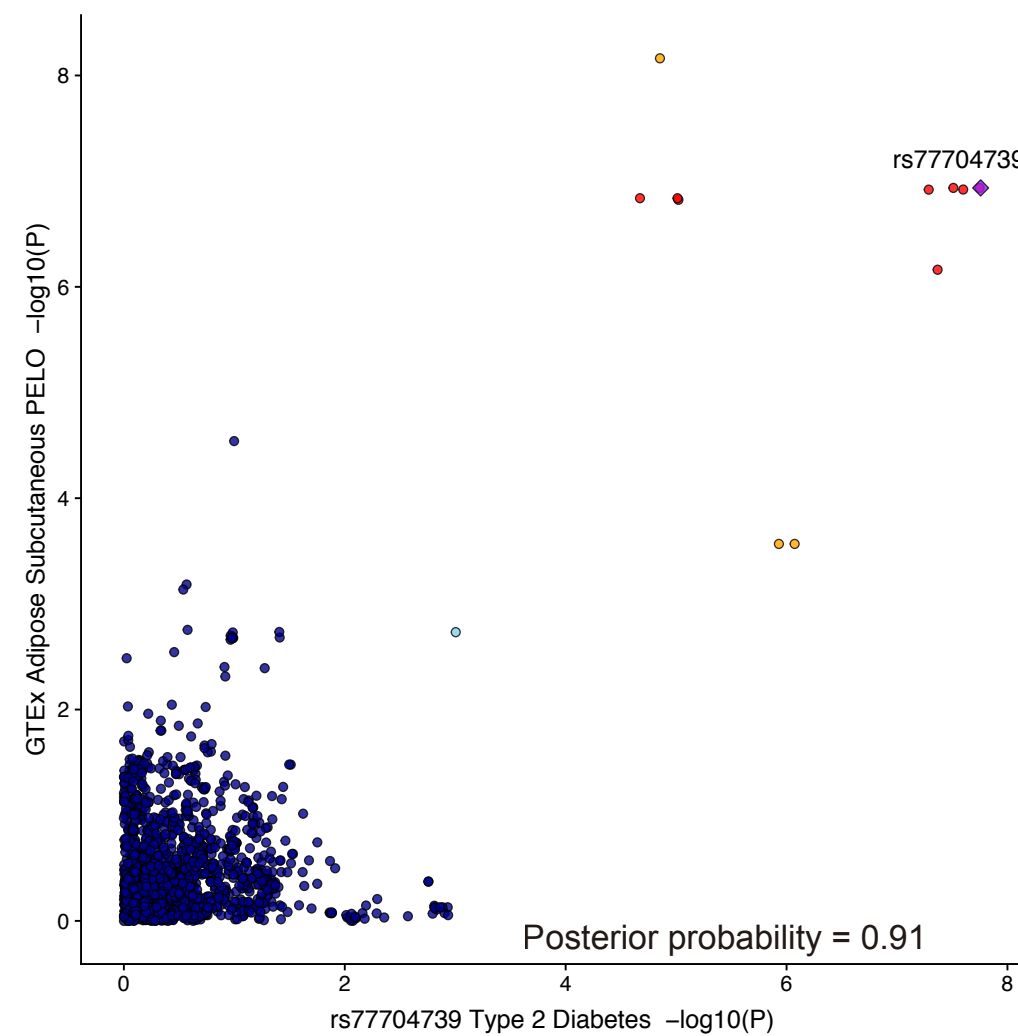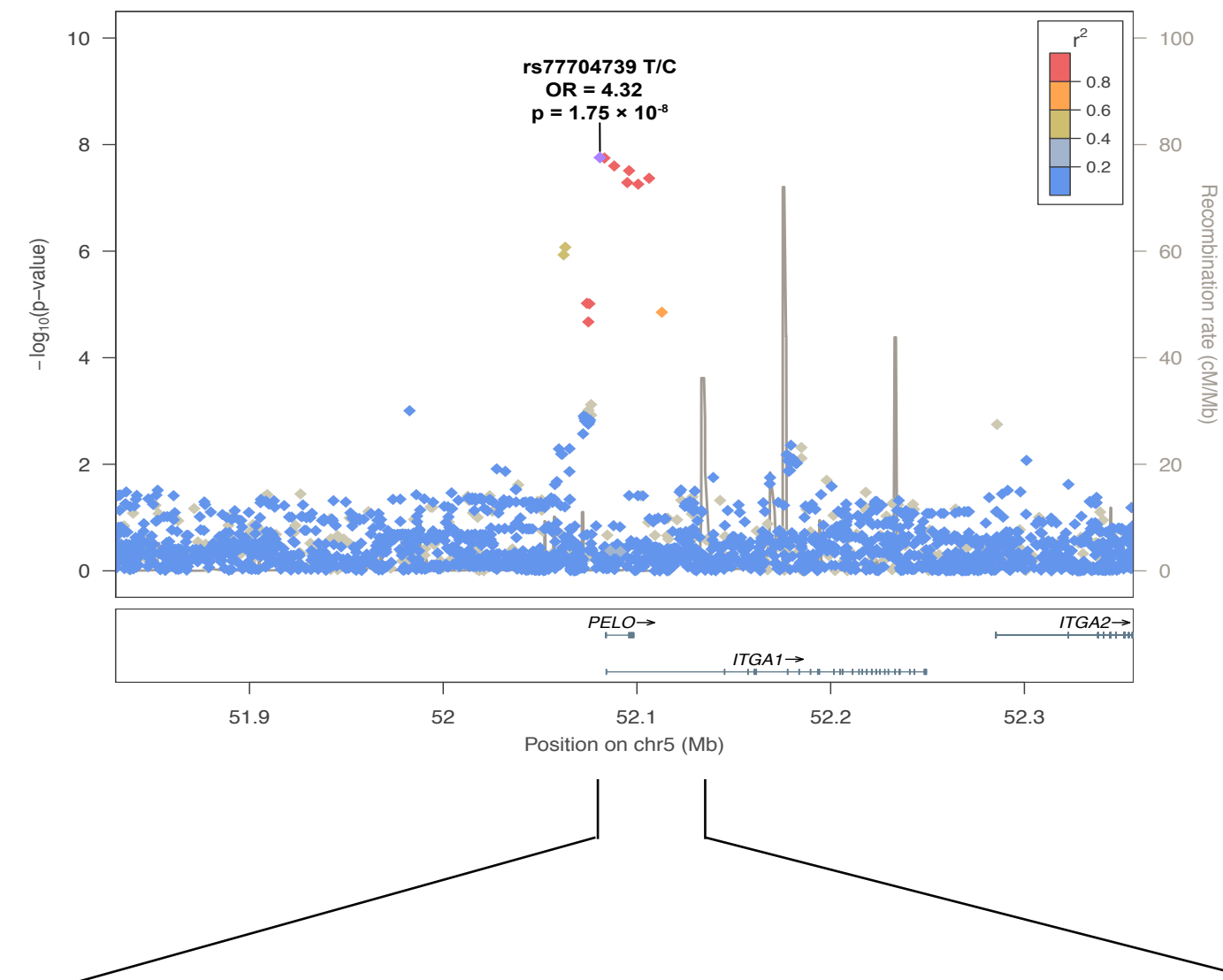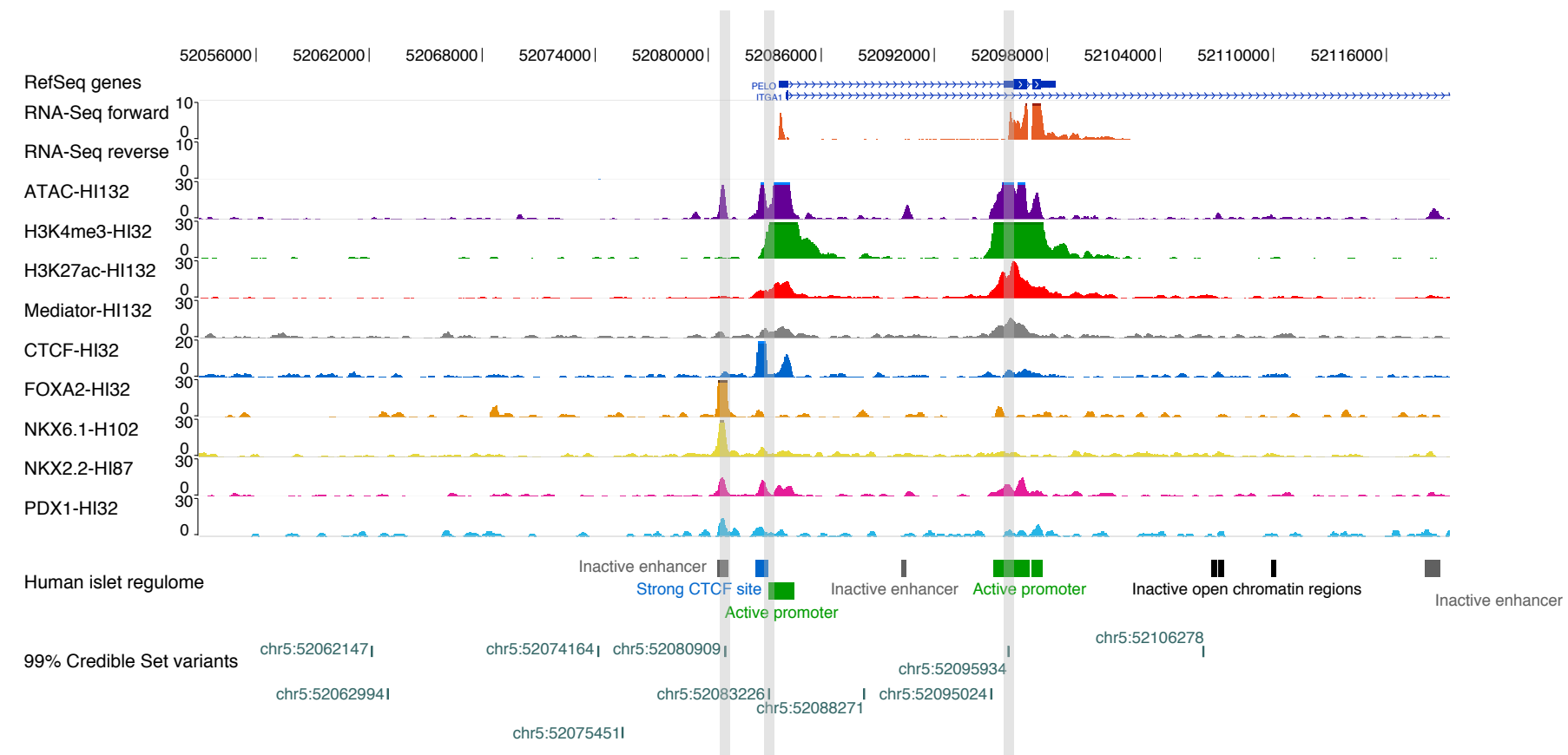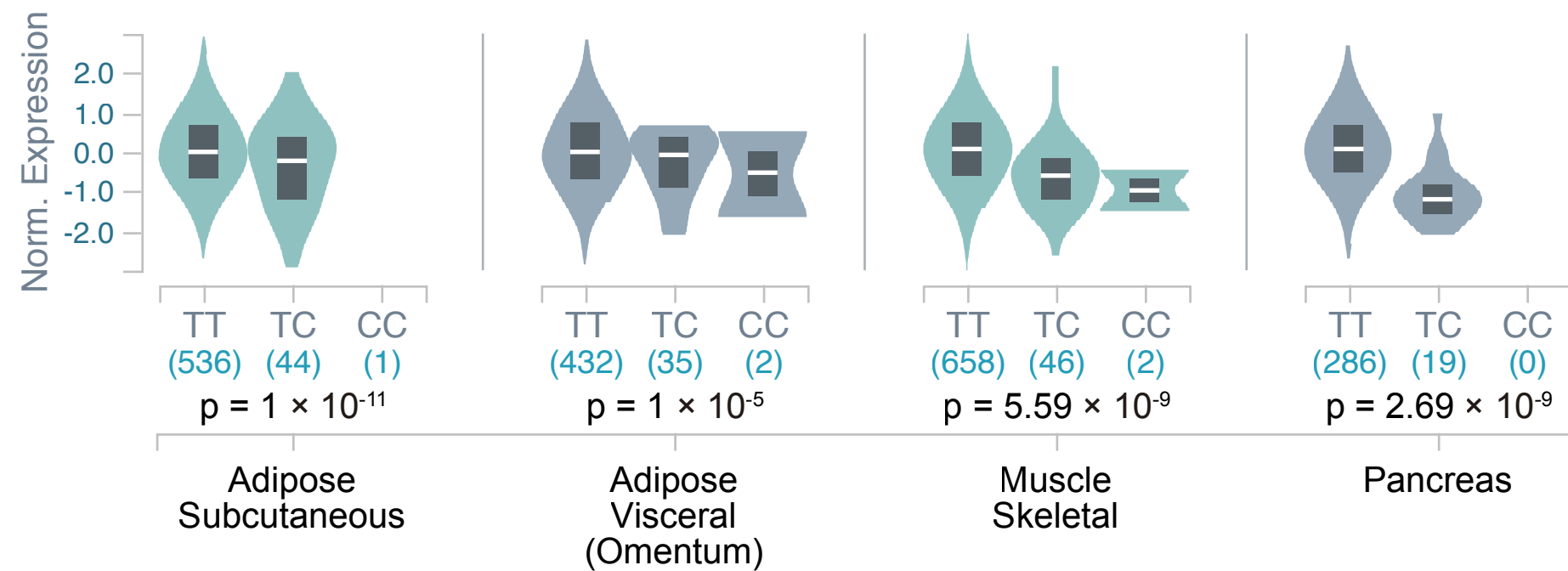on with four reference panels (info score ≥ 0.7), and combining them, colored according to MAF and varianat type (SNP vs alternative forms of variation, such as indels). As shown in the bar plot, combining the results from the four reference panels increased the final set of variants for association testing when compared with the results for each of the panels alone (GoNL, UK10K, 1000G Phase 3 or HRC), especially in the low and rare frequency spectrum. For example, we covered up to 5.5 M rare variants (0.01 > MAF > 0.001) by combining panels, while only 2,3 M, 2,9 M, 3,2 M and 3,8 M of rare variants were imputed independently with GoNL, UK10K, 1000G phase 3 and HRC, respectively. **b** Comparison of the contribution of each reference panel in the combined results. Each bar represents the number of variants that had the best imputation accuracy for a given reference panel. As shown in the figure, although the HRC panel showed overall higher imputation scores, as it provided around 10 of the final 16 M variants, the contribution of the other reference panels, primarily with non-SNP variants, was substantial. Indels seen in the bar plot for HRC correspond to genotyped indels. All variants with info score < 0.7, MAF < 0.001 and HWE for controls $p < 1.0 \times 10^{-6}$ were filtered. **c** Venn Diagram illustrating the loci that identified by each reference panel. New loci are depicted in bold. As shown in this figure, only 67 of the 94 GWAS significant loci were identified by all four reference panels, while 27 of them (28.7%) were only identified by one, two or three of the four panels.

**Figure 2. Results from the analysis of additive and non-additive inheritance models. a** The Venn Diagram shows the number of loci that were identified when analyzing multiple inheritance models. As seen in the Venn Diagram, the strongest association for 37 of the 94 associated loci was non-additive. Moreover, the analysis of non-additive models was crucial for the identification of 14 novel (in bold) associated loci. **b** Power calculation of the rs201654520 indel in *CACNB4* associated with cardiovascular disease. The results show that the additive-based test would require a population sample size of 370,646 individuals to find this recessive association, while the population sample size needed for the recessive model was 21,021. **c** Power calculation of the rs77704739 variant near the *PELO* gene associated with type 2 diabetes. The results show that the additive-based test would require a population sample size of 188,637 individuals to find this recessive association, while the population sample size needed for the recessive model is 67,611. **d** Power calculation of the rs557998486 indel near

the *THUMPD2* gene associated with age-related macular degeneration. The results show that the additive-based test would require a population sample size of 6,493,419 individuals to find this recessive association, while the population sample size for the recessive model is 157,450.

**Figure 3. Functional characterization of the rs77704739 recessive association near the *PELO* gene. a** Signal plot for chromosome 5 region surrounding rs77704739. Each point represents a variant, with its *p*-value from the discovery stage on a −log10 scale in the y axis. The x axis represents the genomic position (hg19). Three credible set variants are located in open chromatin sites in human pancreatic islets, one of them classified as an active promoter and one highly bounded by pancreatic islet specific transcription factors, such as PDX1, NKX2.2, NKX6.1 and FOXA2. **b** Colocalization plots from LocusCompare for the rs77704739 variant in adipose subcutaneous tissue. As seen in the plots, the signals from both eQTL data and the recessive T2D association results colocalize. c Violin plot from GTEx showing that the recessive rs77704739 variant significantly modifies the expression of *PELO* gene in subcutaneous and visceral adipose tissue, skeletal muscle and pancreas. GTEx V7 was used for colocalization analyses, whereas GTEx V8 was used to generate the violin plots.

# Table 1. New associations from the GERA cohort analysis

| Phenotype (Cases/Controls) | CHR | Nearest Gene | Position | rsID | Alleles | MAF | Lowest P-value Model | Additive Model OR (CI 95%) | Additive Model P-value | Lowest P-value Model OR (CI 95%) | Lowest P-value Model P-value | Dominance Deviation P-value |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Allergic Rhinitis (13,936/42,701) | 3 | *LINC02044* | 112,911,615 | rs2399472 | C/T | 0.073 | Additive | 1.17 (1.10-1.23) | $1.55 \times 10^{-8}$ | 1.17 (1.10-1.23) | $1.55 \times 10^{-8}$ | $6.66 \times 10^{-1}$ |
| | 8 | *DLC1* | 13,164,746 | rs10112506 | A/G | 0.390 | Dominant | 0.94 (0.91-0.97) | $8.61 \times 10^{-5}$ | 0.89 (0.86-0.93) | $1.54 \times 10^{-8}$ | $2.86 \times 10^{-4}$ |
| Asthma (9,209/47,428) | 5 | *ETF1* | 137,858,067 | rs154073 | C/T | 0.429 | Recessive | 1.09 (1.06-1.13) | $6.06 \times 10^{-5}$ | 1.18 (1.12-1.25) | $4.23 \times 10^{-9}$ | $9.28 \times 10^{-3}$ |
| | 9 | *PTCH1* | 98,344,866 | rs67053006 | C/G | 0.139 | Additive | 0.87 (0.83-0.91) | $4.14 \times 10^{-8}$ | 0.87 (0.83-0.91) | $4.14 \times 10^{-8}$ | $8.10 \times 10^{-1}$ |
| Cancer (17,131/39,506) | 13 | *TEX29* | 112,115,591 | rs138646839 | C/T | 0.005 | Genotypic | 1.68 (1.39-2.03) | $1.45 \times 10^{-7}$ | 1.60 (1.32-1.96) / > 10 (1.01->10)* | $3.54 \times 10^{-9}$ | - |
| | 18 | *DSC3* | 28,442,343 | rs2014497 | A/G | 0.008 | Additive | 1.50 (1.30-1.72) | $2.44 \times 10^{-8}$ | 1.50 (1.30-1.72) | $2.44 \times 10^{-8}$ | $6.00 \times 10^{-1}$ |
| Cardiovascular (15,009/41,628) | 1 | *DCLRE1B* | 114,448,752 | rs10858023 | C/T | 0.350 | Dominant | 1.09 (1.06-1.12) | $3.26 \times 10^{-8}$ | 1.14 (1.09-1.19) | $2.11 \times 10^{-9}$ | $1.94 \times 10^{-2}$ |
| | 2 | *CACNB4* | 152,912,244 | rs201654520 | CT/C | 0.017 | Recessive | 1.10 (0.98-1.22) | $1.10 \times 10^{-1}$ | 19.02 (5.50-65.84) | $4.32 \times 10^{-9}$ | $4.36 \times 10^{-6}$ |
| Major Depression Disorder (7,264/49,373) | 12 | *CRAT8* | 128,551,715 | rs1455286248 | GT/G | 0.281 | Heterodominant | 0.94 (0.90-0.98) | $3.00 \times 10^{-3}$ | 1.18 (1.12-1.25) | $3.15 \times 10^{-9}$ | $1.10 \times 10^{-6}$ |
| Type 2 Diabetes (6,967/49,670) | 5 | *PELO* | 52,080,909 | rs77704739 | T/C | 0.036 | Recessive | 1.15 (1.05-1.26) | $2.80 \times 10^{-3}$ | 4.32 (2.70-6.92) | $1.75 \times 10^{-8}$ | $1.92 \times 10^{-7}$ |
| Hemorrhoids (9,129/47,508) | 13 | *LMO7* | 76,281,808 | rs186102686 | C/T | 0.004 | Heterodominant | 1.98 (1.58-2.48) | $2.18 \times 10^{-8}$ | 1.99 (1.59-2.49) | $2.03 \times 10^{-8}$ | - |
| Hernia Abdominopelvic (6,291/50,346) | 1 | *LOC102723886* | 219,762,581 | rs2494196 | C/A | 0.274 | Additive | 1.13 (1.08-1.18) | $2.03 \times 10^{-8}$ | 1.13 (1.08-1.18) | $2.03 \times 10^{-8}$ | $6.87 \times 10^{-1}$ |
| | 4 | *STIM2* | 27,019,359 | rs113180595 | T/C | 0.004 | Heterodominant | 2.17 (1.69-2.78) | $1.59 \times 10^{-8}$ | 2.18 (1.70-2.8) | $1.27 \times 10^{-8}$ | - |
| Hypertension Disease (28,391/28,246) | 2 | *LNPK* | 176,532,019 | rs1446802 | A/G | 0.500 | Recessive | 1.07 (1.04-1.09) | $1.66 \times 10^{-6}$ | 1.13 (1.08-1.17) | $4.42 \times 10^{-9}$ | $6.85 \times 10^{-3}$ |
| | 15 | *LINC00928* | 90,081,905 | rs28792763 | G/A | 0.462 | Dominant | 0.94 (0.91-0.96) | $4.14 \times 10^{-6}$ | 0.88 (0.84-0.92) | $4.42 \times 10^{-8}$ | $4.80 \times 10^{-3}$ |
| | 17 | *HIC1* | 1,959,826 | rs112963849 | C/A | 0.082 | Additive | 1.15 (1.10-1.21) | $1.71 \times 10^{-8}$ | 1.15 (1.10-1.21) | $1.71 \times 10^{-8}$ | $8.01 \times 10^{-1}$ |
| Iron Deficiency Anemia (2,439/54,198) | 7 | *LOC102723427* | 67,292,424 | rs79798837 | C/T | 0.118 | Dominant | 0.77 (0.70-0.85) | $1.69 \times 10^{-7}$ | 0.74 (0.66-0.83) | $3.80 \times 10^{-8}$ | $8.92 \times 10^{-2}$ |
| Macular Degeneration (3,685/52,952) | 2 | *THUMPD2* | 40,010,523 | rs557998486 | T/TG | 0.009 | Recessive | 1.07 (0.81-1.41) | $6.28 \times 10^{-1}$ | 10.5** | $2.75 \times 10^{-8}$ | - |
| Osteoporosis (5,399/51,238) | 22 | *LOC100507657* | 27,772,054 | rs139959245 | C/T | 0.007 | Additive | 1.91 (1.53-2.37) | $4.79 \times 10^{-8}$ | 1.91 (1.53-2.37) | $4.79 \times 10^{-8}$ | - |
| Psychiatric (8,624/48,013) | 2 | *PRKCE* | 46,278,720 | rs12712961 | T/A | 0.452 | Additive | 1.10 (1.06-1.14) | $1.66 \times 10^{-8}$ | 1.10 (1.06-1.14) | $1.66 \times 10^{-8}$ | $2.57 \times 10^{-1}$ |
| Peripheral Vascular Disease (4,301/52,336) | 11 | *HIPK3* | 33,391,655 | rs80274406 | A/G | 0.091 | Genotypic | 1.06 (0.98-1.15) | $1.76 \times 10^{-1}$ | 1.17 (1.07-1.27) / 0.26 (0.13-0.53)* | $4.26 \times 10^{-9}$ | $6.32 \times 10^{-6}$ |
| | 19 | *SNAR-A12* | 48,403,215 | rs2932761 | A/G | 0.289 | Genotypic | 0.97 (0.93-1.02) | $3.04 \times 10^{-1}$ | 1.11 (1.03-1.18 / 0.76 (0.66-0.87)* | $3.55 \times 10^{-8}$ | $1.35 \times 10^{-8}$ |
| Acute reaction to Stress (4,314/52,323) | 2 | *NUP35* | 184,407,101 | rs577242570 | T/G | 0.004 | Additive | 2.33 (1.77-3.08) | $4.56 \times 10^{-8}$ | 2.33 (1.77-3.08) | $4.56 \times 10^{-8}$ | - |
| Varicose Veins (2,483/54,154) | 3 | *DYNC1LI1* | 32,652,184 | rs62250779 | G/A | 0.073 | Genotypic | 1.17 (1.05-1.3) | $5.60 \times 10^{-3}$ | 1.29 (1.16-1.45) / 0.13 (0.03-0.60)* | $2.13 \times 10^{-9}$ | $9.58 \times 10^{-4}$ |
| | 8 | *RDH10-AS1* | 74,284,818 | rs2383896 | A/G | 0.479 | Additive | 1.17 (1.11-1.24) | $5.00 \times 10^{-8}$ | 1.17 (1.11-1.24) | $5.00 \times 10^{-8}$ | $9.88 \times 10^{-1}$ |
| | 13 | *SLITRK5* | 88,346,617 | rs117798068 | T/C | 0.011 | Heterodominant | 2.03 (1.63-2.53) | $1.59 \times 10^{-8}$ | 2.07 (1.66-2.59) | $8.41 \times 10^{-9}$ | - |

CHR = Chromosome, Position = Position hg19, Alleles = Non-effect Allele / Effect Allele, MAF=Minor Allele Frequency, OR= Odds Ratio, CI= Confidence Interval

* Odds Ratio and confidence interval for heterozygous / Odds Ratio and confidence interval for effect allele homozygous calculated using the method het+hom from SNPTEST

** Odds Ratio calculated using the Recessive Allele Frequency-Based Test (RAFT) [53]

**Table 2. Replication of new associations with UK Biobank**

| CHR | rsID (Alleles) (MAF) | Best Model | Phenotype (Cases/Controls) | Stage 1. Discovery Additive OR (CI 95%) | P-value | Best Model OR (CI 95%) | P-value | Field (Cases/Controls or Sample Size) | Stage 2. Replication Additive OR (CI 95%) | P-value | Lowest p-value model OR (CI 95%) | P-value | Stage 1 + Stage 2. Meta-analysis Additive OR (CI 95%) | P-value | Lowest p-value model OR (CI 95%) | P-value |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 18 | rs2014497 (A/G) (0.008) | Additive | Cancer (17,131/39,506) | 1.50 (1.30-1.72) | $2.44×10^{-8}$ | 1.50 (1.30-1.72) | $2.44×10^{-8}$ | Self-reported: chronic lymphocytic (237/360,904) | 2.13 (1.14-3.97) | $3.50×10^{-2}$ | 2.13 (1.14-3.97) | $3.50×10^{-2}$ | 1.52 (1.33-1.74) | $1.60×10^{-9}$ | 1.52 (1.33-1.74) | $1.60×10^{-9}$ |
| | | | | | | | | Self-reported: kidney/renal cell cancer (473/360,668) | 1.75 (1.07-2.86) | $4.25×10^{-2}$ | 1.75 (1.07-2.86) | $4.25×10^{-2}$ | 1.51 (1.32-1.73) | $1.49×10^{-9}$ | 1.51 (1.32-1.73) | $1.49×10^{-9}$ |
| | | | | | | | | C69 Malignant neoplasm of eye and adnexa (146/361,048) | 2.51 (1.19-5.3) | $3.56×10^{-2}$ | 2.51 (1.19-5.3) | $3.56×10^{-2}$ | 1.52 (1.33-1.75) | $1.95×10^{-9}$ | 1.52 (1.33-1.75) | $1.95×10^{-9}$ |
| 1 | rs2494196 (C/A) (0.274) | Additive | Hernia Abdominopelvic (6,291/50,346) | 1.13 (1.08-1.18) | $2.03×10^{-8}$ | 1.13 (1.08-1.18) | $2.03×10^{-8}$ | Self-reported: umbilical hernia (328/360,813) | 1.42 (1.21-1.67) | $2.31×10^{-5}$ | 1.42 (1.21-1.67) | $2.31×10^{-5}$ | 1.15 (1.10-1.19) | $5.35×10^{-11}$ | 1.15 (1.10-1.19) | $5.35×10^{-11}$ |
| | | | | | | | | K40 Inguinal hernia (13,365/347,829) | 1.09 (1.06-1.12) | $3.95×10^{-10}$ | 1.09 (1.06-1.12) | $3.95×10^{-10}$ | 1.10 (1.08-1.12) | $7.78×10^{-17}$ | 1.10 (1.08-1.12) | $7.78×10^{-17}$ |
| | | | | | | | | K41 Femoral hernia (475/360,719) | 1.44 (1.26-1.64) | $1.24×10^{-7}$ | 1.44 (1.26-1.64) | $1.24×10^{-7}$ | 1.16 (1.11-1.21) | $2.26×10^{-12}$ | 1.16 (1.11-1.21) | $2.26×10^{-12}$ |
| | | | | | | | | K42 Umbilical hernia (2,623/358,571) | 1.29 (1.22-1.37) | $1.14×10^{-17}$ | 1.29 (1.22-1.37) | $1.14×10^{-17}$ | 1.19 (1.15-1.22) | $2.94×10^{-22}$ | 1.19 (1.15-1.22) | $2.94×10^{-22}$ |
| | | | | | | | | K43 Ventral hernia (2,470/358,724) | 1.18 (1.11-1.25) | $1.77×10^{-7}$ | 1.18 (1.11-1.25) | $1.77×10^{-7}$ | 1.15 (1.11-1.19) | $1.99×10^{-14}$ | 1.15 (1.11-1.19) | $1.99×10^{-14}$ |
| 2 | rs557998486 (T/TG) (0.009) | Recessive | Macular Degeneration (3,685/52,952) | 1.07 (0.81-1.41) | $6.28×10^{-1}$ | 10.5* | $2.75×10^{-8}$ | Eye problems/disorders: Macular degeneration (2,726/115,164) | 0.98 (0.72-1.32) | $8.81×10^{-1}$ | 7.58 (1.54-37.32) | $4.1×10^{-2}$ | 1.01(0.82-1.24)** | $7.91×10^{-1}$*** | 26.51(7.57-92.85)** | $3.29×10^{-8}$*** |
| 5 | rs77704739 (T/C) (0.036) | Recessive | Type 2 Diabetes (6,967/49,670) | 1.15 (1.05-1.26) | $2.80×10^{-3}$ | 4.32 (2.70-6.92) | $1.75×10^{-8}$ | Self-reported: diabetes (14,114/347,027) | 1.03 (0.97-1.09) | $3.87×10^{-1}$ | 1.88 (1.35-2.6) | $4.95×10^{-4}$ | 1.06 (1.01-1.12) | $1.78×10^{-2}$ | 2.46 (1.88-3.21) | $4.68×10^{-11}$ |

CHR = Chromosome, Position = Position hg19, Alleles = Non-effect Allele / Effect Allele, MAF= Minor Allele Frequency, OR= Odds Ratio

* Odds Ratio calculated using the Recessive Allele Frequency-Based Test (RAFT)

** Obtained through a mega-analysis with UK Biobank using the "expected" method from SNPTEST

*** Obtained using METAL method "SAMPLESIZE" to combine the p-values taking into account the sample size and direction of effect.