# ARTICLE

Check for updates

# The impact of non-additive genetic associations on age-related complex diseases

Marta Guindo-Martínez [1,18], Ramon Amela[1,18], Silvia Bonàs-Guarch [1,2,3], Montserrat Puiggròs [1], Cecilia Salvoro [1], Irene Miguel-Escalada [1,2,3], Caitlin E. Carey[4,5], Joanne B. Cole [6,7,8,9], Sina Rüeger[10], Elizabeth Atkinson [4,5,11], Aaron Leong[8,12], Friman Sanchez[1], Cristian Ramon-Cortes [1], Jorge Ejarque [1], Duncan S. Palmer[4,5,17], Mitja Kurki[10], FinnGen Consortium*, Krishna Aragam[11,13,14], Jose C. Florez[6,7,15], Rosa M. Badia [1], Josep M. Mercader [1,6,7,15,19] ✉ & David Torrents[1,16,19] ✉

Genome-wide association studies (GWAS) are not fully comprehensive, as current strategies typically test only the additive model, exclude the X chromosome, and use only one reference panel for genotype imputation. We implement an extensive GWAS strategy, GUIDANCE, which improves genotype imputation by using multiple reference panels and includes the analysis of the X chromosome and non-additive models to test for association. We apply this methodology to 62,281 subjects across 22 age-related diseases and identify 94 genome-wide associated loci, including 26 previously unreported. Moreover, we observe that 27.7% of the 94 loci are missed if we use standard imputation strategies with a single reference panel, such as HRC, and only test the additive model. Among the new findings, we identify three novel low-frequency recessive variants with odds ratios larger than 4, which need at least a three-fold larger sample size to be detected under the additive model. This study highlights the benefits of applying innovative strategies to better uncover the genetic architecture of complex diseases.

[1] Barcelona Supercomputing Center (BSC), Barcelona, Spain. [2] Regulatory Genomics and Diabetes, Centre for Genomic Regulation, The Barcelona Institute of Science and Technology, Barcelona, Spain. [3] CIBER de Diabetes y Enfermedades Metabólicas Asociadas, Madrid, Spain. [4] Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA, USA. [5] Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General Hospital, Boston, MA, USA. [6] Programs in Metabolism and Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA. [7] Diabetes Unit and Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, USA. [8] Harvard Medical School, Boston, MA, USA. [9] Division of Endocrinology and Center for Basic and Translational Obesity Research, Boston Children's Hospital, Boston, MA, USA. [10] Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland. [11] Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA. [12] Department of Medicine, Massachusetts General Hospital, Boston, MA, USA. [13] Cardiology Division, Massachusetts General Hospital, Boston, MA, USA. [14] Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA, USA. [15] Department of Medicine, Harvard Medical School, Boston, MA, USA. [16] Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain. [17] Present address: GENOMICS plc, Oxford, UK. [18] These authors contributed equally: Marta Guindo-Martínez, Ramon Amela. [19] These authors jointly supervised this work: Josep M. Mercader, David Torrents. *A full list of members and their affiliations appears in the Supplementary Information. ✉email: mercader@broadinstitute.org; david.torrents@bsc.es

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 they may contribute to the pathophysiology of complex diseases. Some of the current limitations of GWAS could be overcome by increasing sample sizes and by applying more comprehensive analytical methods with improved imputation strategies[3]. Besides, whole-exome sequencing (WES) and whole-genome sequencing (WGS) datasets, such as the NHLBI Trans-Omics for Precision Medicine (TOPMed) Project[4], and the UK Biobank[5], are rapidly growing and expanding the range of variants to be tested in genetic association studies. However, together with the general tendency of increasing sample sizes, this imposes additional methodological and computational challenges. These can require scientists to restrict and simplify the analysis by limiting it to autosomal chromosomes, to a single reference panel for imputation, and to a single (additive) inheritance model for association testing, leaving a relevant fraction of the genetic architecture of the disease unexplored[6].

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[7] 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[8,9]. The importance of evaluating non-additive inheritance models is well reported in the context of Mendelian diseases[10] and occasionally for complex traits as well, such as the recessive effects of the *FTO* locus in obesity[11], and in or near *ITGA1*[12], *TBC1D4*[13], and *CDKAL1*[11,14] genes in type 2 diabetes, as well as the known non-additive effects of HLA haplotypes in autoimmune diseases[15] and ulcerative colitis[16]. 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 association strategies that should include, among other improvements, the analysis of non-additive inheritance models.

In this work, and to fill this gap and to determine the prevalence and contribution of the different inheritance patterns involved in the genetic architecture of complex diseases, we design and implement a comprehensive strategy for genetic association analysis that combines imputation from multiple reference panels with association testing under five different inheritance models across multiple phenotypes. We apply this strategy to the Kaiser Permanente Research Program on Genes, Environment and Health: A Genetic Epidemiology Research on Adult Health and Aging (GERA) cohort[17], which includes 62,281 subjects from European ancestry and 22 diseases. Finally, we release here both the summary statistics for all the models of inheritance as well as the complete methodology, provided to the community as an easy-to-use and standalone pipeline. This pipeline allows the analysis of existing and newly generated GWAS data with better efficiency and more comprehensive testing, improving the chances of variant discovery.

## Results

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

**Genotype imputation and association testing using multiple reference panels.** After applying genetic quality control to the GERA cohort (see "Methods"), we retained 56,637 individuals with European ancestry for further downstream analysis (Supplementary Data 1). To cover the maximum number and type of genetic variants, we next applied an imputation strategy with four reference panels: the Genome of the Netherlands (GoNL)[18,19], the UK10K Project[20], the 1000 Genomes Project (1000G) phase 3[21], and Haplotype Reference Consortium (HRC)[22], and imputed 11.2 M, 11.4 M, 13.1 M, and 11.7 M high-quality imputed variants (IMPUTE2[23] info score ≥0.7 and minor allele frequency [MAF] ≥ 0.001) with each panel, respectively. After combining the results of the four reference panels by choosing, for each variant, the panel that provided the highest imputation accuracy, we retained a total of 16,059,686 variants covering all the autosomes and the X chromosome (Fig. 1a). Using this strategy we imputed 2.6 M and 5.5 M high-quality, low-frequency (0.05> MAF > 0.01) and rare variants (0.01 > MAF > 0.001), respectively, as well as 1.6 M indels. Note that as many as 684,393 common variants (MAF ≥ 0.05), 255,106 low-frequency, 1.7 M rare, and all indels (1.6 M) would be missed if only the HRC reference panel was used. This highlights the benefit of combining different reference panels for comprehensive association testing (Fig. 1b).

To evaluate our imputation strategy, we used sequenced data from UK10K, which includes 3781 sequenced genomes, to build an in silico array of 599,208 variants. We then imputed the in silico array using 1000G, GoNL and HRC as reference panels and compared the imputed genotype dosages with the sequenced genotype dosages (allelic dosage $R^2$, see "Methods" and Supplementary Fig. 1). Our results empirically demonstrate that our IMPUTE2-info threshold of 0.7 is a good cutoff for well-predicted genotypes (Supplementary Fig. 2). Moreover, the combination of the results from the different panels based on IMPUTE2-info values outperforms single reference panels in terms of % of variants with allelic dosage $R^2 \geq 0.5$ in autosomes (Fig. 1c) and in the X chromosome (Fig. 1d) for both SNPs and indels (Supplementary Fig. 3).

To evaluate that our strategy is able to identify previously known loci that are well-powered in our cohort, we performed a power comparison under the additive model for type 2 diabetes and age-related macular degeneration, since both are well-defined diseases with available summary statistics. The power comparison between our additive results for type 2 diabetes and those from the DIAMANTE consortium[24] showed that while we had an 80% power to find two variants with a *p*-value of $5.0 \times 10^{-8}$ in
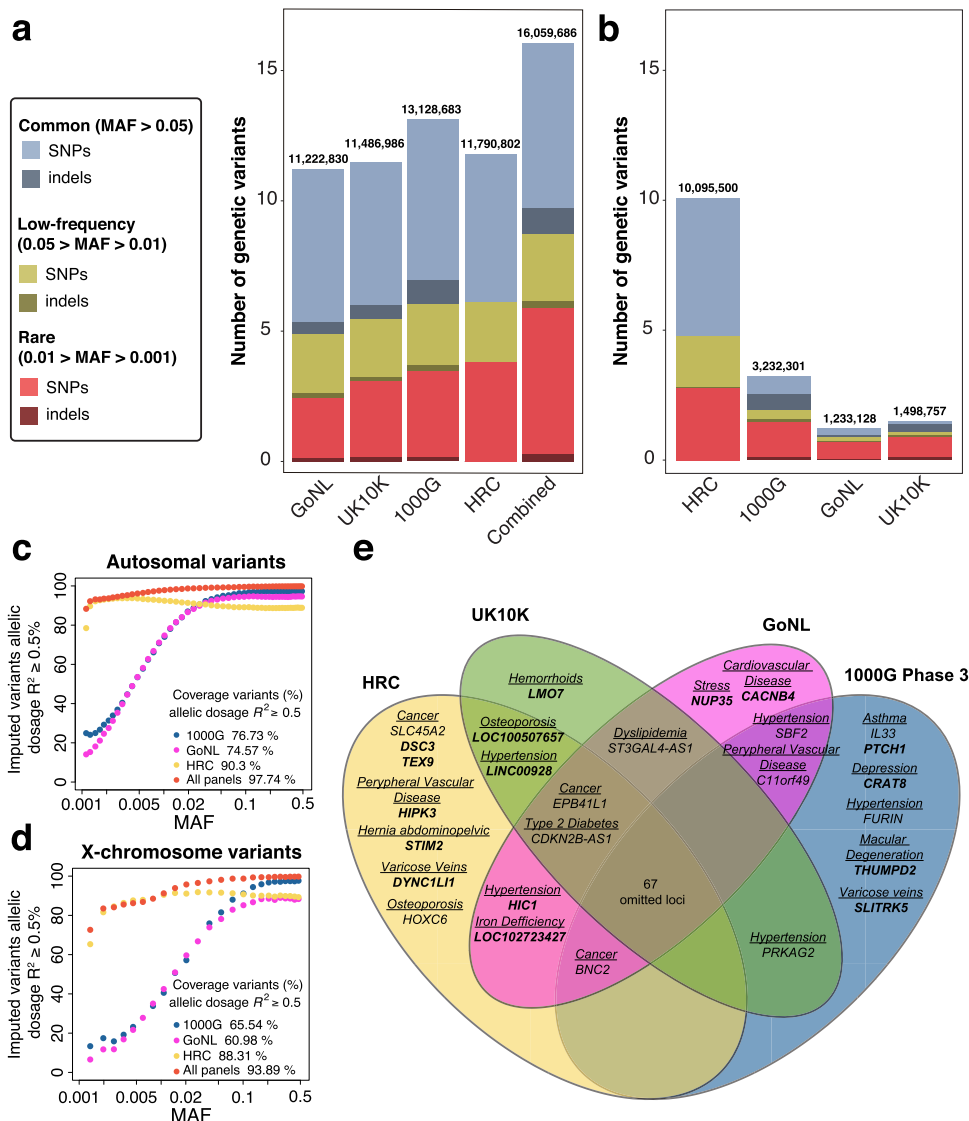
**Fig. 1 Graphical representation illustrating the benefits of combining the results from different reference panels. a** Comparison of the number of variants after the imputation with four reference panels (info score ≥ 0.7), and combining them, colored according to MAF and variant type (SNP vs alternative 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** Percentage of high-quality imputed variants (IMPUTE2-info score ≥ 0.7) with an allelic dosage $R^2 \geq 0.5$ between sequenced genotypes in UK10K samples vs variants imputed in the same UK10K samples using 1000G phase 3, GoNL, and HRC reference panels for the autosomes. The percentage of high-quality imputed variants with allelic dosage $R^2$ values (y axis) are represented across several MAF ranges (x-axis) for each of the reference panels and the combined panels imputed results. The combination of the three reference panels outperforms the single reference panels with 97.74% of variants with $R^2 \geq 0.5$. **d** Percentage of variants in the X chromosome with an IMPUTE2-info score ≥ 0.7 and with an allelic dosage $R^2 \geq 0.5$ for UK10K imputed genotypes across MAF ranges for 1000G phase 3, GoNL, and HRC reference panels and the combined results. The combination of the results from the three panels outperforms single reference panels with 93.89% of variants with allelic dosage $R^2 \geq 0.5$. **e** Venn Diagram illustrating the loci 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.

our dataset, we found four variants with genome-wide significance (Supplementary Data 2). For age-related macular degeneration[25], we had 80% power to find six variants with a $p$-value of $5.0 \times 10^{-8}$ of which we identified four of them with genome-wide significance. The remaining two variants were nominally significant in our study (Supplementary Data 3). These analyses suggest that our strategy is robust to find previously known variants given sufficient power.

We next tested all the 16 M variants for association with the 22 conditions available in the GERA cohort and five different inheritance models (Supplementary Figs. 4–25). This analysis identified 94 independent loci associated with 17 phenotypes at a usual genome-wide significance level ($p < 5.0 \times 10^{-8}$) of which 63 for 14 phenotypes were also experiment-wide significant ($p < 2.0 \times 10^{-8}$) after considering correction for the different models of inheritance (see "Methods") (Supplementary Data 4). According to the GWAS Catalog, 68 of the 94 genome-wide significant loci had been previously reported to be associated with the same disease (Supplementary Data 5), whereas 26 of them correspond to previously unreported loci with associations across 16 phenotypes (Table 1).

Of these 26 new loci, 16 correspond to common, 3 to low-frequency, and 7 to rare variants. Only a fraction of the 26 new loci would have been genome-wide significant by using individual imputation panels (Fig. 1e), namely 19/26 using HRC, 14/26 using 1000G Phase 3, 13/26 using UK10K or 14/26 using GoNL. In addition, the lead marker for three of the novel signals was an indel, not covered by HRC, further confirming the benefits of combining multiple panels with our approach.

**Identification of recessive variants with large effects.** The implementation of refined GWAS strategies not only increases the number of associated variants, but also allows the identification of loci with large impact on the disease. Among the variants that were not detected under the additive model, and hence are expected to be missed by the majority of current GWAS, we highlight three variants with large recessive effects. First, an intronic indel in the *CACNB4* gene, rs201654520, associated with a nearly 20-fold increase in risk for cardiovascular disease (MAF = 0.017, OR [CI 95%] = 19.0 [5.5–65.8], $p = 4.3 \times 10^{-8}$). *CACNB4* encodes the β4 subunit of the voltage-dependent calcium channel. This subunit contributes to the flux of calcium ions into the cell by increasing peak calcium current and triggering muscle contraction. Interestingly, an intronic single nucleotide polymorphism (SNP) within *CACNB4*, rs150793926, was associated with idiopathic dilated cardiomyopathy in African Americans[26], but this variant is not in linkage disequilibrium (LD) with rs201654520 (LD $r^2$ [27] = 0.0016 for European ancestry and LD $r^2$ = 0.0 for African ancestry).

A second recessive variant with large effect, rs77704739, near the *PELO* gene, is associated with a fourfold risk for type 2 diabetes (MAF = 0.036, OR [CI 95%] = 4.3 [2.7–6.9], $p = 1.75 \times 10^{-8}$). An independent signal that is about 112 K base pairs away (rs870992, LD $r^2$ = 0.0009) was previously associated with type 2 diabetes in the Greenlandic population, also with a recessive effect[12]. To provide insights into the underlying molecular mechanisms in disease, we interrogated comprehensive catalogs of genetic effects on gene expression: eQTLGen Consortium[28] and GTEx[29]. The rs77704739 variant was significantly associated with gene expression of *PELO* in multiple tissues, including diabetes-relevant tissues such as adipose tissue, skeletal muscle, and pancreas. Colocalization analyses showed a probability higher than 0.8 in several tissues, including subcutaneous adipose tissue and skeletal muscle, suggesting this gene as the effector transcript (Fig. 2a, b, and Supplementary Data 6). In addition, we found that the lead variants in the *PELO* locus overlap with active promoter annotations in human pancreatic islets and open chromatin sites highly bounded by islet-specific transcription factors[30,31] (Fig. 2c).

Third, a rare indel, rs557998486, located near the *THUMPD2* gene, is associated with age-related macular degeneration (MAF = 0.009, OR = 10.5, $p = 2.75 \times 10^{-8}$). Interestingly, the fact that we found no SNPs in LD with this lead indel further confirms the

benefits of multiple reference panel imputation strategies that include alternative forms of variation. The lead indel rs557998486 overlaps DNAse I hypersensitivity sites in retinal and iris cell lines[32], highlighting a candidate open chromatin region that is also predicted to be an enhancer assigned to the *THUMPD2* gene according to GeneHancer[33]. One of the variants with the highest LD with rs557998486 (rs116649730, LD $r^2 = 0.32$) is associated with reduced expression of its nearest gene, *THUMPD2* (z-score = −4.85, $p = 1.25 \times 10^{-6}$), according to eQTLGen Consortium data.

Our empirical evaluation also demonstrates that our imputation approach is accurate for the new variants, including those with a large recessive effect, and that the combination of multiple reference panels increases the certainty of the imputed genotypes (Table 1 and Supplementary Data 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 and 69[5]. 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 Data 8). 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[34]) 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).

The variant rs77704739 variant near *PELO* was associated with type 2 diabetes (OR-recessive [95% CI] = 1.9 [1.4–2.6], $p = 4.95 \times 10^{-4}$) and metformin use (OR-recessive [95% CI] = 2.3 [1.6–3.4], $p = 3.8 \times 10^{-5}$) in the UK Biobank[5] (Supplementary Data 8), also only under the recessive model.
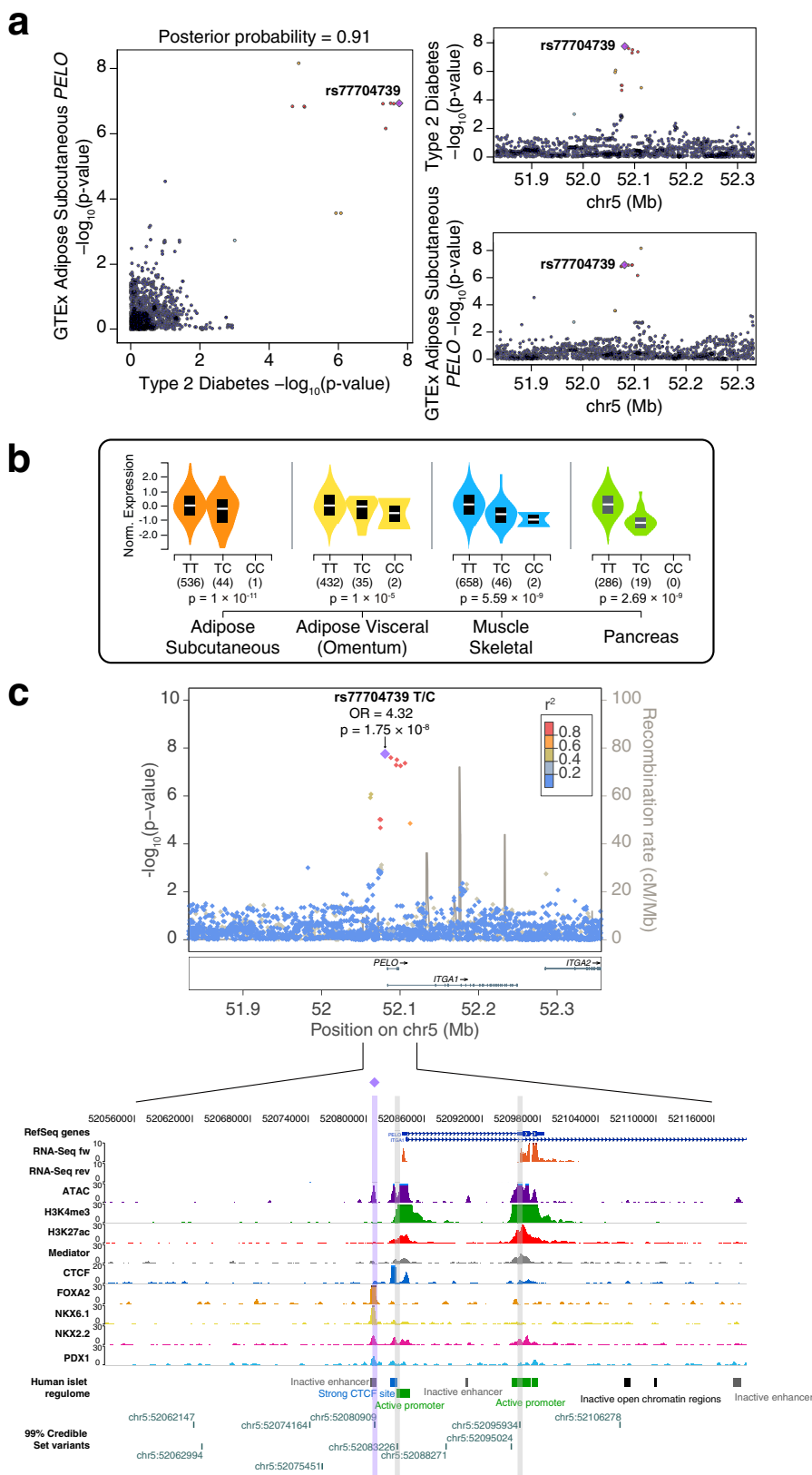
We further sought replication of the association within the *CACNB4* gene with cardiovascular disease in FinnGen, a cohort of ~218 K Finnish individuals with an average age of 63, as it includes individuals with a higher average age (63 vs 56 in UK Biobank) and the risk of developing a cardiovascular disease is well-known to increase with age[35]. 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 Fig. 26a). Despite the 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 Fig. 26). 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 Data 8). 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[36].

We also sought replication of the recessive association of rs557998486 near *THUMPD2* gene with macular degeneration in FinnGen. rs557998486 was associated with increased risk of macular degeneration in UK Biobank under the recessive model

**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%) | P-value | Lowest p-value model OR (CI 95%) | P-value | Dominance deviation P-value | Best panel empirical $R^2$ [a] |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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) | $1.55 \times 10^{-8}$ | $6.66 \times 10^{-1}$ | 1.000 |
| Asthma (9209/47,428) | 8 | DLC1 | 13,164,746 | rs10112506 | A/G | 0.390 | Dominant | 0.94 (0.91–0.97) | $8.61 \times 10^{-6}$ | 0.89 (0.86–0.93) | $1.54 \times 10^{-8}$ | $2.86 \times 10^{-4}$ | 0.998 |
| | 5 | ETF1 | 137,858,067 | rs154073 | C/T | 0.429 | Recessive | 1.09 (1.06–1.13) | $6.06 \times 10^{-8}$ | 1.18 (1.12–1.25) | $4.23 \times 10^{-9}$ | $9.28 \times 10^{-3}$ | 0.991 |
| | 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) | $4.14 \times 10^{-8}$ | $8.10 \times 10^{-1}$ | –[b] |
| 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)[c] | $3.54 \times 10^{-8}$ | $6.00 \times 10^{-1}$ | 0.802 |
| | 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) | $2.44 \times 10^{-8}$ | $6.00 \times 10^{-1}$ | 0.988 |
| Cardiovascular (15,009/41,628) | 1 | DCLRE1B | 114,448,752 | rs10858023 | C/T | 0.350 | Dominant | 1.09 (1.06–112) | $3.26 \times 10^{-8}$ | 1.14 (1.09–1.19) | $2.11 \times 10^{-9}$ | $1.94 \times 10^{-2}$ | 0.996 |
| | 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) | $4.32 \times 10^{-8}$ | $4.36 \times 10^{-6}$ | 0.973 |
| Major depression disorder (7264/49,373) | 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) | $3.15 \times 10^{-9}$ | $1.10 \times 10^{-6}$ | –[b] |
| Type 2 diabetes (6967/49,670) | 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) | $1.75 \times 10^{-8}$ | $1.92 \times 10^{-7}$ | 0.998 |
| Hemorrhoids (9129/47,508) | 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) | $2.03 \times 10^{-8}$ | – | 0.933 |
| Hernia abdominopelvic (6291/50,346) | 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) | $2.03 \times 10^{-8}$ | $6.87 \times 10^{-1}$ | 0.997 |
| | 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) | $1.27 \times 10^{-8}$ | – | 0.647 |
| 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) | $4.42 \times 10^{-8}$ | $6.85 \times 10^{-3}$ | 1.000 |
| | 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) | $4.42 \times 10^{-8}$ | $4.80 \times 10^{-3}$ | 0.907 |
| | 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) | $1.71 \times 10^{-8}$ | $8.01 \times 10^{-1}$ | 0.826 |
| Iron deficiency anemia (2439/54,198) | 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) | $3.80 \times 10^{-8}$ | $8.92 \times 10^{-2}$ | 0.948 |
| Macular degeneration (3685/52,952) | 2 | THUMPD2 | 40,010,523 | rs557998486 | T/TG | 0.009 | Recessive | 1.07 (0.81–1.41) | $6.28 \times 10^{-1}$ | 10.5[d] | $2.75 \times 10^{-8}$ | – | 0.865 |
| Osteoporosis (5399/51,238) | 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) | $4.79 \times 10^{-8}$ | – | 0.851 |
| Psychiatric (8624/48,013) | 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) | $1.66 \times 10^{-8}$ | $2.57 \times 10^{-1}$ | 0.994 |
| Peripheral Vascular disease (4301/52,336) | 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)[c] | $4.26 \times 10^{-8}$ | $6.32 \times 10^{-6}$ | 0.923 |
| | 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)[c] | $3.55 \times 10^{-8}$ | $1.35 \times 10^{-3}$ | 0.998 |
| Acute reaction to stress (4314/52,323) | 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) | $4.56 \times 10^{-8}$ | – | 0.875 |
| Varicose veins (2483/54,154) | 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)[c] | $2.13 \times 10^{-8}$ | $9.58 \times 10^{-4}$ | 0.939 |
| | 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) | $5.00 \times 10^{-8}$ | $9.88 \times 10^{-1}$ | 0.995 |
| | 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) | $8.41 \times 10^{-9}$ | – | 0.752 |

CHR chromosome, Position position hg19, Alleles non-effect allele/effect allele, MAF minor allele frequency, OR odds ratio, CI confidence interval.
aEmpirical r-squared correlation ($R^2$) between imputed and sequenced allele dosage for the best panel from our in silico analysis using an array of UK10K genotypes as a backbone and imputing with 1000G, HRC, and GoNL.
bThis variant is not present in UK10K.
cOdds ratio and confidence interval for heterozygous/odds ratio and confidence interval for effect allele homozygous calculated using the method het+hom from SNPTEST.
dOdds ratio calculated using the recessive allele frequency-based test (RAFT)[6].

(OR [CI 95%] = 7.6 [1.5–37.3], $p = 4.1 \times 10^{-2}$), eye surgery (beta [CI 95%] = 1.6 [0.6–2.6], $p = 1.17 \times 10^{-3}$) (Supplementary Data 8), and with increased C-reactive protein, a known biomarker for macular degeneration[37,38] (beta [CI 95%] = 1.1

[0.7–1.5], $p = 1.15 \times 10^{-4}$) (Supplementary Data 9). In FinnGen this variant was not significantly associated although it showed the same direction of effect. However, the meta-analysis did not reach the genome-wide significance (rs557998486 $p = 9.6 \times 10^{-6}$)

**Fig. 2 Functional characterization of the rs77704739 recessive association near the *PELO* gene. a** 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. **b** Violin plot from GTEx showing that the recessive rs77704739 variant significantly modifies the expression of *PELO* gene in subcutaneous ($n = 581$ independent samples) and visceral adipose tissue ($n = 469$ independent samples), skeletal muscle ($n = 706$ independent samples) and pancreas ($n = 305$ independent samples). The box plots have lines extending from the boxes (whiskers) indicating variability outside the upper and lower quartiles. GTEx V7 was used for colocalization analyses, whereas GTEx V8 was used to generate the violin plots. **c** 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.

and had a high heterogeneity (heterogeneity $I^2 = 87.1$, heterogeneity $p = 4.3 \times 10^{-4}$).

**Detection ranges of the different inheritance models.** Our findings provide an empirical overview of the detection range of five different inheritance models and show how each of them captures a fraction of the genetic variants associated with complex traits. Compared to current GWAS that usually only consider additive allelic effects, we found three different scenarios. Among all the 94 associated loci identified, 12 showed genome-wide significance only under the additive model, 62 under both additive and non-additive models, and 20 showed genome-wide significance only when non-additive tests were applied (Fig. 3a). To further classify these variants, we tested whether any of the 62 variants associated with both additive and non-additive models deviate from additivity through a dominance deviation test[11]. Eleven of these 62 variants (17.7%) showed significant deviation from additivity (dominance deviation test $p < 0.05$). However, variants not showing a significant deviation from additivity may become significant for other models with larger sample sizes. Altogether, the dominance deviation test over the 93 autosomal loci identified 62 (66%) additive and 24 (25.5%) non-additive associations, and 8 undetermined. Based on the smallest GWAS $p$-value, we further classified non-additive associations into 9 recessive, 13 dominant, 8 heterodominat, and 7 genotypic (Supplementary Data 4).

We also observed that each of the available models for association testing has a different range of detection. To identify the 94 genome-wide associated loci, the additive test, as expected, was the most sensitive model (74 loci), followed by the genotypic (59 loci), the dominant (56 loci), the recessive (43 loci), and the heterodominant (32 loci). When considering known loci, 48 of the 68 previously reported loci were identified by more than one model in our analysis, and almost half of these (22 loci) with all five models. In contrast, of the 26 newly discovered variants, only 8 were identified with multiple models, whereas the majority of them (18 loci), were detected only with the additive (6 loci), the genotypic (4 loci), the recessive (4 loci), and the dominant (3 loci) model. Of note, 13 out of 26 (50%) novel loci were only identified by non-additive models.

To further investigate to what extent the additive model captures non-additive signals, and how much this depends on sample size, we carried out power calculations on loci that we identified only under a non-additive model, such as rs201654520 within *CACNB4* gene and rs77704739 near the *PELO* gene. These power calculations showed that the additive test would require a population sample size of at least 370,646 individuals to detect the recessive association of rs201654520 in *CACNB4* (Fig. 3b), and at least 188,637 individuals to capture the recessive signal of rs77704739 near the *PELO* gene (Fig. 3c), while the population sample size required for the recessive model was only 21,021 and 67,611, respectively. In this study, we were able to identify both associations with a modest sample size by using the most well-suited disease model.

**The GUIDANCE framework.** The complete methodology described here and used for the analysis of the GERA cohort was integrated into a framework, called GUIDANCE. GUIDANCE allows the analysis of genome-wide genotyped data in a single execution in distributed computing infrastructures without the need for extensive computational expertise or constant user intervention. The GUIDANCE workflow requires quality-controlled genotyped data as input and provides association results, graphical outputs, and statistical summaries. Integrating state-of-the-art tools with in-house code written in Java, bash, and R[39], GUIDANCE efficiently performs large-scale GWAS, including (1) the pre-phasing of haplotypes, (2) the imputation of genotypes using multiple reference panels, (3) the association testing for different inheritance models and the integration of the results from different panels, (4) a cross-phenotype analysis when more than one phenotype is available in the cohort (Supplementary Data 10), and, finally, (5) the generation of summary statistics tables and graphic representations of the results (Supplementary Fig. 27), for both the autosomes and the X chromosome. While GUIDANCE can be executed as a standalone compact program it can also be used in modules (Supplementary Fig. 28), which makes it adaptable to existing frameworks and provides a higher level of control to users.

GUIDANCE runs in distributed computing platforms, including the cloud, without requiring a broad background in these environments. This is feasible since GUIDANCE was implemented on top of the COMP Superscalar Programming Framework (COMPSs)[40]. The GUIDANCE workflow was implemented as a sequential Java program containing the calls to the GWAS tools, encapsulated in Java methods, and selected as tasks, while COMPSs controls the execution of those tasks on the underlying distributed infrastructure. The source code, the pre-compiled binaries, and the documentation to use GUIDANCE are available at http://cg.bsc.es/guidance.

## Discussion
Current genetic association studies are undergoing fundamental strategic and methodological changes to identify new associations. The gradual incorporation of WES and WGS, together with the continuous efforts to increase the sample size, impose computational and methodological challenges that cannot always be overcome. In addition, the inclusion of the X chromosome, and non-additive models during association testing, despite increasing the computational burden, can also lead to the identification of new associations. In this study, we demonstrate the value of applying a comprehensive GWAS strategy, including denser imputation strategies, the X chromosome, and non-additive association tests, to an existing large-scale genetic resource, the GERA cohort. We show that by applying more innovative imputation protocols we increased the number and the type of variants tested for association, including low-frequency and rare SNPs as well as alternative forms of variation, such as indels. Our analysis of the GERA cohort shows that between 13 and 20 of the genome-wide significant associations (14–21%) would not have

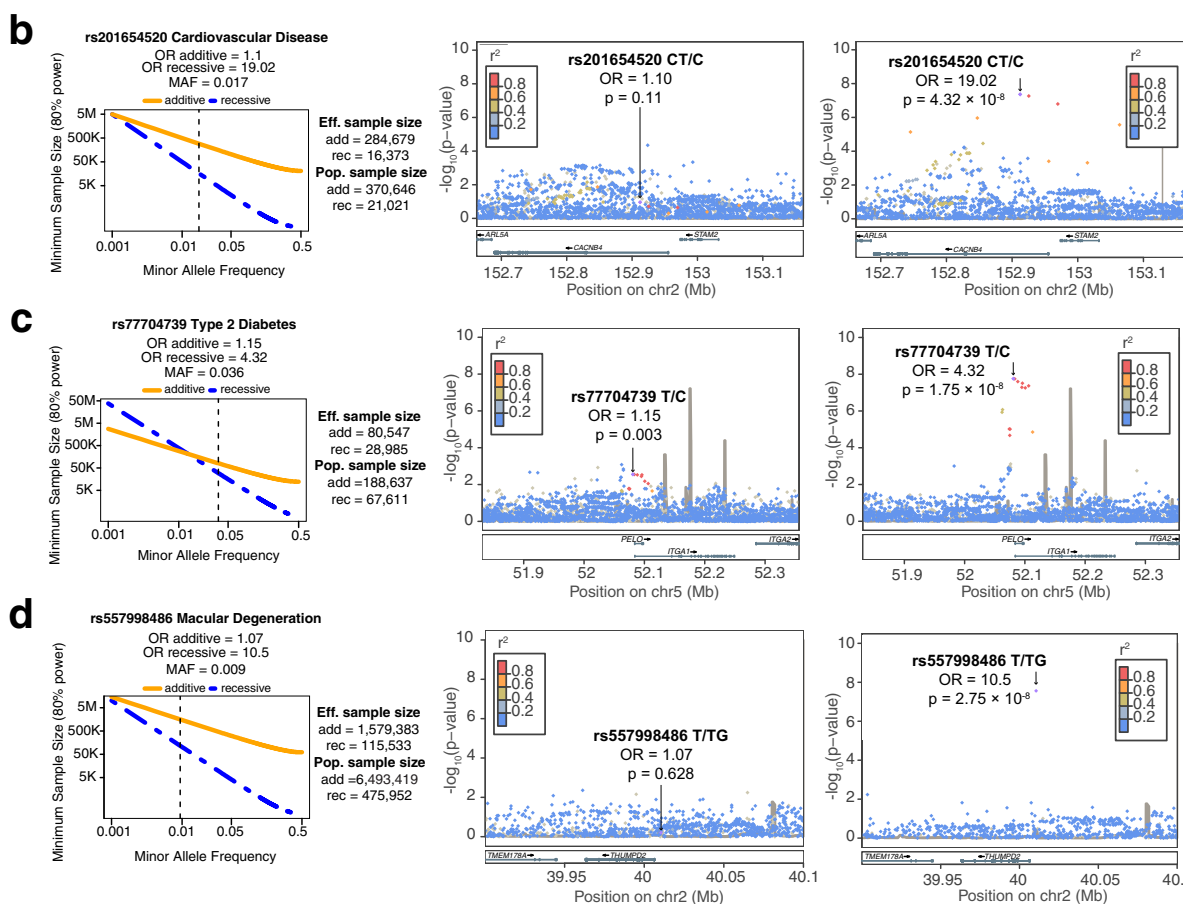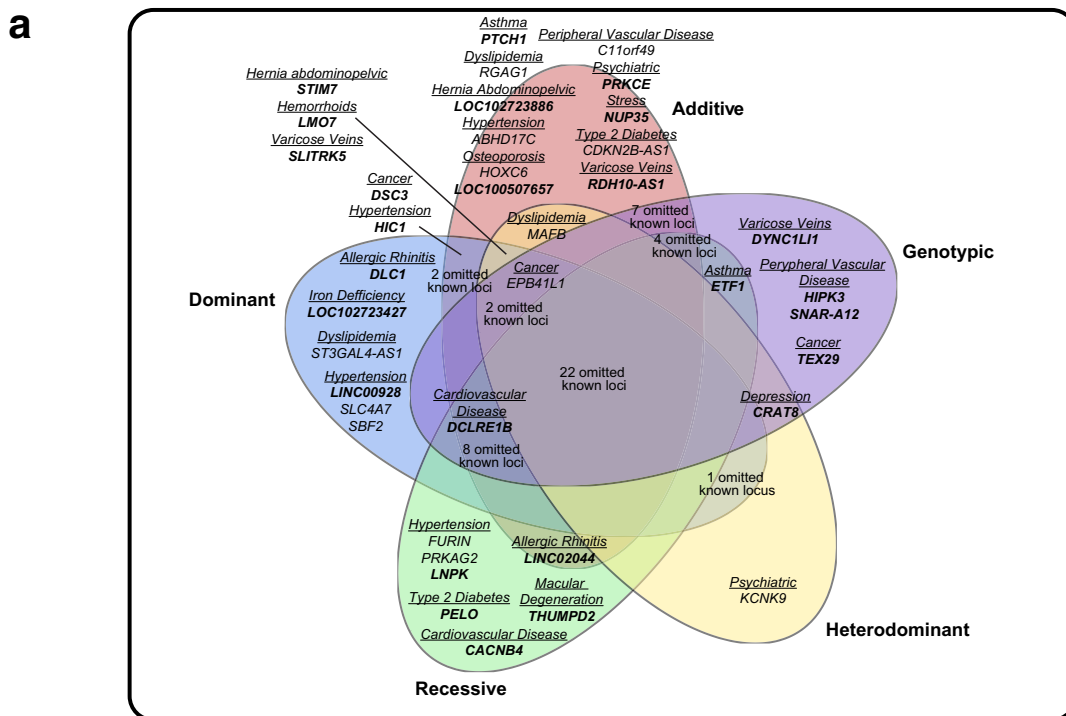**Table 2 Replication of new associations with UK Biobank.**

| CHR | rsID (alleles) (MAF) | Best model | Stage 1. GERA Discovery | | | | | Stage 2. UK Biobank replication | | | | | Stage 1 + Stage 2. Meta-analysis | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Phenotype (cases/controls) | Additive OR (CI 95%) | P-value | Best model OR (CI 95%) | P-value | Field (cases/controls or sample size) | Additive OR (CI 95%) | P-value | Lowest p-value model OR (CI 95%) | P-value | 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) | $2.44 \times 10^{-8}$ | 1.50 (1.30–1.72) | $2.44 \times 10^{-8}$ | Self-reported: chronic lymphocytic 360,904 | 2.13 (1.14–3.97) | $3.50 \times 10^{-2}$ | 2.13 (1.14–3.97) | $3.50 \times 10^{-2}$ | 1.52 (1.33–1.74) | $1.60 \times 10^{-9}$ | 1.52 (1.33–1.74) | $1.60 \times 10^{-9}$ |
| 1 | rs2494196 (C/A) (0.274) | Additive | Hernia abdominopelvic (6291/50,346) | 1.13 (1.08–1.18) | $2.03 \times 10^{-8}$ | 1.13 (1.08–1.18) | $2.03 \times 10^{-8}$ | Self-reported: umbilical hernia (328/360,813) | 1.42 (1.21–1.67) | $2.31 \times 10^{-5}$ | 1.42 (1.21–1.67) | $2.31 \times 10^{-5}$ | 1.15 (1.10–1.19) | $5.35 \times 10^{-11}$ | 1.15 (1.10–1.19) | $5.35 \times 10^{-11}$ |
| | | | | | | | | K40 Inguinal hernia (13,365/347,829) | 1.09 (1.06–1.12) | $3.95 \times 10^{-10}$ | 1.09 (1.06–1.12) | $3.95 \times 10^{-10}$ | 1.10 (1.08–1.12) | $7.78 \times 10^{-17}$ | 1.10 (1.08–1.12) | $7.78 \times 10^{-17}$ |
| | | | | | | | | K41 Femoral hernia (475/360,719) | 1.44 (1.26–1.64) | $1.24 \times 10^{-7}$ | 1.44 (1.26–1.64) | $1.24 \times 10^{-7}$ | 1.16 (1.11–1.21) | $2.26 \times 10^{-12}$ | 1.16 (1.11–1.21) | $2.26 \times 10^{-12}$ |
| | | | | | | | | K42 Umbilical hernia (2623/358,571) | 1.29 (1.22–1.37) | $1.14 \times 10^{-17}$ | 1.29 (1.22–1.37) | $1.14 \times 10^{-17}$ | 1.19 (1.15–1.22) | $2.94 \times 10^{-22}$ | 1.19 (1.15–1.22) | $2.94 \times 10^{-22}$ |
| | | | | | | | | K43 Ventral hernia (2470/358,724) | 1.18 (1.11–1.25) | $1.77 \times 10^{-7}$ | 1.18 (1.11–1.25) | $1.77 \times 10^{-7}$ | 1.15 (1.11–1.19) | $1.99 \times 10^{-14}$ | 1.15 (1.11–1.19) | $1.99 \times 10^{-14}$ |
| 2 | rs557998486 (T/TG) (0.009) | Recessive | Macular degeneration (3685/52,952) | 1.07 (0.81–1.41) | $6.28 \times 10^{-1}$ | 10.5[a] | $2.75 \times 10^{-8}$ | Eye problems/disorders: Macular degeneration (2726/115,164) | 0.98 (0.72–1.32) | $8.81 \times 10^{-1}$ | 7.58 (1.54–37.32) | $4.1 \times 10^{-2}$ | 1.01 (0.82–1.24)[b] | $7.91 \times 10^{-1c}$ | 26.51 (7.57–92.85)[b] | $3.29 \times 10^{-8c}$ |
| 5 | rs77704739 (T/C) (0.036) | Recessive | Type 2 diabetes (6967/49,670) | 1.15 (1.05–1.26) | $2.80 \times 10^{-3}$ | 4.32 (2.70–6.92) | $1.75 \times 10^{-8}$ | Self-reported: diabetes (14,114/347,027) | 1.03 (0.97–1.09) | $3.87 \times 10^{-1}$ | 1.88 (1.35–2.6) | $4.95 \times 10^{-4}$ | 1.06 (1.01–1.12) | $1.78 \times 10^{-2}$ | 2.46 (1.88–3.21) | $4.68 \times 10^{-11}$ |

CHR chromosome, Position position hg19, Alleles non-effect allele/effect allele, MAF minor allele frequency, OR odds ratio.
[a]Odds ratio calculated using the recessive allele frequency-based test (RAFT).
[b]Obtained through a mega-analysis with UK Biobank using the expected method from SNPTEST.
[c]Obtained using METAL method SAMPLESIZE to combine the p-values taking into account the sample size and direction of effect.

been identified when using a single reference panel. Likewise, our analysis in the GERA cohort demonstrates that 21% of the associations would be missed by only testing the additive model. Overall, 27.6% of associations would not have been identified by applying the commonly used HRC and additive model

association testing. The inclusion of the X chromosome within the study allowed us to identify an intronic variant (rs67648651) associated with dyslipidemia in the *RGAG1* gene among the previously known associated loci since a variant in LD with our top variant (rs5985471, LD $r^2 = 0.92$ for European ancestry) has

**Fig. 3 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 13 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 475,952.

already been associated with low-density lipoprotein cholesterol[20]. This supports the value of analyzing the X chromosome, which is not considered in the majority of GWAS.

In this study, we show the potential of identifying large effect recessive associations by maximizing the use of current reference panels and testing different inheritance models, as exemplified by the associations with type 2 diabetes, cardiovascular disease, and macular degeneration with variants near *PELO*, *CACNB4*, and *THUMPD2*, respectively. This strategy opens new avenues for future analyses in large-scale biobanks, as demonstrated with our power calculations, which show that even the largest available GWAS meta-analyses or biobanks would not have enough power to identify these associations using only the additive model. For example, the *CACNB4* gene, associated with cardiovascular disease, would require a sample size equivalent to 370,000 individuals when using the additive test, 17 times larger than the required sample size under a recessive analysis. After considering all the supporting evidence illustrated with many examples in this study, the results suggest that these new associations deserve future validations and follow-up analysis. Therefore, this study demonstrates the importance of a comprehensive analysis including non-additive models when performing GWAS, which can increase the discovery not only on GWAS relying on genotyping array data, but also on WES or WGS association studies.

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

To easily apply this strategy to genetic studies we present GUIDANCE, a standalone and easy-to-use application that allows an efficient and comprehensive GWAS analysis in different computing platforms, such as cloud and high-performance computing architectures. GUIDANCE is designed to allocate an unlimited number of reference panels. This can be useful for GWAS performed in specific diverse populations where the addition of a population-specific reference panel alongside the commonly used ones could be an advantage. This feature can also be useful when incorporating reference panels that have a better ascertainment of other type of variation, such as structural variants, with the most commonly used ones. In addition, it is possible to launch the execution by steps, incorporating the previously obtained results, avoiding repeating all the computations. In a moment where the community is facing computational and methodological challenges due to the growing complexity and size of genetic datasets, the availability of robust and

complete analysis platforms can improve the efficiency of genetic studies, standardizing analysis strategies among meta-analysis cohorts to ensure consistency.

Finally, to share our results with the community and to promote the analysis of non-additive inheritance models in GWAS, a public searchable database including additive and non-additive summary statistics for 16 M of variants and 22 phenotypes is available at the Common Metabolic Diseases Knowledge Portal (https://cmdkp.org/), and full summary statistics both at the Common Metabolic Diseases Knowledge Portal and at http://cg.bsc.es/guidance.

## Methods

**Empirical evaluation of the imputation strategy.** To empirically evaluate our imputation strategy, we extracted the corresponding genotypes in UK10K for each variant present in the genotyping array used in GERA to build an in silico array of 599,208 genotyped variants for the 3781 UK10K individuals. We then imputed the data using the 599,208 variants as a backbone, using 1000G, GoNL, and HRC as reference panels. We used IMPUTE2 to impute the genotypes. We also combined the imputed results based on the highest imputation accuracy according to the IMPUTE2-info values as we did with the GERA cohort. We transformed the IMPUTE2 genotype probabilities into genotype dosages and calculated the r-squared correlation ($R^2$) between the imputed genotype dosages and the sequenced genotype dosages from UK10K, termed allelic dosage $R^2$ (Supplementary Fig. 1).

**GUIDANCE workflow description.** By combining and integrating state-of-the-art GWAS analysis tools into the COMP Superscalar programming Framework (COMPSs), we developed GUIDANCE, a standalone application that performs haplotype phasing, genome-wide imputation, association testing, and PheWAS analysis of large GWAS datasets (Supplementary Fig. 27).

As shown in Supplementary Fig. 27, GUIDANCE's workflow starts with quality-controlled genotype data and ends with providing association results, graphical outputs, and statistical summaries.

Once everything is settled in the GUIDANCE configuration file, GUIDANCE performs an efficient two-stage imputation procedure, by pre-phasing the genotypes into whole haplotypes followed by genotype imputation itself[41]. SHAPEIT2[42] or EAGLE2[43] and IMPUTE2[23] or MINIMAC4[44] can be used for pre-phasing and genotype imputation, respectively. In addition, GUIDANCE accepts one or multiple reference panels, allowing the integration of the results obtained from all panels by selecting for each variant the genotypes from the reference panel that provides the highest imputation accuracy according to the IMPUTE2-info score or MINIMAC2 $r^2$ (Supplementary Fig. 29). GUIDANCE also performs a post-imputation quality control to eliminate low-quality imputed variants under the basis of the IMPUTE2-info score or MINIMAC2 $r^2$ and the MAF.

After genotype imputation and post-imputation quality control, GUIDANCE applies SNPTEST for association testing, where additive, dominant, recessive, heterodominant, and genotype models can be analyzed. Here, the user can decide to include several covariates for the association test, such as principal components to adjust for population stratification, or any other confounders. GUIDANCE also allows testing for multiple phenotypes or for a single phenotype with different covariates in the same execution. After association testing, variants are filtered by the deviation from Hardy–Weinberg equilibrium (HWE) *p*-value. Finally, GUIDANCE generates summary reports for each trait with all the inheritance models tested in the association and the corresponding graphical representation, i.e., Manhattan and quantile–quantile (Q–Q) plots (Supplementary Fig. 4–25), also providing a matrix identifying cross-phenotype associations (Supplementary Data 10).

GUIDANCE can be executed as a standalone compact program or as independent modules (see Supplementary Fig. 28 for a list of independent modules) to facilitate the use of GUIDANCE into existing frameworks.

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

## The analysis of GERA cohort

*GERA cohort description.* GERA cohort data was obtained through dbGaP under accession phs000674.v1.p1. For further information about the specific phenotypes (ICD-9-CM codes) included in GERA, visit its website on dbGaP (https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/GetPdf.cgi?id=phd004308). The Resource for Genetic Epidemiology Research on Aging (GERA) Cohort was created by a RC2 Grand Opportunity grant that was awarded to the Kaiser Permanente Research Program on Genes, Environment, and Health (RPGEH) and the UCSF Institute for Human Genetics (AG036607; Schaefer/Risch, PIs). The RC2 project enabled genome-wide SNP genotyping (GWAS) to be conducted on a cohort of over 100 K adults who were members of the Kaiser Permanente Medical Care Plan, Northern California Region (KPNC), and participating in its RPGEH. The resulting GERA cohort is composed of 42% of males, 58% of females, and ranges in age from 18 to over 100 years old with an average age of 63 years at the time of the RPGEH survey (2007). Nineteen percent of the individuals are from non-European ancestry, while 81% are described as white non-Hispanic participants. After an explicit requirement of consent by email, data from 78,486 participants were deposited in dbGaP, with similar demographic characteristics to those of the initial genotyped cohort.

*Quality control.* A subset of 62,281 subjects of European ancestry underwent quality control analyses. A three-step quality control protocol was applied using PLINK[45,46], and included two stages of SNP removal and an intermediate stage of sample exclusion.

The exclusion criteria for genetic markers consisted of proportion of missingness $\geq 0.05$, 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 $\geq 0.125$ from which we removed the individual with the highest proportion of missingness), sample call rates $\geq 0.02$, and population structure showing more than four standard deviations within the distribution of the study population according to the first seven principal components (Supplementary Fig. 30). After QC, 56,637 subjects remained for the analysis (Supplementary Data 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. For the haplotype phasing of chromosome X, we used the --chrX flag required for SHAPEIT to only phase female samples since males only have one X chromosome, and to impute missing data in male samples. We did the same imputing genotypes using IMPUTE2 as it requires the -chrX flag alongside the sample file with the sex information. In the association test, we used -method newml from SNPTEST to ignore samples with missing sex or males encoded wrongly (males should be coded 0/1, as homozygote females), and to assume a model of full X inactivation. Hence, the logistic regression model assumes a complete inactivation of one allele in females and equal effect size between males and females. For heterogeneity between males and females, and to allow a complete inactivation of the X chromosome in females, we used -stratify_on sex to separate the effects and the baselines of males and females, accounting for hemizygosity for males, while for females, we followed an autosomal model.

Finally, we excluded variants with HWE controls $p < 1 \times 10^{-6}$ and with a case count for homozygous for the alternative allele below three in the final results for the recessive and genotypic model, as we observed a trend for genomic inflation and deflation in the recessive and genotypic model before removing these variants.

*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[47], all the genome-wide significant variants were collapsed into ranges (500 kb) that define each associated locus.

To distinguish between known or new associated regions, for each top variant we looked for any proxy variant with an LD $r^2 > 0.35$ in the GWAS catalog (accession 5 September 2019) associated with the same phenotype or a related one (for example, bone mineral density, cholesterol levels or diastolic/systolic blood pressure phenotypes for osteoporosis, dyslipidemia or hypertension, respectively). HLA regions at chromosome 6 were excluded since the particularities of these

regions required further detailed studies on their LD pattern. Proxies were selected using LDlink (https://ldlink.nci.nih.gov/)[48].

We defined an experiment-wide significant $p$-value cutoff of $p < 2.0 \times 10^{-8}$ by applying the Bonferroni correction for 2.5 effective test ($5.0 \times 10^{-8}/2.5$ effective test). This factor of 2.5 was obtained from a simulation study when four genetic models (additive, dominant, recessive, and genotypic) are used[49] since the genetic models are not independent. However, a new simulation study including the heterodominant model should be done for a more accurate effective number of tests.

## Replication with UK biobank

*Phenotype curation.* UK Biobank participants agreed to provide detailed information about their lifestyle, environment, and medical history, to donate biological samples (for genotyping and for biochemical assays), to undergo measures, and to have their health followed (http://www.ukbiobank.ac.uk/).

When collecting and analyzing a wide range of phenotypes from the UK Biobank, a central challenge was the curation and harmonization of the vast array of categorizations, variable scalings, and follow-up responses. Fortunately, to this end, the PHEnome Scan ANalysis Tool (or PHESANT: https://github.com/MRCIEU/PHESANT)[50] performs much of the transformations and recodings required to generate meaningful, interpretable phenotypes.

We have made further adjustments based on user feedback, owing to the value of transparency in generating our phenotype guidelines. Applying these changes to the PHESANT source code, phenotypes were parsed using our modified version (github.com/astheeggeggs/PHESANT) on a virtual machine on the Google Cloud Platform.

We first restricted to the subset of European individuals, before passing the resultant phenotypic data to PHESANT. The 'variable list' file and 'data-coding' file, whose formats were defined in the original version of PHESANT were updated as new phenotypes were added in the latest UK Biobank release. Recodings of variables, and inherent orderings of categorical variables, are defined in the 'data-coding' file. The 'Excluded' column of the 'variable list' file defines the collection of variables that we do not wish to interrogate.

A high-level overview of the PHESANT pipeline, our defaults, and the associated short flags for the phenomescan.r code are displayed in Supplementary Fig. 31. In addition to the inverse-rank normalization applied to the collection of continuous phenotypes, we also consider the raw version of the continuous phenotype, with no transformation applied to the data.

Curation of the ICD10 codes was carried out separately for computational efficiency. For the ICD10 phenotype, individuals are assigned a vector of ICD10 diagnoses. We truncated these codes to two digits, and assigned each individual to either case or control status for that ICD10 code in turn by checking if their vector contains that code. Throughout, we assumed the data contained no missingness, so the sum of cases and controls throughout was the number of individuals in our 'European' subset of the UK Biobank data. As in the PHESANT categorical (multiple) phenotypes, ICD10 code case/control phenotypes were removed if <50 individuals had the diagnosis.

*Association testing and meta-analysis for UK Biobank phenotypes.* We performed the association testing for the curated phenotypes as implemented in SNPTEST for additive, dominant, recessive, heterodominant, and genotypic inheritance models, as it has been described in the "Analyzing GERA cohort using GUIDANCE" section. For all genotypic variants identified in the discovery stage, we assigned the recessive model after we identified it as the underlying model.

After the association testing, we filtered and ordered all the phenotypes based on the $p$-value for the best model of inheritance obtained from the GERA cohort analysis, with special consideration to equivalent phenotypes or related traits.

With the association testing results of both GERA cohort and UK Biobank, we meta-analyzed the results using METAL[51]. We use the inverse variance-weighted fixed-effect model for all the variants except for the rs557998486 variant associated with macular degeneration, since its beta, calculated with the em method from SNPTEST, was inflated. Therefore, we performed a sample size-based meta-analysis, which converts the direction of the effect and the $p$-value into a $z$-score.

For biomarkers, only the results from the first visit were taken into account since <10% of the cases were present in the second visit.

*Association testing and meta-analysis with FinnGen.* We used SAIGE[52] 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[53]. 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[54], 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[54]. The approximate Bayes factor (ABF) for each variant was estimated as

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

where

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

$$z = \frac{\beta}{\text{SE}}. \quad (3)$$

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

$$\text{Posterior probability}_i = \frac{\text{ABF}_i}{T}, \quad (4)$$

where $\text{ABF}_i$ corresponds to the approximate Bayes' factor for the marker $i$, and $T$ represents the sum of all the ABF values enclosed in the interval. As commonly employed by SNPTEST, this calculation assumes that the prior of the $\beta$ is a Gaussian with mean 0 and variance 0.04.

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

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

To determine whether any identified overlap between GERA GWAS loci and eQTLGen or GTEx eQTLs was due to a true shared association signal, we performed a colocalization analysis. Colocalization was assessed by a Bayesian test using summary statistics from both studies[55]: summary statistics from the *cis* eQTLGen and GTEx were downloaded from the eQTLGen website and GTEx portal, respectively. The test was performed using the R package coloc v3.2-1[55–57]. The test provided a posterior probability for the GWAS locus and the eQTL to share the same causal variant(s).

We integrated available epigenomic datasets to examine the role of human pancreatic islet transcriptional regulation underlying rs77704739 association with type 2 diabetes. We used the WashU EpiGenome Browser (http://epigenomegateway.wustl.edu/browser/, last access on July 2019) and previously published RNA-seq, ATAC-seq, and ChIP-seq assays of H3K4me3, H3K27ac, Mediator, CTCF, and islet transcription factors (FOXA2, MAFB, NKX2.2, NKX6.1, and PDX1) in human pancreatic islets[30,31] and islet regulome annotations[31].

*Comparison of power calculation under different inheritance models.* Power calculations were performed using libraries epiR[58] and GeneticsDesign[59]. For each variant, the power was computed across different allele frequencies and sample sizes. Frequencies of homozygous for different allele frequencies were estimated assuming Hardy–Weinberg equilibrium. The sample size needed to achieve 80% power was plotted against the allele frequency. For the additive model, we chose the observed odds ratio for the additive model, whereas the observed odds ratio for the recessive model was chosen for the recessive model.

**Reporting summary.** Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability
The complete summary statistics from this study have been deposited and are available to download at the Common Metabolic Diseases Knowledge Portal (https://cmdkp.org/) and at http://cg.bsc.es/gera_summary_stats.

## Code availability
GUIDANCE[60] is available at https://gitlab.bsc.es/computational-genomics/guidance (https://doi.org/10.5281/zenodo.4446121). A zipped folder containing the code used to generate the results presented in the article can be found in GitLab (https://gitlab.bsc.es/computational-genomics/guidance) and Zenodo (https://zenodo.org/record/4446121#.YASVmWhKiUk) and can be cited with the DOI 10.5281/zenodo.4446121. It contains example configuration files, explanations on how to execute the framework, a containerized version of the pipeline, and the source code. More detailed instructions can be found on its website (http://cg.bsc.es/guidance/).

## References
1. Welter, D. et al. The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res.* **42**, D1001–D1006 (2014).
2. Manolio, T. A. et al. Finding the missing heritability of complex diseases. *Nature* **461**, 747–753 (2009).
3. Bonas-Guarch, S. et al. Re-analysis of public genetic data reveals a rare X-chromosomal variant associated with type 2 diabetes. *Nat. Commun.* **9**, 321 (2018).
4. Taliun, D. et al. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature* **590**, 290–299 (2021).
5. Bycroft, C. et al. The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
6. Tam, V. et al. Benefits and limitations of genome-wide association studies. *Nat. Rev. Genet.* **20**, 467–484 (2019).
7. Zhu, Z. et al. Dominance genetic variation contributes little to the missing heritability for human complex traits. *Am. J. Hum. Genet.* **96**, 377–385 (2015).
8. Salanti, G. et al. Underlying genetic models of inheritance in established type 2 diabetes associations. *Am. J. Epidemiol.* **170**, 537–545 (2009).
9. Lettre, G., Lange, C. & Hirschhorn, J. N. Genetic model testing and statistical power in population-based association studies of quantitative traits. *Genet. Epidemiol.* **31**, 358–362 (2007).
10. Antonarakis, S. E. & Beckmann, J. S. Mendelian disorders deserve more attention. *Nat. Rev. Genet.* **7**, 277–282 (2006).
11. Wood, A. R. 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).
12. 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).
13. Moltke, I. et al. A common Greenlandic TBC1D4 variant confers muscle insulin resistance and type 2 diabetes. *Nature* **512**, 190–193 (2014).
14. Steinthorsdottir, V. et al. A variant in CDKAL1 influences insulin response and risk of type 2 diabetes. *Nat. Genet.* **39**, 770–775 (2007).
15. Lenz, T. L. et al. Widespread non-additive and interaction effects within HLA loci modulate the risk of autoimmune diseases. *Nat. Genet.* **47**, 1085–1090 (2015).
16. 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).
17. Hoffmann, T. J. et al. Imputation of the rare HOXB13 G84E mutation and cancer risk in a large population-based cohort. *PLoS Genet.* **11**, e1004930 (2015).
18. Boomsma, D. I. et al. The Genome of the Netherlands: design, and project goals. *Eur. J. Hum. Genet.* **22**, 221–227 (2014).
19. Genome of the Netherlands Consortium. Whole-genome sequence variation, population structure and demographic history of the Dutch population. *Nat. Genet* **46**, 818–825 (2014).
20. UK10K Consortium. et al. The UK10K project identifies rare variants in health and disease. *Nature* **526**, 82–90 (2015).
21. 1000 Genomes Project Consortium. et al. An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56–65 (2012).
22. McCarthy, S. et al. A reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.* **48**, 1279–1283 (2016).
23. Howie, B. N., Donnelly, P. & Marchini, J. A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genet.* **5**, e1000529 (2009).
24. Mahajan, A. et al. Fine-mapping type 2 diabetes loci to single-variant resolution using high-density imputation and islet-specific epigenome maps. *Nat. Genet.* **50**, 1505–1513 (2018).
25. Fritsche, L. G. et al. A large genome-wide association study of age-related macular degeneration highlights contributions of rare and common variants. *Nat. Genet.* **48**, 134–143 (2016).
26. Xu, H. et al. A genome-wide association study of idiopathic dilated cardiomyopathy in African Americans. *J. Pers. Med.* **8**, 11 (2018).
27. Hill, W. G. & Robertson, A. Linkage disequilibrium in finite populations. *Theor. Appl Genet* **38**, 226–231 (1968).
28. Võsa, U. et al. Unraveling the polygenic architecture of complex traits using blood eQTL metaanalysis. Preprint at bioRxiv https://doi.org/10.1101/447367 (2018).
29. GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* **45**, 580–585 (2013).
30. Pasquali, L. et al. Pancreatic islet enhancer clusters enriched in type 2 diabetes risk-associated variants. *Nat. Genet.* **46**, 136–143 (2014).

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

32. Consortium, E. P. The ENCODE (ENCyclopedia Of DNA Elements) Project. *Science* **306**, 636–640 (2004).

33. Fishilevich, S. et al. GeneHancer: genome-wide integration of enhancers and target genes in GeneCards. *Database* **2017**, https://doi.org/10.1093/database/bax028 (2017).

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

35. Yazdanyar, A. & Newman, A. B. The burden of cardiovascular disease in the elderly: morbidity, mortality, and costs. *Clin. Geriatr. Med.* **25**, 563–577 (2009). vii.

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

37. Han, X. et al. Using Mendelian randomization to evaluate the causal relationship between serum C-reactive protein levels and age-related macular degeneration. *Eur. J. Epidemiol.* **35**, 139–146 (2020).

38. Molins, B., Romero-Vazquez, S., Fuentes-Prior, P., Adan, A. & Dick, A. D. C-reactive protein as a therapeutic target in age-related macular degeneration. *Front. Immunol.* **9**, 808 (2018).

39. R Foundation for Statistical Computing. *R: A Language and Environment for Statistical Computing* (2019).

40. Lordan, F. et al. ServiceSs: an interoperable programming framework for the Cloud. *J. Grid Comput.* **12**, 67–91 (2014).

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

42. Delaneau, O., Zagury, J. F. & Marchini, J. Improved whole-chromosome phasing for disease and population genetic studies. *Nat. Methods* **10**, 5–6 (2013).

43. Loh, P. R. et al. Reference-based phasing using the Haplotype Reference Consortium panel. *Nat. Genet.* **48**, 1443–1448 (2016).

44. Das, S. et al. Next-generation genotype imputation service and methods. *Nat. Genet.* **48**, 1284–1287 (2016).

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

46. Chang, C. C. et al. Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 7 (2015).

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

48. Machiela, M. J. & Chanock, S. J. LDlink: a web-based application for exploring population-specific haplotype structure and linking correlated alleles of possible functional variants. *Bioinformatics* **31**, 3555–3557 (2015).

49. Mercader, J. M. et al. Altered brain-derived neurotrophic factor blood levels and gene variability are associated with anorexia and bulimia. *Genes Brain Behav.* **6**, 706–716 (2007).

50. Millard, L. A. C., Davies, N. M., Gaunt, T. R., Davey Smith, G. & Tilling, K. Software Application Profile: PHESANT: a tool for performing automated phenome scans in UK Biobank. *Int. J. Epidemiol.* **47**, 29–35 (2018).

51. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190–2191 (2010).

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

53. Lumley, T. rmeta: Meta-analysis. R package version 3.0 (2018).

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

55. Giambartolomei, C. et al. A Bayesian framework for multiple trait colocalization from summary association statistics. *Bioinformatics* **34**, 2538–2545 (2018).

56. Wallace, C. Statistical testing of shared genetic control for potentially related traits. *Genet. Epidemiol.* **37**, 802–813 (2013).

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

58. Stevenson, M. et al. epiR: tools for the analysis of epidemiological. R package version 1.0, 2 edn (2019).

59. Warnes, G., Duffy, D., Man, M., Qiu, W. & Lazarus, R. GeneticsDesign: functions for designing genetics studies. R package version 1.52.0 edn (2019).

60. Guindo-Martínez, M. et al. The impact of non-additive genetic associations on age-related complex diseases. https://doi.org/10.5281/zenodo.4446121 (2021).

61. Lim, E. T. et al. A novel test for recessive contributions to complex diseases implicates Bardet-Biedl syndrome gene BBS10 in idiopathic type 2 diabetes and obesity. *Am. J. Hum. Genet.* **95**, 509–520 (2014).

## Acknowledgements

## Author contributions

M.G.-M., J.M.M., and D.T. conceived and planned the main analyses. M.G.-M., J.M.M., and D.T. wrote and edited the manuscript. S.B.-G. designed and performed the quality control. M.G.-M. performed the main analysis, and with collaboration of R.A., M.P., C.R.-C., F.S., and J.E. R.M.B. developed GUIDANCE on top of COMPSs. S.B.-G. and I.M.-E. performed the functional characterization. C.S. performed the dominance deviation test and the gene expression analysis. J.M.M., C.E.C., J.B.C., E.A., A.L., K.A., D.P., and J.C.F. contributed with UK Biobank data and analysis. S.R. and M.K. contributed with FinnGen data and analysis. J.M.M. and D.T. supervised the study. All authors reviewed and approved the final manuscript.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1038/s41467-021-21952-4.

**Correspondence** and requests for materials should be addressed to J.M.M. or D.T.

**Peer review information** *Nature Communications* thanks the anonymous reviewers for their contribution to the peer review of this work. Peer reviewer reports are available.

**Reprints and permission information** is available at http://www.nature.com/reprints

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.