

1 Components of genetic associations across 2,138 2 phenotypes in the UK Biobank highlight novel adipocyte 3 biology

4
5 Yosuke Tanigawa^{1*}, Jiehan Li^{2,3*}, Johanne Marie Justesen^{1,2,4}, Heiko Horn^{5,6},
6 Matthew Aguirre^{1,7}, Christopher DeBoever^{1,8}, Chris Chang⁹, Balasubramanian Narasimhan^{1,10},
7 Kasper Lage^{5,6,11}, Trevor Hastie^{1,10}, Chong Yon Park², Gill Bejerano^{1,7,12,13}, Erik Ingelsson^{2,3*+},
8 Manuel A. Rivas^{1*+}

- 9 1. Department of Biomedical Data Science, School of Medicine, Stanford University, Stanford, CA, USA.
- 10 2. Department of Medicine, Division of Cardiovascular Medicine, Stanford University, Stanford, CA, USA.
- 11 3. Stanford Cardiovascular Institute, Stanford University, Stanford, CA 94305.
- 12 4. Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health and Medical Sciences,
13 University of Copenhagen, Copenhagen, Denmark
- 14 5. Department of Surgery, Massachusetts General Hospital, Harvard Medical School, Boston, MA, USA.
- 15 6. Broad Institute of MIT and Harvard, Cambridge, MA, USA.
- 16 7. Department of Pediatrics, Stanford University School of Medicine, Stanford University, Stanford, CA, USA.
- 17 8. Department of Genetics, Stanford University, Stanford, CA, USA.
- 18 9. Grail, Inc., 1525 O'Brien Drive, Menlo Park, CA, USA.
- 19 10. Department of Statistics, Stanford University, Stanford, CA, USA.
- 20 11. Institute for Biological Psychiatry, Mental Health Center Sct. Hans, University of Copenhagen, Roskilde, Denmark.
- 21 12. Department of Developmental Biology, Stanford University, Stanford, CA, USA.
- 22 13. Department of Computer Science, Stanford University, Stanford, CA, USA.

23
24 *These authors contributed equally

25 +Corresponding authors

26 Abstract

27 Population-based biobanks with genomic and dense phenotype data provide opportunities for
28 generating effective therapeutic hypotheses and understanding the genomic role in disease
29 predisposition. To characterize latent components of genetic associations, we applied truncated
30 singular value decomposition (DeGAs) to matrices of summary statistics derived from genome-
31 wide association analyses across 2,138 phenotypes measured in 337,199 White British
32 individuals in the UK Biobank study. We systematically identified key components of genetic
33 associations and the contributions of variants, genes, and phenotypes to each component. As
34 an illustration of the utility of the approach to inform downstream experiments, we report
35 putative loss of function variants, rs114285050 (*GPR151*) and rs150090666 (*PDE3B*), that
36 substantially contribute to obesity-related traits, and experimentally demonstrate the role of
37 these genes in adipocyte biology. Our approach to dissect components of genetic associations
38 across the human phenome will accelerate biomedical hypothesis generation by providing
39 insights on previously unexplored latent structures.

40 Introduction

41 Human genetic studies have been profoundly successful at identifying regions of the genome
42 contributing to disease risk^{1,2}. Despite these successes, there are challenges to translating
43 findings to clinical advances, much due to the extreme polygenicity and widespread pleiotropy
44 of complex traits, which are the presence of genetic effects of a variant across multiple
45 phenotypes and multiple variants across a single phenotype³⁻⁵. In retrospect, this is not
46 surprising given that most common diseases are multifactorial. However, it remains unclear
47 exactly which factors, acting alone or in combination, contribute to disease risk and how those
48 factors are shared across diseases. With the emergence of sequencing technologies, we are
49 increasingly able to pinpoint alleles, possibly rare and with large effects, which may aid in
50 therapeutic target prioritization⁶⁻¹³. Furthermore, large population-based biobanks, such as the
51 UK Biobank, have aggregated data across tens of thousands of phenotypes¹⁴. Thus, an
52 opportunity exists to characterize the phenome-wide landscape of genetic associations across
53 the spectrum of genomic variation, from coding to non-coding, and rare to common.

54 Singular value decomposition (SVD), a mathematical approach developed by differential
55 geometers¹⁵, can be used to combine information from several (likely) correlated vectors to form
56 basis vectors, which are guaranteed to be orthogonal and to explain maximum variance in the
57 data, while preserving the linear structure that helps interpretation. In the field of human
58 genetics, SVD is routinely employed to infer genetic population structure by calculating principal
59 components using the genotype data of individuals¹⁶.

60 To address the pervasive polygenicity and pleiotropy of complex traits, we propose an
61 application of truncated SVD (TSVD), a reduced rank approximation of SVD¹⁷⁻¹⁹, to characterize
62 the underlying (latent) structure of genetic associations using summary statistics computed for
63 2,138 phenotypes measured in the UK Biobank population cohort¹⁴. We applied our novel
64 approach, referred to as DeGAs – Decomposition of Genetic Associations – to assess
65 associations among latent components, phenotypes, variants, and genes. We highlight its
66 application to body mass index (BMI), myocardial infarction (MI), and gallstones, motivated by
67 high polygenicity in anthropometric traits, global burden, and economic costs, respectively. We
68 assess the relevance of the inferred key components through GREAT genomic region ontology
69 enrichment analysis²⁰ and functional experiments. The results from DeGAs applied to protein-
70 truncating variants (PTV) dataset indicated strong associations of targeted PTVs to obesity-
71 related traits, while phenome-wide association analyses (PheWAS) uncovered the differential
72 region-specific regulation of our top candidates in fat deposition. For these reasons, we
73 prioritized adipocytes as our experimental model system for the follow-up functional studies of
74 our candidate genes. Given that the roles of adipocytes in regulating metabolic fitness have
75 been established at both local and systemic levels of pathology associated with obesity, it is
76 likely that the differentiation and function of adipocytes may shape the effects of our candidate
77 genes at the cellular and molecular level.

78 Results

79 DeGAs method overview

80 We generated summary statistics by performing genome-wide association studies (GWAS) of
81 2,138 phenotypes from the UK Biobank (Fig. 1a, Supplementary Tables S1-S2). We performed
82 variant-level quality control, which includes linkage-disequilibrium (LD) pruning and removal of
83 variants in the MHC region, to focus on 235,907 variants for subsequent analyses. Given the
84 immediate biological consequence, subsequent downstream implications, and medical
85 relevance of coding variants and predicted protein-truncating variants (PTVs), commonly
86 referred to as loss-of-function variants^{12,21,22}, we performed separate analyses on three variant
87 sets: (1) all directly-genotyped variants, (2) coding variants, and (3) PTVs (Supplementary Fig.
88 S1). To eliminate unreliable estimates of genetic associations, we selected associations with p-
89 values < 0.001 , and standard error of beta value or log odds ratio of less than 0.08 and 0.2,
90 respectively, for each dataset. The Z-scores of these associations were aggregated into a
91 genome- and phenome-wide association summary statistic matrix W of size $N \times M$, where N
92 and M denote the number of phenotypes and variants, respectively. N and M were 2,138 and
93 235,907 for the “all” variant group; 2,064 and 16,135 for the “coding” variant group; and 628 and
94 784 for the PTV group. The rows and columns of W correspond to the GWAS summary
95 statistics of a phenotype and the phenome-wide association study (PheWAS) of a variant,
96 respectively. Given its computational efficiency compared to the vanilla SVD, we applied TSVD
97 to each matrix and obtained a decomposition into three matrices $W = USV^T$ (U: phenotype, S:
98 variance, V: variant). This reduced representation of $K = 100$ components altogether explained
99 41.9% (all), 62.8% (coding) and 75.5% (PTVs) of the variance in the original summary statistic
100 matrices (Fig. 1b, Methods, Supplementary Fig. S2).

101 In DeGAs framework, we employ these latent components characterized from a densely
102 phenotyped population-based cohort to investigate the genetics of common complex traits (Fig.
103 1c). To characterize each latent component and identify the relevant component given a
104 phenotype, a gene, or a variant, or vice versa; we defined five different quantitative scores:
105 phenotype squared cosine score, phenotype and variant contribution score, variant contribution
106 score, and gene contribution score. The squared cosine scores quantify the relative importance
107 of component for a given phenotype or gene, and are defined based on the squared distance of
108 a component from the origin on the latent space. Contribution scores quantify relative
109 importance of a phenotype, variant, or gene to a given component and is defined based on the
110 squared distance of a phenotype, variant, or gene from the origin²³ (Fig. 1d, Methods). Using
111 scores, DeGAs identifies the key latent components for a given complex trait and annotated
112 them with the driving phenotypes, genes, and variants (Fig. 1c, Methods). We performed
113 biological characterization of DeGAs components with the genomic region enrichment analysis
114 tool (GREAT)²⁰ followed by functional experiments in adipocytes (Fig. 1e).

115 Characterization of latent structures of DeGAs

116 The PCA plots show the projection of phenotypes and variants onto DeGAs latent components.
117 (Fig. 2a-b). For the variant PCA plot, we overlay biplot annotation as arrows to interpret the
118 direction of the components (Fig. 2b). Overall, we find that the first five DeGAs components can
119 be attributed to: 1) fat-free mass that accounts for the “healthy part” of body weight²⁴ (32.7%,
120 Supplementary Table S3) and two intronic variants in *FTO* (rs17817449: contribution score of
121 1.15% to PC1, rs7187961: 0.41%); and a genetic variant proximal to *AC105393.1* (rs62106258:
122 0.46%); 2) whole-body fat mass (61.5%) and the same three *FTO* and *AC105393.1* variants
123 (rs17817449: 0.97%, rs7187961: 0.28%, rs62106258: 0.27%); 3) bioelectrical impedance
124 measurements (38.7%), a standard method to estimate body fat percentage^{25,26}, and genetic
125 variants proximal to *ACAN* (rs3817428: 0.64%), *ADAMTS3* (rs11729800: 0.31%), and
126 *ADAMTS17* (rs72770234: 0.29%); 4) eye meridian measurements (80.9%), and two intronic
127 variants in *WNT7B* (rs9330813: 5.73%, rs9330802: 1.14%) and a genetic variant proximal to
128 *ATXN2* (rs653178: 0.96%); and 5) bioelectrical impedance and spirometry measures (45.4%
129 and 26.0%, respectively) and genetic variants proximal to *FTO* (rs17817449: 0.17%), *ADAMTS3*
130 (rs11729800: 0.11%), and *PSMC5* (rs13030: 0.11%) (Fig. 2c-d, Supplementary Table S4).

131 To highlight the ability of DeGAs to capture related set of phenotypes, genes, and
132 variants in genetic associations, we applied TSVD to the missing-value imputed and Z-score
133 transformed phenotype matrix and characterized the first 100 latent components (Methods).
134 Using the individual and phenotype PCA plots, we found a fewer number of components
135 explains most of the variance and several phenotypes, such as traffic intensity of the nearest
136 major road and creatinine (enzymatic) in urine, are dominantly driving the top phenotypic PCs
137 (Supplementary Fig. S4-S5). We applied GWAS for each of the decomposed phenotypes
138 (Supplementary Fig. S6). Through the genetic correlation analysis with the derived summary
139 statistics, we found non-zero genetic correlations among the phenotypic PCs (Supplementary
140 Fig. S7-S8).

141 Applying DeGAs components for BMI, MI, and gallstones

142 To illustrate the application of DeGAs in characterizing the genetics of complex traits, we
143 selected three phenotypes, BMI, MI, and gallstones given the large contribution of
144 anthropometric traits on the first five components, that ischemic heart diseases is a leading
145 global fatal and non-fatal burden, and that gallstones is a common condition with severe pain
146 and large economic costs where polygenic risk factors are largely unknown^{27,28}. We identified
147 the top three key components for these three phenotypes with DeGAs using the “all” variants
148 dataset.

149 For BMI, we find that the top three components of genetic associations (PC2, PC1, and
150 PC30) altogether explained over 69% of the genetic associations (47%, 18%, and 4%,
151 respectively, Supplementary Fig. S3a). The top two components (PC2 and PC1) corresponded
152 to components of body fat (PC2) and fat-free mass measures (PC1), as described above. PC30
153 was driven by fat mass (28.7%) and fat-free mass (6.8%), but also by non-melanoma skin
154 cancer (7.72%) – linked to BMI in epidemiological studies²⁹ – and childhood sunburn (7.61%)
155 (Fig. 3a, Supplementary Table S4).

156 For MI, a complex disease influenced by multiple risk factors³⁰, we found that the top
157 components were attributed to genetics of lipid metabolism (PC22, high-cholesterol, statin
158 intake, and *APOC1*), alcohol intake (PC100), and sleep duration and food intake (PC83, 25.2%)
159 that collectively corresponded to 36% of the genetic associations (Fig. 3a, Supplementary Fig.
160 S3b, S9-S10, Supplementary Table S4).

161 Cholelithiasis is a disease involving the presence of gallstones, which are concretions
162 that form in the biliary tract, usually in the gallbladder³¹. We found that the top components
163 contributing to gallstones corresponded to associations with fresh fruit (PC72) and water intake
164 (PC64), as well as bioelectrical impedance of whole body (PC67) corresponding to 51% of
165 genetic associations altogether (Fig. 3a, Supplementary Fig. S3c, S9, S11, Supplementary
166 Table S4). We confirmed the robustness of these results with respect to the selection of number
167 of components, K (Methods, Supplementary Fig. S12-S16).

168 Biological characterization of DeGAs components

169 To provide biological characterization of the key components, we applied the genomic region
170 enrichment analysis tool (GREAT)²⁰ to dissect the biological relevance of the identified
171 components with both coding and non-coding variants. Given the coverage of the manually
172 curated knowledge of mammalian phenotypes, we focused on the mouse genome informatics
173 (MGI) phenotype ontology and set $p = 5 \times 10^{-6}$ as the Bonferroni-corrected statistical
174 significance threshold (Method)³². For each key component, we applied GREAT and found an
175 enrichment for the mouse phenotypes consistent with the phenotypic description of our
176 diseases of interest²⁰. The top component for BMI, identified as the body fat measures
177 component (PC2), showed enrichment of several anthropometric terms including abnormally
178 short feet (brachypodia) (MP:0002772, binomial fold = 9.04, $p = 1.3 \times 10^{-23}$), increased birth
179 weight (MP:0009673, fold = 6.21, $p = 1.3 \times 10^{-11}$), and increased body length (MP:0001257,
180 binomial fold = 3.01, $p = 1.3 \times 10^{-36}$) (Fig. 3B, Supplementary Table S5). For MI, we found
181 enrichment of cardiac terms, such as artery occlusion (PC22, MP:0006134, fold = 15.86, $p =$
182 1.14×10^{-25}) and aortitis (PC22, MP:0010139, aorta inflammation, fold = 9.36, $p =$
183 3.41×10^{-31}) (Supplementary Fig. S17, Supplementary Table S6). Similarly, for gallstones, the
184 top enrichment was for abnormal circulating phytosterol level (PC72, MP:0010075, fold = 11.54,
185 $p = 5.51 \times 10^{-11}$), which is known to be involved in gallstone development³³ (Supplementary
186 Fig. S18, Supplementary Table S7).

187 To test the specificity of the enriched ontology terms while considering the correlation
188 structure within ontology terms, we took the top five enriched terms for each DeGAs component,
189 obtained the list of genes annotated with these top terms, and measured their pairwise gene set
190 similarity across 100 DeGAs components using Jaccard index (Methods). Jaccard index is a set
191 similarity measure ranges between zero and one where one means the complete match and
192 zero means complete disjoint of the two sets. We found the median of the pairwise similarity to
193 be 0.029 (Supplementary Fig. S19).

194 Coding and protein truncating variants

195 Given the challenges with interpreting genetic associations across thousands of possibly
196 correlated phenotypes and diverse variant functional categories, we applied DeGAs to coding
197 variant-phenotype associations and PTV associations. For the coding dataset, we identified
198 PC2 and PC1 as the top two key components for BMI, with 51% and 14% of phenotype squared
199 contribution scores, respectively (Supplementary Fig. S20). The major drivers of these two
200 components include fat mass measurements (55.2% of phenotype contribution score for PC2),
201 fat-free mass measurements (33.3%, PC1), genetic variants on *MC4R* (3.7% gene contribution
202 score for PC2), and *ZFAT* (3.4% gene contribution score for PC1) (Supplementary Fig. S21-
203 S22, Supplementary Table S4).

204 Predicted PTVs are a special class of protein-coding genetic variants with possibly
205 strong effects on gene function^{9,12,21,34}. More importantly, strong effect PTV-trait associations
206 can uncover promising drug targets, especially when the direction of effect is consistent with
207 protection of human disease. Using the PTV dataset, we identified PC1 and PC3 as the top two
208 key components for BMI, with 28% and 12% of phenotype squared contribution scores,
209 respectively (Supplementary Fig. S23). The major drivers of PC1 were weight-related
210 measurements, including left and right leg fat-free mass (5.0% and 3.7% of phenotype
211 contribution score for PC1, respectively), left and right leg predicted mass (4.9% each), weight
212 (4.6%), and basal metabolic rate (4.6%), whereas the drivers of PC3 included standing height
213 (13.7%), sitting height (8.1%), and high reticulocyte percentage (6.4%) (Fig. 4a, Supplementary
214 Table S4). Top contributing PTVs to PC1 included variants in *PDE3B* (19.0%), *GPR151*
215 (12.3%), and *ABTB1* (8.5%), whereas PC3 was driven by PTVs on *TMEM91* (8.6%), *EML2-AS1*
216 (6.7%), and *KIAA0586* (6.0%) (Fig. 4b, Supplementary Table S4).

217 Based on stop-gain variants in *GPR151* (rs114285050) and *PDE3B* (rs150090666)
218 being key contributors to the top two components of genetic associations for PTVs and BMI
219 (Fig. 4c), we proceeded to detailed phenome-wide association analysis (PheWAS) assessing
220 associations of these PTVs with anthropometric phenotypes. PheWAS analysis of these
221 variants confirmed strong associations with obesity-related phenotypes including waist
222 circumference (*GPR151*, marginal association beta = -0.065, $p = 2.5 \times 10^{-8}$), whole-body fat
223 mass (*GPR151*, beta = -0.069, $p = 1.4 \times 10^{-7}$), trunk fat mass (*GPR151*, beta = -0.071, $p =$
224 1.5×10^{-7}), hip circumference (*PDE3B*, beta = 0.248, $p = 1.8 \times 10^{-11}$), right leg fat-free mass
225 (*PDE3B*, beta = 0.129, $p = 4.2 \times 10^{-8}$) and body weight (*PDE3B*, beta = 0.177, $p = 4.6 \times 10^{-8}$)
226 (Fig. 4d, Supplementary Fig. S24, Supplementary Table S8-9). Among 337,199 White British
227 individuals, we found 7,560 heterozygous and 36 homozygous carriers of the *GPR151* variant
228 and 947 heterozygous carriers of *PDE3B* variants. To assess the effect of the PTVs on BMI, a
229 commonly-used measure of obesity, we performed univariate linear regression analysis with
230 age, sex, and the first four genetic PCs as covariates and found that heterozygous and carriers
231 of *GPR151* PTVs showed 0.324 kg/m² lower BMI than the average UK Biobank participant ($p =$
232 4.13×10^{-7}). We did not find evidence of association with homozygous carriers (N = 28; $p =$
233 0.665), presumably due to lack of power (Supplementary Fig. S25). Heterozygous carriers of
234 *PDE3B* PTVs showed 0.647 kg/m² higher BMI ($p = 2.09 \times 10^{-4}$) than the average UK Biobank
235 participant (Supplementary Fig. S26).

236 Functional experiments for candidate genes in cellular models of 237 adipocytes

238 We sought to illustrate the potential application of DeGAs in prioritizing therapeutic targets using
239 functional follow-up experiments. Several of our most interesting findings were observed from
240 strong associations between PTVs and obesity-related traits. Variants in *GPR151* and *PDE3B*
241 are the two strongest contributors, albeit in opposite directions, to the top component (PC1)
242 driving the genetic associations between PTVs and BMI (Fig. 4 a-c). In addition to BMI, a simple
243 indicator of overall body fat level, PheWAS studies have suggested strong correlations between
244 regional body fat distribution and these two PTVs, with *GPR151* being more considerably
245 associated with waist circumference and trunk fat (Fig. 4d), while *PDE3B* was more notably
246 related to hip circumference and lower-body fat (Fig. S24). Regional fat deposition is more
247 accurately reflected by the local development and function of adipocytes in terms of size,
248 number and lipid content. In order to explore how these two candidates that regulate body fat
249 composition differently, we chose to study their impacts on biological characteristics of
250 adipocytes. Specifically, the expression and function of *PDE3B* and *GPR151* were evaluated in
251 mouse 3T3-L1 and human Simpson-Golabi-Behmel Syndrome (SGBS) cells, two well-
252 established preadipocyte models used for studying adipocyte differentiation (i.e. adipogenesis)
253 and function^{35,36}.

254 First, we demonstrated that both genes were expressed in preadipocytes, but showed
255 different expression patterns when cells were transforming into mature adipocytes: *PDE3B*
256 increased dramatically during both mouse and human adipogenesis, while *GPR151* maintained
257 a low expression level throughout the differentiation (Fig. 5a-b). Next, to explore the causal
258 relationships between gene expression and adipogenesis, we introduced short interfering RNA
259 (siRNA) against *Pde3b* and *Gpr151*, respectively, into 3T3-L1 preadipocytes and monitored the
260 impact of gene knockdown on conversion of preadipocytes to adipocytes. Knockdown of
261 *Gpr151* (Fig. 5c) drastically impaired adipocyte differentiation, as evidenced by lowered
262 expression of adipogenesis markers (*Pparg*, *Cebpa* and *Fabp4*) (Fig. 5d), as well as the
263 reduced formation of lipid-containing adipocytes (Fig. 5e-f). Further, to test the functional
264 capacity of the fat cells lacking *Gpr151*, we performed a lipolysis assay - an essential metabolic
265 pathway of adipocytes and thus, a key indicator of adipocyte function - on mature adipocytes
266 derived from preadipocytes transfected with either scrambled siRNA (scRNA) or si*Gpr151*. Not
267 surprisingly, *Gpr151*-deficient lipid-poor adipocytes showed dramatically lower lipolysis, along
268 with impaired capability of responding to isoproterenol (ISO), a β -adrenergic stimulus of lipolysis
269 (Fig. 5g). These data suggest that *GPR151* knockdown in adipocyte progenitor cells may block
270 their conversion into mature adipocytes; thus, preventing the expansion of adipose tissue.
271 These results are directionally consistent with our DeGAs and univariate regression analysis
272 showing that *GPR151* PTVs are associated with lower obesity and fat mass, especially central
273 obesity (e.g. waist circumference and trunk fat mass) (Fig. 4d).

274 To further analyze the functional impact of *GPR151* in adipocytes, we generated an
275 overexpression model of *GPR151* by infecting 3T3-L1 preadipocytes with virus expressing Flag-
276 tagged human *GPR151* driven by either EF1 α or aP2 promotor (Supplementary Fig. S27a).
277 Overexpression of *GPR151* by both constructs were confirmed at the gene and protein levels
278 (Supplementary Fig. S27b-d). However, despite the substantial effect of *Gpr151* knockdown on

279 adipogenesis (Fig. 5), overexpression of *GPR151* in preadipocytes failed to influence adipocyte
280 differentiation significantly, as shown by similar levels of adipogenic markers compared to the
281 non-infected controls (Supplementary Fig. S27e-f). To eliminate the potential masking effects of
282 any unperturbed cells in the partially infected cell population, we specifically selected *GPR151*-
283 overexpressing cells by staining Flag-*GPR151* positive cells with APC-conjugated flag antibody
284 and sorted APC+ and APC- cells from the differentiating adipocyte cultures (Supplementary Fig.
285 S27g-l). In both EF1 α - and aP2-driven *GPR151* overexpression models, *GPR151* mRNA levels
286 were enriched in APC+ cells compared to APC- cells (Supplementary Fig. S27m-n). However,
287 APC+ cells expressed genes characteristics of differentiating adipocytes in a similar level to that
288 of APC- cells (Supplementary Fig. S27m-n). These data conclude that overexpression of
289 *GPR151* in preadipocytes cannot further enhance adipogenesis, suggesting that the
290 endogenous level of *GPR151* in preadipocytes may be sufficient to maintain the normal
291 differentiation potential of preadipocytes. Although *GPR151* is predominantly expressed in the
292 brain, especially in hypothalamic neurons that control appetite and energy expenditure³⁷, we
293 identified for the first time that the *GPR151* protein is present in both subcutaneous and visceral
294 adipose tissue from mice (SAT and VAT), albeit in a very low level (Supplementary Fig. S27o).
295 Together with our gain- and loss-of-function studies of *GPR151* in preadipocyte models, we
296 infer that the regulatory role of *GPR151* in body weight may involve both central and peripheral
297 effects. The minimal but indispensable presence of *GPR151* in adipose progenitor cells in
298 generating lipid-rich adipocytes seems to represent an important mechanism by which *GPR151*
299 promotes obesity.

300 In contrast to *GPR151*, knockdown of *Pde3b* in 3T3-L1 preadipocytes (Supplementary
301 Fig. S28a) showed no significant influence on adipogenesis and lipolysis (under either basal or
302 β -adrenergic stimulated conditions), as compared to scRNA-transfected controls
303 (Supplementary Fig. S28b-e). Since PDE3B is expressed primarily in differentiated adipocytes
304 (Fig. 5a-b), future research efforts should be concentrated on studying the metabolic role of
305 PDE3B in mature adipocytes. As an essential enzyme that hydrolyzes both cAMP and cGMP,
306 PDE3B is known to be highly expressed in tissues that are important in regulating energy
307 homeostasis, including adipose tissue³⁸. *Pde3b* whole-body knockout in mice reduces the
308 visceral fat mass³⁹ and confers cardioprotective effects⁴⁰. There is a growing body of evidence
309 that cardiometabolic health is linked to improved body fat distribution (i.e. lower visceral fat,
310 higher subcutaneous fat) in a consistent direction⁴¹. Our PheWAS analysis suggests that
311 *PDE3B* PTVs have the strongest association with subcutaneous and lower-body adiposity (e.g.
312 hip and leg fat mass) (Supplementary Fig. S24). Therefore, understanding the fat depot-specific
313 metabolic effects of *PDE3B* may help uncover the mechanism underlying the positive
314 relationship of *PDE3B* PTVs with peripheral fat accumulation and favorable metabolic profiles.

315 Discussion

316 We developed DeGAs, an application of TSVD, to decompose genome-and phenome-wide
317 summary statistic matrix from association studies of thousands of phenotypes for systematic
318 characterization of latent components of genetic associations and advanced the understanding
319 on polygenic and pleiotropic architecture of complex traits. Applying DeGAs, we identified key
320 latent components characterized with disease outcomes, risk factors, comorbidity structures,

321 and environmental factors, with corresponding sets of genes and variants, providing insights on
322 their context specific functions. We demonstrated the robustness of the results by applying
323 DeGAs with different parameters. With additional biological characterization of latent
324 components using GREAT, we find component-specific enrichment of relevant phenotypes in
325 mouse phenotype ontology. This replication across species highlights the ability of DeGAs to
326 capture functionally relevant sets of both coding and non-coding variants in each component.

327 Our comparison of DeGAs to an alternative approach – decomposition of individual
328 phenotype data followed by GWAS – highlights the ability of DeGAs to curve out biomedically
329 relevant genetic signals as latent components. As an illustration of in-depth analysis of genetic
330 variants with different functional consequences, we reported applications of DeGAs for different
331 functional categories.

332 In DeGAs, we provided multiple ways to investigate the biological relevance of latent
333 components, including quantitative scores and ontology enrichment analysis. These metrics are
334 useful to annotate and interpret latent components, which are otherwise just mathematical
335 objects in a high-dimensional space. For example, we found a significant contribution of
336 anthropometric traits among the top 5 components, which may reflect the pervasive polygenicity
337 of these traits^{42,43} or phenotype selection in the UK Biobank study – anthropometric traits are
338 measured for most of the participants and their association signals are strong and stable. By
339 leveraging the ability of TSVD to efficiently summarize most of the variance in the input
340 association statistic matrix, DeGAs provides a systematic way to interpret polygenic and
341 pleiotropic genetic architecture of common complex traits.

342 Given that DeGAs is applied on summary statistics and does not require individual level
343 data, there is substantial potential to dissect genetic components of the human phenome when
344 applied to data from population-based biobanks around the globe^{14,44–47}. In fact, we are the first
345 to develop a computational method that can jointly analyze genetics of thousands of phenotypes
346 from a densely phenotyped population. As a proof of concept, we report novel potential
347 therapeutic targets against obesity or its complications based on combination of quantitative
348 results from DeGAs, phenome-wide analyses in the UK Biobank, and functional studies in
349 adipocytes. Due to the difference of phenotype and variant selection, it is possible that the latent
350 structure discovered from DeGAs can be different if one takes GWAS summary statistics from a
351 disparate GWAS study. However, DeGAs is capable of identifying the most relevant
352 components for a given input dataset using quantitative scores. In fact, our analysis for the three
353 datasets – “all”, coding, and PTVs – identified different PCs for each trait of our interest, but
354 their characterization with contribution scores showed consistent results.

355 Taken together, we highlight the directional concordance of our experimental data with
356 the quantitative results from DeGAs and PTV-phenotype associations: *GPR151* inhibition may
357 reduce total body and central fat, while deletion of *PDE3B* may favor subcutaneous, rather than
358 visceral, fat deposition; both are expected to have beneficial effects on cardiometabolic health.
359 Although these two genes were recently reported to be associated with obesity in another
360 recent study based on the UK Biobank⁴⁸, we are the first to experimentally identify *GPR151* as a
361 promising therapeutic target to treat obesity, partly due to its requisite role in regulating
362 adipogenesis. We also suggest *PDE3B* as a potential target of adipocyte-directed therapy. In
363 this study, we focused on evaluating the functional effects of these genes on adipocyte function
364 and development. We do not exclude the contribution nor the importance of other tissues or

365 mechanisms underlying body weight changes. Indeed, some lines of evidence support
366 additional effects of *GPR151* on obesity via the central nervous system – possibly on appetite
367 regulation³⁷.

368 The resource made available with this study, including the DeGAs app, an interactive
369 web application in the Global Biobank Engine⁴⁹, provides a starting point to investigate genetic
370 components, their functional relevance, and new therapeutic targets. These results highlight the
371 benefit of comprehensive phenotyping on a population and suggest that systematic
372 characterization and analysis of genetic associations across the human phenome will be an
373 important part of efforts to understand biology and develop novel therapeutic approaches.

374 Methods

375 Study population

376 The UK Biobank is a population-based cohort study collected from multiple sites across the
377 United Kingdom. Information on genotyping and quality control has previously been described¹⁴.
378 In brief, study participants were genotyped using two similar arrays (Applied Biosystems UK
379 BiLEVE Axiom Array (807,411 markers) and the UK Biobank Axiom Array (825,927 markers)),
380 which were designed for the UK Biobank study. The initial quality control was performed by the
381 UK Biobank analysis team and designed to accommodate the large-scale dataset of ethnically
382 diverse participants, genotyped in many batches, using two similar novel arrays¹⁴.

383 Genotype data preparation

384 We used genotype data from the UK Biobank dataset release version 2¹⁴ and the hg19 human
385 genome reference for all analyses in the study. To minimize the variabilities due to population
386 structure in our dataset, we restricted our analyses to include 337,199 White British individuals
387 based on the following five criteria reported by the UK Biobank in the file “ukb_sqc_v2.txt”:

- 388 1. self- reported white British ancestry (“in_white_British_ancestry_subset” column)
- 389 2. used to compute principal components (“used_in_pca_calculation” column)
- 390 3. not marked as outliers for heterozygosity and missing rates (“het_missing_outliers”
391 column)
- 392 4. do not show putative sex chromosome aneuploidy (“putative_sex_chromo-
393 some_aneuploidy” column)
- 394 5. have at most 10 putative third-degree relatives (“excess_relatives” column).

395
396 We annotated variants using the VEP LOFTEE plugin (<https://github.com/konradjk/loftee>) and
397 variant quality control by comparing allele frequencies in the UK Biobank and gnomAD
398 (gnomad.exomes.r2.0.1.sites.vcf.gz) as previously described¹².

399 We focused on variants outside of major histocompatibility complex (MHC) region
400 (chr6:25477797-36448354) and performed LD pruning using PLINK with “--indep 50 5 2”.
401 Furthermore, we selected variants according to the following rules:

- 402 • Missingness of the variant is less than 1%.

- 403 ● Minor-allele frequency is greater than 0.01%.
 - 404 ● The variant is in the LD-pruned set.
 - 405 ● Hardy-Weinberg disequilibrium test p-value is greater than 1.0×10^{-7} .
 - 406 ● Manual cluster plot inspection. We investigated cluster plots for subset of our variants
 - 407 and removed 11 variants that has unreliable genotype calls as previously described¹².
 - 408 ● Passed the comparison of minor allele frequency with gnomAD dataset as previously
 - 409 described¹².
- 410 These variant filters are summarized in Supplementary Fig. S1.

411 Phenotype data preparation

412 We organized 2,138 phenotypes from the UK Biobank in 11 distinct groups (Supplementary
413 Table 1). We included phenotypes with at least 100 cases for binary phenotypes and 100
414 individuals with non-missing values for quantitative phenotypes. For disease outcome
415 phenotypes, cancer, and family history, we used the same definitions as previously described¹².
416 We used specific data fields and data category from the UK Biobank to define the phenotypes in
417 the following categories as well as 19 and 42 additional miscellaneous binary and quantitative
418 phenotypes: medication, imaging, physical measurements, assays, and binary and quantitative
419 questionnaire (Supplementary Table 1-2).

420 Some phenotype information from the UK Biobank contains three instances, each of
421 which corresponds to (1) the initial assessment visit (2006-2010), (2) first repeat assessment
422 visit (2012-2013), and (3) imaging visit (2014-). For binary phenotype, we defined "case" if the
423 participants are classified as case in at least one of their visits and "control" otherwise. For
424 quantitative phenotype, we took a median of non-NA values. In total, we defined 1,196 binary
425 phenotypes and 943 quantitative phenotypes.

426 Genome-wide association analyses of 2,138 phenotypes

427 Association analyses for single variants were applied to the 2,138 phenotypes separately. For
428 binary phenotypes, we performed Firth-fallback logistic regression using PLINK v2.00a (17 July
429 2017) as previously described^{12,50}. For quantitative phenotypes, we applied generalized linear
430 model association analysis with PLINK v2.00a (20 Sep. 2017). We applied quantile
431 normalization for phenotype (--pheno-quantile-normalize option), where we fit a linear model
432 with covariates and transform the phenotypes to normal distribution $N(0, 1)$ while preserving the
433 original rank. We used the following covariates in our analysis: age, sex, types of genotyping
434 array, and the first four genotype principal components computed from the UK Biobank.

435 To test the effects of population stratification correction on the association analysis, we
436 performed additional GWAS with age, sex, types of array, and the first ten genotype principal
437 components as covariates for the five quantitative traits and five binary traits. For each pair of
438 GWAS summary statistics with four and ten genotype principal components, we computed the
439 genetic correlations and confirmed that the two GWAS run yielded the almost identical results
440 (Supplementary Table S10).

441 Summary statistic matrix construction and variant filters

442 We constructed three Z-score summary statistic matrices. Each element of the matrix
443 corresponds to summary statistic for a particular pair of a phenotype and a variant. We imposed
444 different sets of variant filters.

- 445 ● Variant quality control filter: Our quality control filter described in the previous section on
446 genotype data preparation.
- 447 ● Non-MHC variant filter: All variants outside of major histocompatibility complex region.
448 With this filter, variants in chr6:25477797-36448354 were excluded from the summary
449 statistic matrix.
- 450 ● Coding-only: With this filter, we subset to include only the variants having the VEP
451 LOFTEE predicted consequence of: missense, stop gain, frameshift, splice acceptor,
452 splice donor, splice region, loss of start, or loss of stop.
- 453 ● PTVs-only: With this filter, we subset to include only the variants having the VEP
454 LOFTEE predicted consequence of: stop gain, frameshift, splice acceptor, or splice
455 donor.

456 By combining these filters, we defined the following sets of variants

- 457 ● All-non-MHC: This is a combination of our variant QC filter and non-MHC filter.
- 458 ● Coding-non-MHC: This is a combination of our variant QC filter, non-MHC filter, and
459 Coding-only filter.
- 460 ● PTVs-non-MHC: This is a combination of our variant QC filter, non-MHC filter, and
461 PTVs-only filter.

462 In addition to phenotype quality control and variant filters, we introduced value-based filters
463 based on statistical significance to construct summary statistic matrices only with confident
464 values. We applied the following criteria for the value filter:

- 465 ● P-value of marginal association is less than 0.001.
- 466 ● Standard error of beta value or log odds ratio is less than 0.08 for quantitative
467 phenotypes and 0.2 for binary phenotypes.

468 With these filters, we obtained the following two matrices:

- 469 ● All-non-MHC dataset that contains 2,138 phenotypes and 235,907 variants. We label
470 this dataset as **“all” dataset**.
- 471 ● “Coding-non-MHC” dataset that contains phenotypes and 784 variants. We label this
472 dataset as **“Coding only” dataset**.
- 473 ● “PTVs-non-MHC” dataset that contains 628 phenotypes and 784 variants. We label this
474 dataset as **“PTVs only” dataset**.

475 The coding-only and PTVs-only datasets contain a fewer number of phenotypes because not all
476 the phenotypes have statistically significant associations with coding variants or PTVs. The
477 effects of variant filters are summarized in Fig. S1. Finally, we transformed the summary
478 statistics to Z-scores so that each vector that corresponds to a particular phenotype has zero
479 mean with unit variance.

480 Truncated singular value decomposition of the summary statistic 481 matrix

482 For each summary statistic matrix, we applied truncated singular value decomposition (TSVD).
483 The matrix, which we denote as W , of size $N \times M$, where N denotes the number of phenotypes
484 and M denotes the number of variants, is the input data. With TSVD, W is factorized into a
485 product of three matrices: U , S , and V^T : $W = USV^T$, where $U = (u_{i,k})_{i,k}$ is an orthonormal
486 matrix of size $N \times K$ whose columns are phenotype (left) singular vectors, S is a diagonal matrix
487 of size $K \times K$ whose elements are singular values, and $V = (v_{j,k})_{j,k}$ is an orthonormal matrix of
488 size $M \times K$ whose columns are variant (right) singular vectors. While singular values in S
489 represent the magnitude of the components, singular vectors in U and V summarizes the
490 strength of association between phenotype and component and variant and component,
491 respectively. With this decomposition, the k -th latent component (principal component, PC k)
492 are represented as a product of k -th column of U , k -th diagonal element in S , and k -th row of
493 V^T . For TSVD on the summary statistics, we used implicitly restarted Lanczos bidiagonalization
494 algorithm (IRLBA)⁵¹ (<https://github.com/bwlewis/irlba>) implemented on SciDB⁵² to compute the
495 first K components in this decomposition.

496 Relative variance explained by each of the components

497 A scree plot (Fig. S1) quantify the variance explained by each component: variance explained
498 by k -th component = $s_k^2 / \text{Var}_{\text{Tot}}(W)$ where, s_k is the k -th diagonal element in the diagonal matrix
499 S and $\text{Var}_{\text{Tot}}(W)$ is the total variance of the original matrix before DeGAs is applied.

500 Selection of number of latent components in TSVD

501 In order to apply TSVD to the input matrix, the number of components should be specified. We
502 apply $K = 100$ for our analysis for all of the datasets. Following a standard practice of keeping
503 components with eigenvalues greater than the average²³, we first computed the expected value
504 of squared eigenvalues under the null model where the distribution of variance explained scores
505 across the full-ranks are uniform. This can be computed with the rank of the original matrix,
506 which is equal to the number of phenotypes in our datasets:

$$507 \quad E[\text{Variance explained by } k\text{-th component under the null}] = 1/(\text{Rank}(W)^2)$$

508 We then compared the eigenvalues characterized from TSVD with the expected value. For all of
509 the three datasets, we found that that of 100-th component is greater than the expectation. This
510 indicates even the 100-th components are informative to represent the variance of the original
511 matrix. In the interest of computational efficiency, we set $K = 100$.

512 To demonstrate the robustness of the DeGAs components with respect to the number of
513 latent components (K), we performed additional analyses with $K = 90$ and $K = 110$, and
514 investigated the first five latent components as well as the top three components for the three
515 phenotypes of our interest.

516 Factor scores

517 From these decomposed matrices, we computed **factor score** matrices for both phenotypes
518 and variants as the product of singular vector matrix and singular values. We denote the one for
519 phenotypes as $F_p = (f_{i,j}^p)_{i,j}$ the one for variants as $F_v = (f_{i,j}^v)_{i,j}$ and defined as follows:

$$520 \quad F_p = US$$
$$521 \quad F_v = VS$$

522 Since these factor scores are mathematically the same as principal components in principal
523 component analysis (PCA), one can investigate the contribution of the phenotypes or variants
524 for specific principal components by simply plotting factor scores²³ (Fig. 2a-b). Specifically,
525 phenotype factor score is the same as phenotype principal components and variant factor score
526 is the same as variant principal components. By normalizing these factor scores, one can
527 compute contribution scores and cosine scores to quantify the importance of phenotypes,
528 variants, and principal components as described below.

529 Scatter plot visualization with biplot annotations

530 To investigate the relationship between phenotype and variants in the TSVD eigenspace, we
531 used a variant of biplot visualization^{53,54}. Specifically, we display phenotypes projected on
532 phenotype principal components ($F_p = US$) as a scatter plot. We also show variants projected on
533 variant principal components ($F_v = VS$) as a separate scatter plot and added phenotype singular
534 vectors (U) as arrows on the plot using sub-axes (Fig. 2b, 4c, S5-6). In scatter plot with biplot
535 annotation, the inner product of a genetic variant and a phenotype represents the direction and
536 the strength of the projection of the genetic association of the variant-phenotype pair on the
537 displayed latent components. For example, when a variant and a phenotype share the same
538 direction on the annotated scatter plot, that means the projection of the genetic associations of
539 the variant-phenotype pair on the displayed latent components is positive. When a variant-
540 phenotype pair is projected on the same line, but on the opposite direction, the projection of the
541 genetic associations on the shown latent components is negative. When the variant and
542 phenotype vectors are orthogonal or one of the vectors are of zero length, the projection of the
543 genetic associations of the variant-phenotype pair on the displayed latent components is zero.
544 Given the high dimensionality of the input summary statistic matrix, we selected relevant
545 phenotypes to display to help interpretation of genetic associations in the context of these traits.

546 Contribution scores

547 To quantify the contribution of the phenotypes, variants, and genes to a given component, we
548 compute **contribution scores**. We first define **phenotype contribution score** and **variant**
549 **contribution score**. We denote phenotype contribution score and variant contribution score for
550 some component k as $\text{cntr}_k^{\text{phe}}(i)$ and $\text{cntr}_k^{\text{var}}(j)$, respectively. They are defined by squaring the
551 left and right singular vectors and normalizing them by Euclidian norm across phenotypes and
552 variants:

$$553 \quad \text{cntr}_k^{\text{phe}}(i) = (u_{i,k})^2$$

554
$$\text{cntr}_k^{\text{var}}(j) = (v_{i,k})^2$$

555 where, i and j denote indices for phenotype and variant, respectively. Because U and V are
556 orthonormal, the sum of phenotype and variant contribution scores for a given component are
557 guaranteed to be one, i.e. $\sum_i \text{cntr}_k^{\text{phe}}(i) = \sum_j \text{cntr}_k^{\text{var}}(j) = 1$.

558 Based on the variant contribution scores for the k -th component, we define the **gene**
559 **contribution score** for some component k as the sum of variant contribution scores for the set
560 of variants in the gene:

561
$$\text{cntr}_k^{\text{gene}}(g) = \sum_{j \in g} \text{cntr}_k^{\text{var}}(j)$$

562 where, g denotes indices for the set of variants in gene g . To guarantee that gene contribution
563 scores for a given component sum up to one, we treat the variant contribution score for the non-
564 coding variants as gene contribution scores. When multiple genes, g_1, g_2, \dots, g_n are sharing the
565 same variants, we defined the gene contribution score for the union of multiple genes rather
566 than each gene:

567
$$\text{cntr}_k^{\text{gene}}(\{g_i \mid i \in [1, n]\}) = \sum_{\{j \mid j \in g_1 \wedge j \in g_2 \wedge \dots \wedge j \in g_n\}} \text{cntr}_k^{\text{var}}(j)$$

568 With these contribution score for a given component, it is possible to quantify the relative
569 importance of a phenotype, variant, or gene to the component. Since DeGAs identifies latent
570 components using unsupervised learning, we interpret each component in terms of the driving
571 phenotypes, variants, and genes, i.e. the ones with large contribution scores for the component.

572 The top 20 driving phenotypes, variants, and genes (based on contribution scores) for
573 the top five TSVD components and the top three key components for our phenotypes of interest
574 are summarized in Supplementary Table S3.

575 We used stacked bar plots for visualization of the contribution score profile for each of
576 the components. We represent phenotypes, genes, or variants with large contribution scores as
577 colored segments and aggregated contributions from the remaining ones as “others” in the plot
578 (Fig. 2c-d, 3a, 4a-b, Supplementary Fig. S4). To help interpretation of the major contributing
579 factors for the key components, we grouped phenotypes into categories, such as “fat”, “fat-free”
580 phenotypes, and showed the sum of contribution scores for the phenotype groups. The list of
581 phenotype groups used in the visualization is summarized in Supplementary Table S3.

582 Squared cosine scores

583 Conversely, we can also quantify the relative importance of the latent components for a given
584 phenotype or variant with **squared cosine scores**. We denote phenotype squared cosine score
585 for some phenotype i and variant squared cosine score for some variant j as $\cos^2_i^{\text{phe}}(k)$ and
586 $\cos^2_j^{\text{var}}(k)$, respectively. They are defined by squaring of the factor scores and normalizing
587 them by Euclidian norm across components:

588
$$\cos^2_i^{\text{phe}}(k) = \frac{(f_{i,k}^p)^2}{\sum_{k'} (f_{i,k'}^p)^2}$$

589
$$\cos^2_j^{\text{var}}(k) = \frac{(f_{j,k}^v)^2}{\sum_{k'} (f_{j,k'}^v)^2}$$

590 By definition, the sum of squared cosine scores across a latent component for a given
591 phenotype or variant equals to one, i.e. $\sum_k \cos^2_i^{\text{phe}}(k) = \sum_k \cos^2_j^{\text{var}}(k) = 1$. While singular
592 values in the diagonal matrix S quantify the importance of latent components for the global latent
593 structure, the phenotype or variant squared cosine score quantifies the relative importance of
594 each component in the context of a given phenotype or a variant. The squared cosine scores for
595 the phenotypes highlighted in the study is summarized in Fig. S3 and Supplementary Fig. S9.

596 Note that squared cosine scores and contribution scores are two complementary scoring
597 metrics to quantify the relationship among phenotypes, components, variants, and genes. It
598 does not necessarily have inverse mapping property. For example, it is possible to see a
599 situation, where for a given phenotype p , phenotype squared cosine score identifies k as the top
600 key component, but phenotype contribution score for k identifies p' ($p' \neq p$) as the top driving
601 phenotype for the component k . This is because the two scores, contribution score and squared
602 cosine score, are both defined by normalizing singular vector and principal component vector
603 matrices, respectively, but with respect to different slices: one for row and the other for column.

604 TSVD of the individual-level phenotypes

605 To characterize the latent components in the raw phenotype data, we first applied median
606 imputation for missing values on the phenotype data followed by Z-score transformation. Using
607 Python scikit-learn package⁵⁵, we applied TSVD on the imputed and normalized phenotype
608 matrix and characterized the first five latent components and visualized the scree plot as well as
609 the phenotype and individual PCs in scatter plots.

610 Genome wide-association analysis for phenotype PCs

611 Using the results of the phenotype decomposition described above, we defined principal
612 components of the individual's phenotype (phenotype PCs) and applied genome-wide
613 association analysis using the same procedure we used for the original quantitative traits. We
614 used R package qqplot to generate Manhattan plot⁵⁶.

615 Genetic correlation of phenotype PCs

616 To compare the results of association analysis of phenotype PCs, we computed genetic
617 correlation using LD score regression⁵⁷. We summarized the estimated genetic correlation (r_g)
618 as heatmap and characterized the median value of absolute value of r_g among the top k
619 phenotype PCs as a function of k .

620 Genomic region enrichment analysis with GREAT

621 We applied the genomic region enrichment analysis tool (GREAT version 4.0.3) to each DeGAs
622 components²⁰. We used the mouse genome informatics (MGI) phenotype ontology, which
623 contains manually curated knowledge about hierarchical structure of phenotypes and genotype-
624 phenotype mapping of mouse³². We downloaded their ontologies on 2017-09-28 and mapped

625 MGI gene identifiers to Ensembl human gene ID through unambiguous one-to-one homology
626 mapping between human and mouse Ensembl IDs. We removed ontology terms that were
627 labelled as “obsolete”, “bad”, or “unknown” from our analysis. As a result, we obtained 709,451
628 mapping annotation spanning between 9,554 human genes and 9,592 mouse phenotypes.

629 For each DeGAs component, we selected the top 5,000 variants according to their
630 variant contribution score and performed enrichment analysis with the default parameter as
631 described elsewhere²⁰. Since we included the non-coding variants in the analysis, we focused
632 on GREAT binomial genomic region enrichment analysis based on the size of regulatory
633 domain of genes and quantified the significance of enrichment in terms of binomial fold
634 enrichment and binomial p-value. Given that we have 9,561 terms in the ontology, we set a
635 Bonferroni-corrected p-value threshold of 5×10^{-6} .

636 To illustrate the results of the genomic region enrichment analysis for the phenotypes of
637 our interest, we made circular bar plots using the R package ggplot2, where each of the key
638 components are displayed in the innermost track with their phenotype squared cosine score to
639 be proportional to their angle, and the resulted significant ontology terms are represented as the
640 bars. To focus on the significant signals with large effect size, we imposed additional filter of
641 binomial fold ≥ 2.0 and binomial p-value threshold of 5×10^{-7} . The binomial fold change is
642 represented as the radius and the binomial p-value is represented as color gradient in a log
643 scale in the plot (Fig. 3b, Supplementary Fig. S7-8, Supplementary Table S5-7).

644 Specificity analysis of GREAT enrichment

645 To test the specificity of the GREAT enrichment of each of the 100 DeGAs components, we
646 computed Jaccard index similarity scores. For each DeGAs latent component, we looked at the
647 GREAT enrichment and took the top five enriched terms sorted by GREAT binomial fold. To
648 measure the similarity between these enriched terms, we identified the set of genes annotated
649 for those terms and computed Jaccard index defined below:

$$650 \quad \text{Similarity}(\text{term set}_A, \text{term set}_B) = \frac{|\text{Gene set}(\text{term set}_A) \cap \text{Gene set}(\text{term set}_B)|}{|\text{Gene set}(\text{term set}_A) \cup \text{Gene set}(\text{term set}_B)|}$$

651 where,

$$652 \quad \text{Gene set}(\text{term set}_A) = \bigcup_{t \in A} \text{Gene set}(\text{term}_t)$$

653 and $\text{Gene set}(\text{term}_t)$ indicates set of genes annotated with term t . We computed all the pair-wise
654 similarity across the top k DeGAs components and summarized their median as a function of k .

655 Quality control of variant calling with intensity plots

656 To investigate the quality of variant calling for the two PTVs highlighted in the study, we
657 manually inspected intensity plots. These plots are available on Global Biobank Engine.

- 658
- 659 • <https://biobankengine.stanford.edu/intensity/rs114285050>
 - <https://biobankengine.stanford.edu/intensity/rs150090666>

660 Phenome-wide association analysis

661 To explore the functional roles of the two PTVs across thousands of potentially correlated
662 phenotypes, we performed a phenome-wide association study (PheWAS). We report the
663 statistically significant ($p < 0.001$) associations with phenotypes with at least 1,000 case count
664 (binary phenotypes) or 1,000 individuals with measurements with non-missing values
665 (quantitative phenotypes) (Fig. 3d, Supplementary Fig. S10). The results of this PheWAS are
666 also available as interactive plots as a part of Global Biobank Engine.

- 667 • <https://biobankengine.stanford.edu/variant/5-145895394>
- 668 • <https://biobankengine.stanford.edu/variant/11-14865399>

669 Univariate regression analysis for the identified PTVs

670 To quantify the effects of the two PTVs on obesity, we performed univariate regression analysis.
671 We extracted individual-level genotype information for the two PTVs with the PLINK2 pgen
672 Python API (<http://www.cog-genomics.org/plink/2.0/>)⁵⁰. After removing individuals with missing
673 values for BMI and genotype, we performed linear regression for BMI
674 (<http://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=21001>) with age, sex, and the first four genomic
675 PCs as covariates:

676
$$\text{BMI} \sim 0 + \text{age} + \text{as.factor}(\text{sex}) + \text{PC1} + \text{PC2} + \text{PC3} + \text{PC4} + \text{as.factor}(\text{PTV})$$

677 where, PC1-4 denotes the first four components of genomic principal components, PTV ranges
678 in 0, 1, or 2 and it indicates the number of minor alleles that the individuals have.

679 Mouse 3T3-L1 cell culture and differentiation

680 3T3-L1 preadipocytes were cultured in Dulbecco's Modified Eagle's Medium (DMEM) containing
681 10% fetal bovine serum (FBS) and antibiotics (100 U/mL of penicillin G and 100 $\mu\text{g}/\text{mL}$ of
682 streptomycin) at 37°C in a humidified atmosphere containing 5% CO₂. To obtain fully
683 differentiated adipocytes, 3T3-L1 preadipocytes were grown into 2-day post-confluence, and
684 then differentiation was induced by using a standard differentiation cocktail containing 0.5 mM of
685 IBMX, 1 μM of dexamethasone, 1 $\mu\text{g}/\text{mL}$ of insulin, and 10% FBS. After 48 h, medium was
686 changed into DMEM supplemented with 10% FBS and 1 $\mu\text{g}/\text{mL}$ of insulin and replenished every
687 48 h for an additional 6 days.

688 Human SGBS cell culture and differentiation

689 SGBS cells were cultured in DMEM/F12 containing 33 μM biotin, 17 μM pantothenate,
690 0.1 mg/mg streptomycin and 100 U/mL penicillin (0F medium) supplemented with 10% FBS in a
691 5% CO₂ incubator. To initiate differentiation, confluent cells were stimulated by 0F media
692 supplemented with 0.01 mg/mL human transferrin, 0.2 nM T3, 100 nM cortisol, 20 nM insulin,
693 250 μM IBMX, 25 nM dexamethasone and 2 μM rosiglitazone. After day 4, the differentiating
694 cells were kept in 0F media supplemented with 0.01 mg/mL human transferrin, 100 nM cortisol,
695 20 nM insulin and 0.2 nM T3 for additional 8-10 days until cells were fully differentiated.

696 siRNA knockdown in 3T3-L1 preadipocytes

697 At 80% confluence, 3T3-L1 preadipocytes were transfected with 50 nM siRNA against
698 *Gpr151* (Origene #SR412988), *Pde3b* (Origene #SR422062), or scrambled negative control
699 (Origene #SR30004) using Lipofectamine™ RNAiMAX Transfection Reagent (Invitrogen)
700 following the manufacturer's protocol. The transfected cells were incubated for 48 h and then
701 subjected to differentiation.

702 Reverse transcription (RT) and qPCR analysis

703 Total RNA was extracted using TRIzol reagent (Invitrogen), following the manufacturer's
704 instruction. RNA was converted to cDNA using High-Capacity cDNA Reverse Transcription Kit
705 (Applied Biosystems). Quantitative PCR reactions were prepared with TaqMan™ Fast
706 Advanced Master Mix (Thermo Fisher Scientific) and performed on ViiA 7 Real-Time PCR
707 System (Thermo Fisher Scientific). All data were normalized to the content of Cyclophilin A
708 (PPIA), as the endogenous control. TaqMan primer information for RT-qPCR is listed below:
709 *GPR151* (Hs00972208_s1), *Gpr151* (Mm00808987_s1), *PDE3B* (Hs00265322_m1), *Pde3b*
710 (Mm00691635_m1), *Pparg* (Mm00440940_m1), *Cebpa* (Mm00514283_s1), *Fabp4*
711 (Mm00445878_m1), *PPIA* (Hs04194521_s1), *Ppia* (Mm02342430_g1).

712 Oil Red O staining and quantification

713 Cells were washed twice with PBS and fixed with 10% formalin for 1 h at room temperature.
714 Cells were then washed with 60% isopropanol and stained for 15 min with a filtered Oil Red O
715 solution (mix six parts of 0.35% Oil Red O in isopropanol with four parts of water). After washing
716 with PBS 4 times, cells were maintained in PBS and visualized by inverted microscope. After
717 taking pictures, Oil Red O stain was extracted with 100% isopropanol and the absorbance was
718 measured at 492 nm by a multi-well spectrophotometer (Bio-Rad).

719 Lipolysis assay

720 Glycerol release into the culture medium was used as an index of lipolysis. Fully differentiated
721 3T3-L1 adipocytes were serum starved overnight and then treated with either vehicle (DMSO)
722 or the lipolytic stimuli isoproterenol (ISO, 10 μ M) for 3 h. The culture medium was collected and
723 the glycerol content in the culture medium was measured using an adipocyte lipolysis assay kit
724 (ZenBio #LIP-1-NCL1). Glycerol release into the culture medium was normalized to the protein
725 content of the cells from the same plate.

726 Overexpression of *GPR151* in 3T3-L1 preadipocytes

727 The *GPR151* construct was obtained from Addgene (#66327). This construct includes a
728 cleavable HA signal to promote membrane localization, a FLAG epitope sequence for cell
729 surface staining followed by codon-optimized human *GPR151* sequence⁵⁸. We PCR-amplified
730 the above sequence with stop codon and assembled it into a lentiviral plasmid (Addgene

731 #85969) with either EF1 α promoter (Addgene # 11154) or aP2 promoter (Addgene # 11424).
732 EF1 α -*GPR151* or aP2-*GPR151* lentiviral plasmid were transfected into human embryonic
733 kidney 293T cells, together with the viral packaging vectors pCMV-dR8.91 and pMD2-G. 72 h
734 after transfection, virus-containing medium was collected, filtered through a 0.45- μ m pore-size
735 syringe filter, and frozen at -80°C. 3T3-L1 preadipocytes at 50% confluence were infected with
736 the lentivirus stocks containing 8 μ g/mL polybrene. Two days after transduction, lentivirus-
737 infected 3T3-L1 preadipocytes were subject to differentiation.

738 Flow cytometry analysis

739 Day 6 differentiating 3T3-L1 adipocytes were collected and washed with ice cold FACS buffer
740 (PBS containing 2% BSA). Cells were first resuspended into FACS staining buffer (BioLegend #
741 420201) at ~1M cells/100 μ l and incubated with anti-mouse CD16/CD32 Fc Block (BioLegend #
742 101319) at room temperature for 10-15 min. Cells were then incubated with APC-conjugated
743 FLAG antibody (BioLegend # 637307) for 20-30 min at room temperature in the dark. Following
744 washing and centrifugation, cells were resuspended in FACS buffer and sorted using a BD
745 InfluxTM Cell Sorter. Cells without FLAG antibody staining were used to determine background
746 fluorescence levels. Cells were sorted based on APC fluorescence and collected directly into
747 TRIzol reagent for RNA extraction.

748 Western Blot Analysis

749 Lysate aliquots containing 50 μ g of proteins were denatured, separated on a 4-10% SDS-
750 polyacrylamide gel, and transferred to nitrocellulose membranes using a Trans-Blot[®] SD Semi-
751 Dry Transfer Cell (Bio-Rad). Membranes were blocked in 5% non-fat milk and incubated
752 overnight at 4 °C with primary antibodies: anti-GPR151 (LSBio # LS-B6760-50) or anti-beta-
753 actin (Cell Signaling #3700). Subsequently, the membranes were incubated for 1 h at room
754 temperature with IRDye[®] 800CW goat-anti-mouse antibody (LI-COR #926-32210). Target
755 proteins were visualized using Odyssey[®] Fc Imaging System (LI-COR).

756 Statistical analysis of functional data

757 Data are expressed as mean \pm SEM. Student's t test was used for single variables, and one-
758 way ANOVA with Bonferroni post hoc correction was used for multiple comparisons using
759 GraphPad Prism 7 software.

760 Acknowledgements

761 This research has been conducted using the UK Biobank Resource under Application Number
762 24983, "Generating effective therapeutic hypotheses from genomic and hospital linkage data"
763 (<http://www.ukbiobank.ac.uk/wp-content/uploads/2017/06/24983-Dr-Manuel-Rivas.pdf>).
764 We thank all the participants in the UK Biobank study. We thank Nasa Sinnott-Armstrong for the
765 insightful discussion. We thank Robert Tibshirani, the members of the Rivas lab, the Ingelsson

766 lab, and the Bejerano lab for helpful feedback. This work was supported by National Human
767 Genome Research Institute (NHGRI) and National Institute of Diabetes and Digestive and
768 Kidney Diseases (NIDDK) of the National Institutes of Health (NIH) under awards
769 R01HG010140 and R01DK106236. The content is solely the responsibility of the authors and
770 does not necessarily represent the official views of the National Institutes of Health. Y.T. is
771 supported by Funai Overseas Scholarship from Funai Foundation for Information Technology
772 and the Stanford University School of Medicine. J.M.J. was funded by grant NNF17OC0025806
773 from the Novo Nordisk Foundation and the Stanford Bio-X Program. M.A.R. and C.D. are
774 supported by Stanford University and a National Institute of Health center for Multi- and Trans-
775 ethnic Mapping of Mendelian and Complex Diseases grant (5U01 HG009080). C.D. is
776 supported by a postdoctoral fellowship from the Stanford Center for Computational,
777 Evolutionary, and Human Genomics and the Stanford ChEM-H Institute. We obtained clip-arts
778 for Fig. 1b from Irasutoya (<https://www.irasutoya.com/>) by following their terms and conditions.
779 The copyright of the original clip-arts belongs to Mr. Takashi Mifune. We would like to thank the
780 Customer Solutions Team from Paradigm4 who helped us implement efficient databases for
781 queries and application of inference methods to the data, and also implemented optimized
782 versions of truncated singular value decomposition.
783

784 Author information

785 Author contributions

786 M.A.R. and E.I. conceived and designed the study. Y.T. and M.A.R. carried out the statistical
787 and computational analyses with advice from J.M.J., H.H., M.A., C.D., B.N., K.L., T.H., G.B., and
788 E.I. J.L., C.Y.P., and E.I. carried out the functional experiments. Y.T., M.A., and C.D. carried out
789 quality control of the data. C.C. optimized and implemented computational methods. Y.T. and
790 M.A.R. developed the DeGAs app in Global Biobank Engine. M.A.R. supervised computational
791 and statistical aspects of the study. E.I. supervised experimental aspects of the study. The
792 manuscript was written by Y.T., J.L., J.M.J., E.I., and M.A.R.; and revised by all the co-authors.
793 All co-authors have approved of the final version of the manuscript.
794

795 Competing financial interests

796 None.
797

798 Data availability:

799 Data is displayed in the Global Biobank Engine (<https://biobankengine.stanford.edu>). Analysis
800 scripts and notebooks are available on GitHub at <https://github.com/rivas-lab/public-resources>.

801 References

- 802 1. Welter, D. *et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait associations.
803 *Nucleic Acids Res.* **42**, 1001–1006 (2014).

- 804 2. Visscher, P. M. *et al.* 10 Years of GWAS Discovery: Biology, Function, and Translation. *Am.*
805 *J. Hum. Genet.* **101**, 5–22 (2017).
- 806 3. The International Schizophrenia Consortium. Common polygenic variation contributes to risk
807 of schizophrenia and bipolar disorder. *Nature* **460**, 748–752 (2009).
- 808 4. Purcell, S. M. *et al.* A polygenic burden of rare disruptive mutations in schizophrenia. *Nature*
809 **506**, 185–190 (2014).
- 810 5. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An Expanded View of Complex Traits: From
811 Polygenic to Omnigenic. *Cell* **169**, 1177–1186 (2017).
- 812 6. Plenge, R. M., Scolnick, E. M. & Altshuler, D. Validating therapeutic targets through human
813 genetics. *Nat. Rev. Drug Discov.* **12**, 581–594 (2013).
- 814 7. Waring, R. Sequence Variations in PCSK9, Low LDL, and Protection against Coronary
815 Heart Disease. *N Engl J Med* **9** (2006).
- 816 8. Cohen, J. *et al.* Low LDL cholesterol in individuals of African descent resulting from frequent
817 nonsense mutations in PCSK9. *Nat. Genet.* **37**, 161–165 (2005).
- 818 9. Rivas, M. A. *et al.* Deep resequencing of GWAS loci identifies independent rare variants
819 associated with inflammatory bowel disease. *Nat. Genet.* **43**, 1066–1073 (2011).
- 820 10. Rivas, M. A. *et al.* A protein-truncating R179X variant in RNF186 confers protection against
821 ulcerative colitis. *Nat. Commun.* **7**, 12342 (2016).
- 822 11. Nejentsev, S., Walker, N., Riches, D., Egholm, M. & Todd, J. A. Rare Variants of IFIH1, a
823 Gene Implicated in Antiviral Responses, Protect Against Type 1 Diabetes. *Science* **324**,
824 387–389 (2009).
- 825 12. DeBoever, C. *et al.* Medical relevance of protein-truncating variants across 337,205
826 individuals in the UK Biobank study. *Nat. Commun.* **9**, 1612 (2018).
- 827 13. Tipney, H. *et al.* The support of human genetic evidence for approved drug indications. *Nat.*
828 *Genet.* **47**, 1–7 (2015).

- 829 14. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
830 *Nature* **562**, 203–209 (2018).
- 831 15. Stewart, G. On the Early History of the Singular Value Decomposition. *SIAM Rev.* **35**, 551–
832 566 (1993).
- 833 16. Novembre, J. *et al.* Genes mirror geography within Europe. *Nature* **456**, 98–101 (2008).
- 834 17. Varah, J. M. On the Numerical Solution of Ill-Conditioned Linear Systems with Applications
835 to Ill-Posed Problems. *SIAM J. Numer. Anal.* **10**, 257–267 (1973).
- 836 18. Hanson, R. J. A Numerical Method for Solving Fredholm Integral Equations of the First Kind
837 Using Singular Values. *SIAM J. Numer. Anal.* **8**, 616–622 (1971).
- 838 19. Hansen, P. C. The truncatedSVD as a method for regularization. *BIT Numer. Math.* **27**,
839 534–553 (1987).
- 840 20. McLean, C. Y. *et al.* GREAT improves functional interpretation of cis-regulatory regions.
841 *Nat. Biotechnol.* **28**, 495–501 (2010).
- 842 21. Rivas, M. A. *et al.* Effect of predicted protein-truncating genetic variants on the human
843 transcriptome. *Science* **348**, 666–669 (2015).
- 844 22. MacArthur, D. G. *et al.* A Systematic Survey of Loss-of-Function Variants in Human Protein-
845 Coding Genes. *Science* **335**, 823–828 (2012).
- 846 23. Abdi, H. & Williams, L. J. Principal component analysis. *Wiley Interdiscip. Rev. Comput.*
847 *Stat.* **2**, 433–459 (2010).
- 848 24. Bigaard, J. *et al.* Body Fat and Fat-Free Mass and All-Cause Mortality. *Obes. Res.* **12**,
849 1042–1049 (2004).
- 850 25. Foster, K. R. & Lukaski, H. C. Whole-body impedance--what does it measure? *Am. J. Clin.*
851 *Nutr.* **64**, 388S-396S (1996).
- 852 26. Talma, H. *et al.* Bioelectrical impedance analysis to estimate body composition in children
853 and adolescents: a systematic review and evidence appraisal of validity, responsiveness,

- 854 reliability and measurement error. *Obes. Rev. Off. J. Int. Assoc. Study Obes.* **14**, 895–905
855 (2013).
- 856 27. Moran, A. E. *et al.* The Global Burden of Ischemic Heart Disease in 1990 and 2010.
857 *Circulation* (2014).
- 858 28. Lammert, F. *et al.* Gallstones. *Nat. Rev. Dis. Primer* **2**, 16024 (2016).
- 859 29. Zhou, D., Wu, J. & Luo, G. Body mass index and risk of non-melanoma skin cancer:
860 cumulative evidence from prospective studies. *Sci. Rep.* **6**, 37691 (2016).
- 861 30. Khera, A. V. & Kathiresan, S. Is Coronary Atherosclerosis One Disease or Many?: Setting
862 Realistic Expectations for Precision Medicine. *Circulation* **135**, 1005–1007 (2017).
- 863 31. Bennion, L. J. & Grundy, S. M. Risk Factors for the Development of Cholelithiasis in Man. *N.*
864 *Engl. J. Med.* **299**, 1161–1167 (1978).
- 865 32. Smith, C. L. & Eppig, J. T. Expanding the mammalian phenotype ontology to support
866 automated exchange of high throughput mouse phenotyping data generated by large-scale
867 mouse knockout screens. *J. Biomed. Semant.* **6**, 11 (2015).
- 868 33. Krawczyk, M. *et al.* Phytosterol and cholesterol precursor levels indicate increased
869 cholesterol excretion and biosynthesis in gallstone disease. *Hepatology* **55**, 1507–1517
870 (2012).
- 871 34. Abul-Husn, N. S. *et al.* A Protein-Truncating HSD17B13 Variant and Protection from Chronic
872 Liver Disease. *N. Engl. J. Med.* **378**, 1096–1106 (2018).
- 873 35. Green, H. & Kehinde, O. An established preadipose cell line and its differentiation in culture
874 II. Factors affecting the adipose conversion. *Cell* **5**, 19–27 (1975).
- 875 36. Wabitsch, M. *et al.* Characterization of a human preadipocyte cell strain with high capacity
876 for adipose differentiation. *Int. J. Obes.* **25**, 8–15 (2001).
- 877 37. Broms, J. *et al.* Monosynaptic retrograde tracing of neurons expressing the G-protein
878 coupled receptor Gpr151 in the mouse brain. *J. Comp. Neurol.* **525**, 3227–3250 (2017).

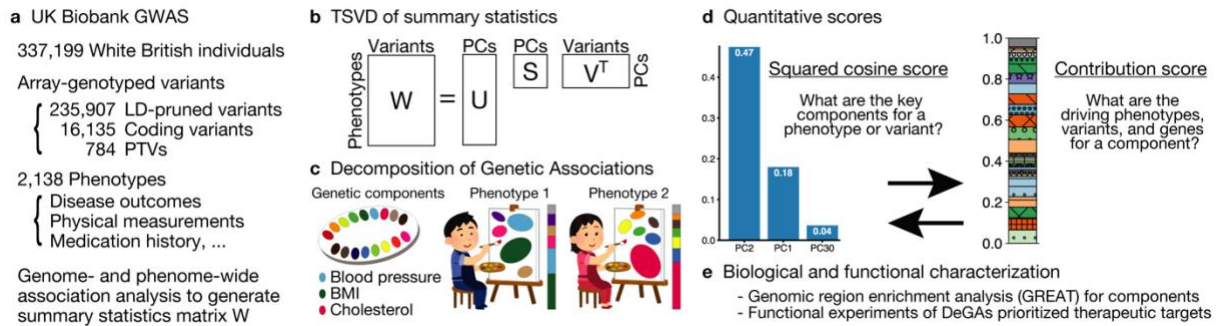
- 879 38. Shakur, Y. *et al.* Regulation and function of the cyclic nucleotide phosphodiesterase (PDE3)
880 gene family. in *Progress in Nucleic Acid Research and Molecular Biology* **66**, 241–277
881 (Elsevier, 2000).
- 882 39. Chung, Y. W. *et al.* White to beige conversion in PDE3B KO adipose tissue through
883 activation of AMPK signaling and mitochondrial function. *Sci. Rep.* **7**, 40445 (2017).
- 884 40. Chung, Y. W. *et al.* Targeted disruption of PDE3B, but not PDE3A, protects murine heart
885 from ischemia/reperfusion injury. *Proc. Natl. Acad. Sci.* 201416230 (2015).
886 doi:10.1073/pnas.1416230112
- 887 41. Emdin, C. A. *et al.* Genetic Association of Waist-to-Hip Ratio With Cardiometabolic Traits,
888 Type 2 Diabetes, and Coronary Heart Disease. *JAMA* **317**, 626–634 (2017).
- 889 42. Locke, A. E. *et al.* Genetic studies of body mass index yield new insights for obesity biology.
890 *Nature* **518**, 197–206 (2015).
- 891 43. Yengo, L. *et al.* Meta-analysis of genome-wide association studies for height and body mass
892 index in ~700,000 individuals of European ancestry. *bioRxiv* 274654 (2018).
893 doi:10.1101/274654
- 894 44. Collins, F. S. & Varmus, H. A New Initiative on Precision Medicine. *N. Engl. J. Med.* **372**,
895 793–795 (2015).
- 896 45. Chen, Z. *et al.* China Kadoorie Biobank of 0.5 million people: survey methods, baseline
897 characteristics and long-term follow-up. *Int. J. Epidemiol.* **40**, 1652–1666 (2011).
- 898 46. Christensen, H., Nielsen, J. S., Sorensen, K. M., Melbye & Brandslund, I. New national
899 Biobank of The Danish Center for Strategic Research on Type 2 Diabetes (DD2). *Clin.*
900 *Epidemiol.* **37** (2012). doi:10.2147/CLEP.S33042
- 901 47. Avlund, K. *et al.* Copenhagen Aging and Midlife Biobank (CAMB): An Introduction. *J. Aging*
902 *Health* **26**, 5–20 (2014).
- 903 48. Emdin, C. A. *et al.* Analysis of predicted loss-of-function variants in UK Biobank identifies
904 variants protective for disease. *Nat. Commun.* **9**, 1613 (2018).

- 905 49. McInnes, G. *et al.* Global Biobank Engine: enabling genotype-phenotype browsing for
906 biobank summary statistics. *Bioinformatics* (2019). doi:10.1093/bioinformatics/bty999
- 907 50. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer
908 datasets. *GigaScience* **4**, 7 (2015).
- 909 51. Baglama, J. & Reichel, L. Augmented Implicitly Restarted Lanczos Bidiagonalization
910 Methods. *SIAM J. Sci. Comput.* **27**, 19–42 (2005).
- 911 52. Brown, P. G. Overview of sciDB: large scale array storage, processing and analysis. in
912 *Proceedings of the 2010 international conference on Management of data - SIGMOD '10*
913 963 (ACM Press, 2010). doi:10.1145/1807167.1807271
- 914 53. Gower, J., Lubbe, S. & Roux, N. *le. Understanding Biplots.* (John Wiley & Sons, Ltd, 2011).
915 doi:10.1002/9780470973196
- 916 54. Gabriel, K. R. The Biplot Graphic Display of Matrices with Application to Principal
917 Component Analysis. *Biometrika* **58**, 453 (1971).
- 918 55. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**,
919 2825–2830 (2011).
- 920 56. Turner, S. D. qqman: an R package for visualizing GWAS results using Q-Q and manhattan
921 plots. *bioRxiv* 005165 (2014). doi:10.1101/005165
- 922 57. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and traits.
923 *Nat. Genet.* **47**, 1236–1241 (2015).
- 924 58. Kroeze, W. K. *et al.* PRESTO-Tango as an open-source resource for interrogation of the
925 druggable human GPCRome. *Nat. Struct. Mol. Biol.* **22**, 362–369 (2015).
926

927 **Figures**

928 **Figure 1**

929



930

931

932

933

934

935

936

937

938

939

940

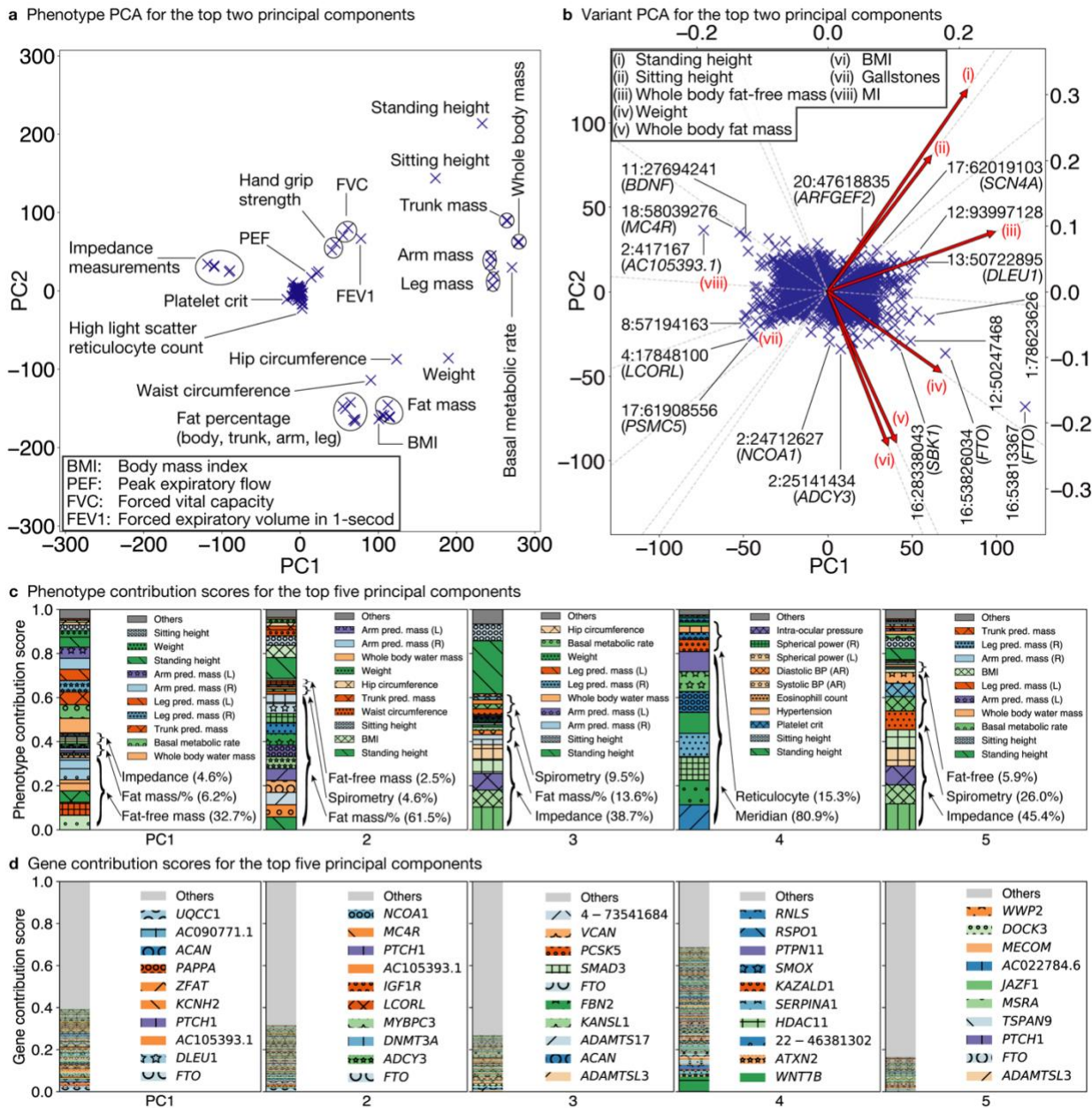
941

942

943

Fig.1 Illustrative study overview. **a** Summary of the UK Biobank genotype and phenotype data used in the study. We included White British individuals and analyzed LD-pruned and quality-controlled variants in relation to 2,138 phenotypes with a minimum of 100 individuals as cases (binary phenotypes) or non-missing values (quantitative phenotypes) (Supplementary Table S1-2). **b** Truncated singular value decomposition (TSVD) applied to decompose genome- and-phenome-wide summary statistic matrix W to characterize latent components. U , S , and V represent resulting matrices of singular values and vectors. **c** Decomposition of Genetic Associations (DeGAs) characterizes latent genetic components, which are represented as different colors on the palette, with an unsupervised learning approach – TSVD, followed by identification of the key components for each phenotype of our interest (painting phenotypes with colors) and annotation of each of the components with driving phenotypes, variants, and genes (finding the meanings of colors). **d** We used the squared cosine score and the contribution score, to quantify compositions and biomedical relevance of latent components. **e** We applied the genomic region enrichment analysis tool (GREAT) for biological characterization of each component and performed functional experiments focusing on adipocyte biology.

944 **Figure 2**

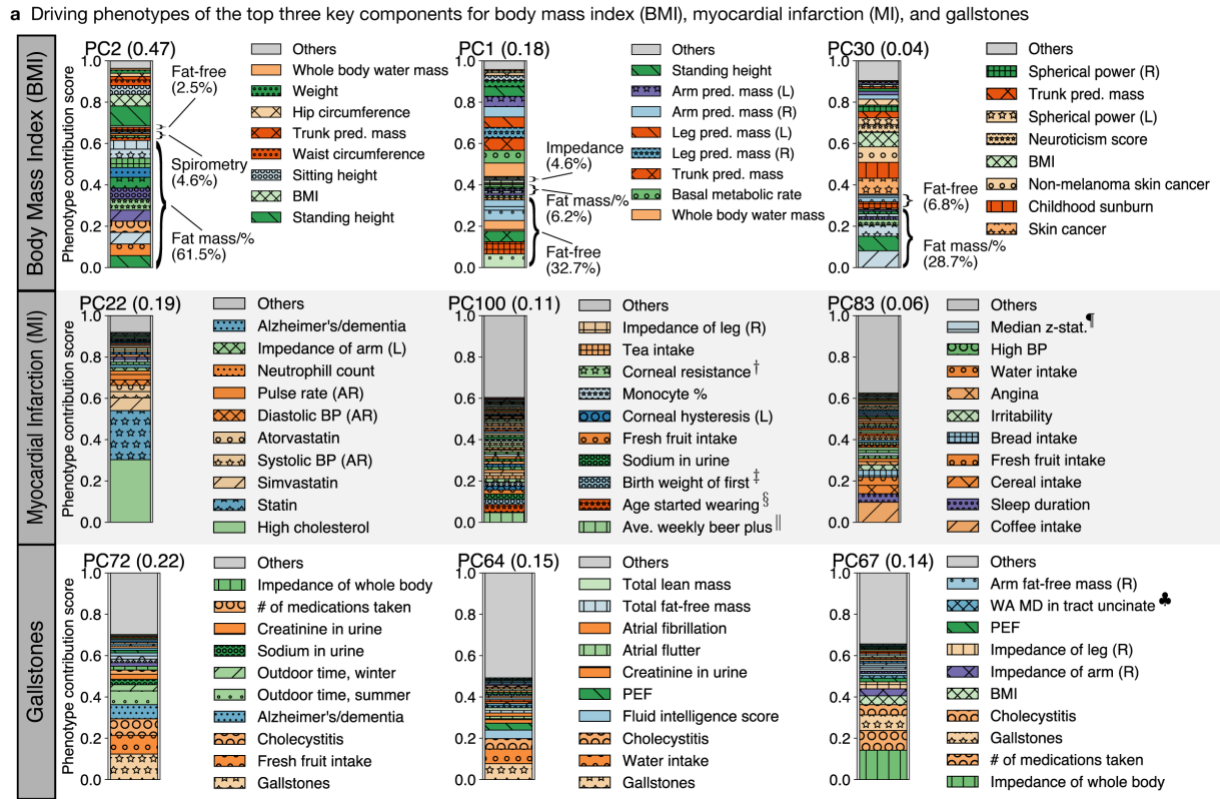


945
946
947
948
949
950
951
952
953
954
955
956
957
958

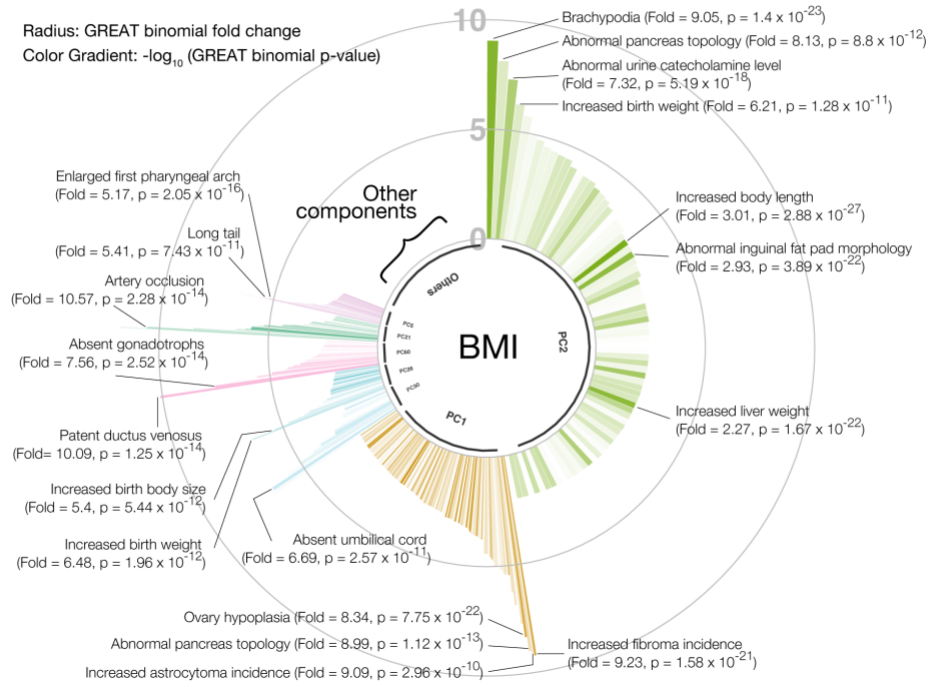
Fig. 2 Characterization of latent structures of genetic associations from genome- and-phenome-wide association summary statistics with DeGAs. **a-b** Components from truncated singular value decomposition (TSVD) corresponds to principal components in the phenotype (**a**) and variant (**b**) spaces. The first two components of all variants, excluding the MHC region, and relevant phenotypes are shown. **b** For variant PCA, we show biplot arrows (red) for selected phenotypes to help interpretation of the direction of principal components (Methods). The variants are labeled based on the genomic positions and the corresponding gene symbols. For example, "16:53813367 (*FTO*)" indicates the variant in gene *FTO* at position 53813367 on chromosome 16. **c-d** Phenotype (**c**) and gene (**d**) contribution scores for the first five components. PC1 is driven by largest part of the body mass that accounts for the "healthy part" (main text) including whole-body fat-free mass and genetic variants on *FTO* and *DLEU1*, whereas PC2 is driven by fat-related measurements, PC3 is driven by bioelectrical impedance measurements, PC4 is driven by eye measurements, and PC5 is driven by bioelectrical impedance and spirometry measurements along with the corresponding genetic variants (main text, Supplementary Table S3-4). Each colored segment represents a phenotype or gene with at least 0.5% and 0.05% of phenotype and gene contribution scores, respectively, and the rest is aggregated as others on the top of the stacked bar plots. The major

959 contributing phenotype groups (Methods, Supplementary Table S3) and additional top 10 phenotypes and the top 10 genes for each
960 component are annotated in **c** and **d**, respectively. pred.: predicted, %: percentage, mass/% mass and percentage, BP: blood
961 pressure, AR: automated reading, L: left, R: right.

962 Figure 3



b Ontology enrichment analysis with the genomic region enrichment analysis tool (GREAT) for body mass index

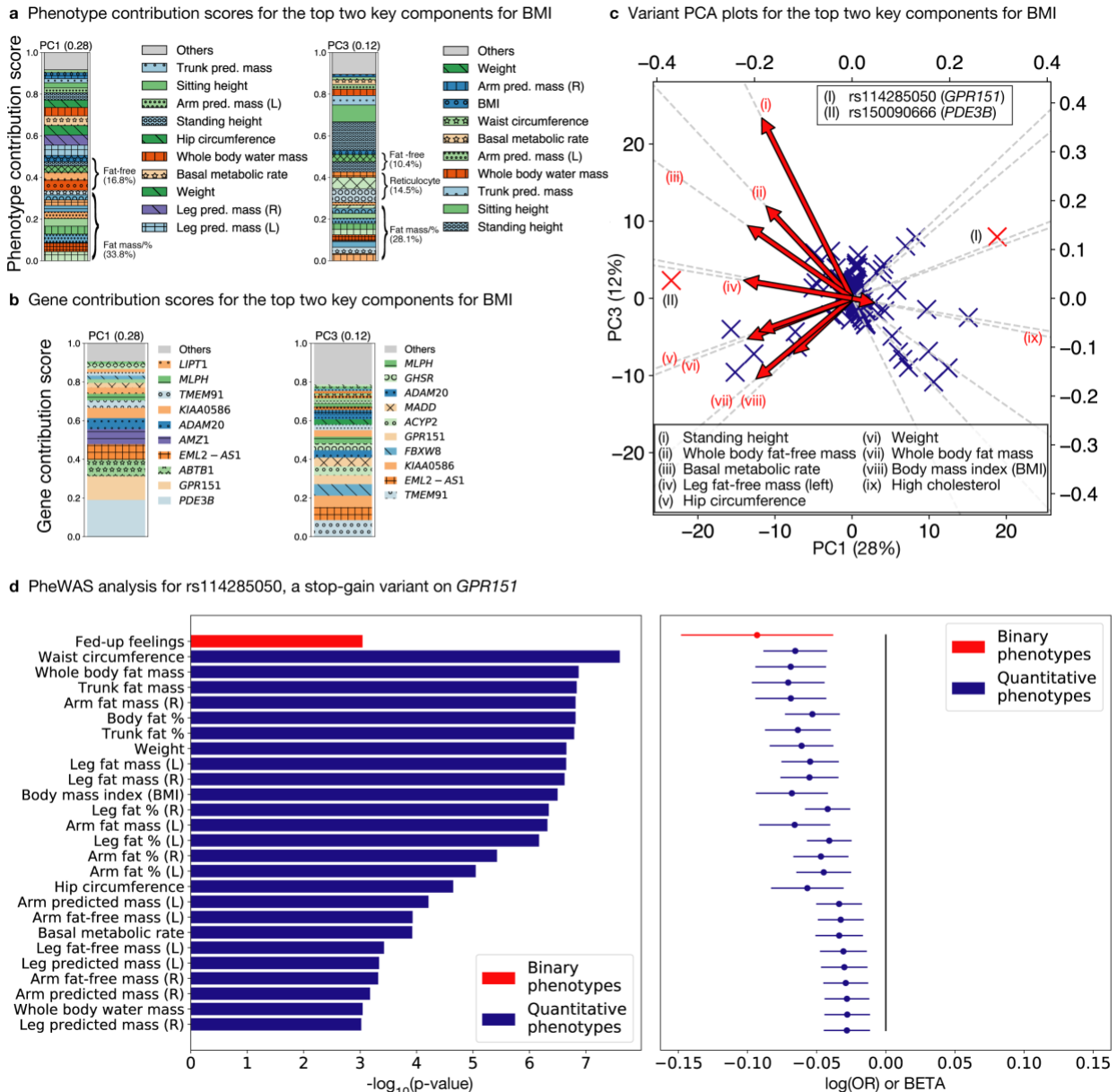


963
964
965
966

Fig.3 The top three key latent components from DeGAs of coding and non-coding variants for body mass index (BMI), myocardial infarction (MI), and gallstones. **a** The top three key components for each phenotype are identified by phenotype squared cosine scores and characterized with the driving phenotypes by phenotype contribution scores (Methods). Each colored segment

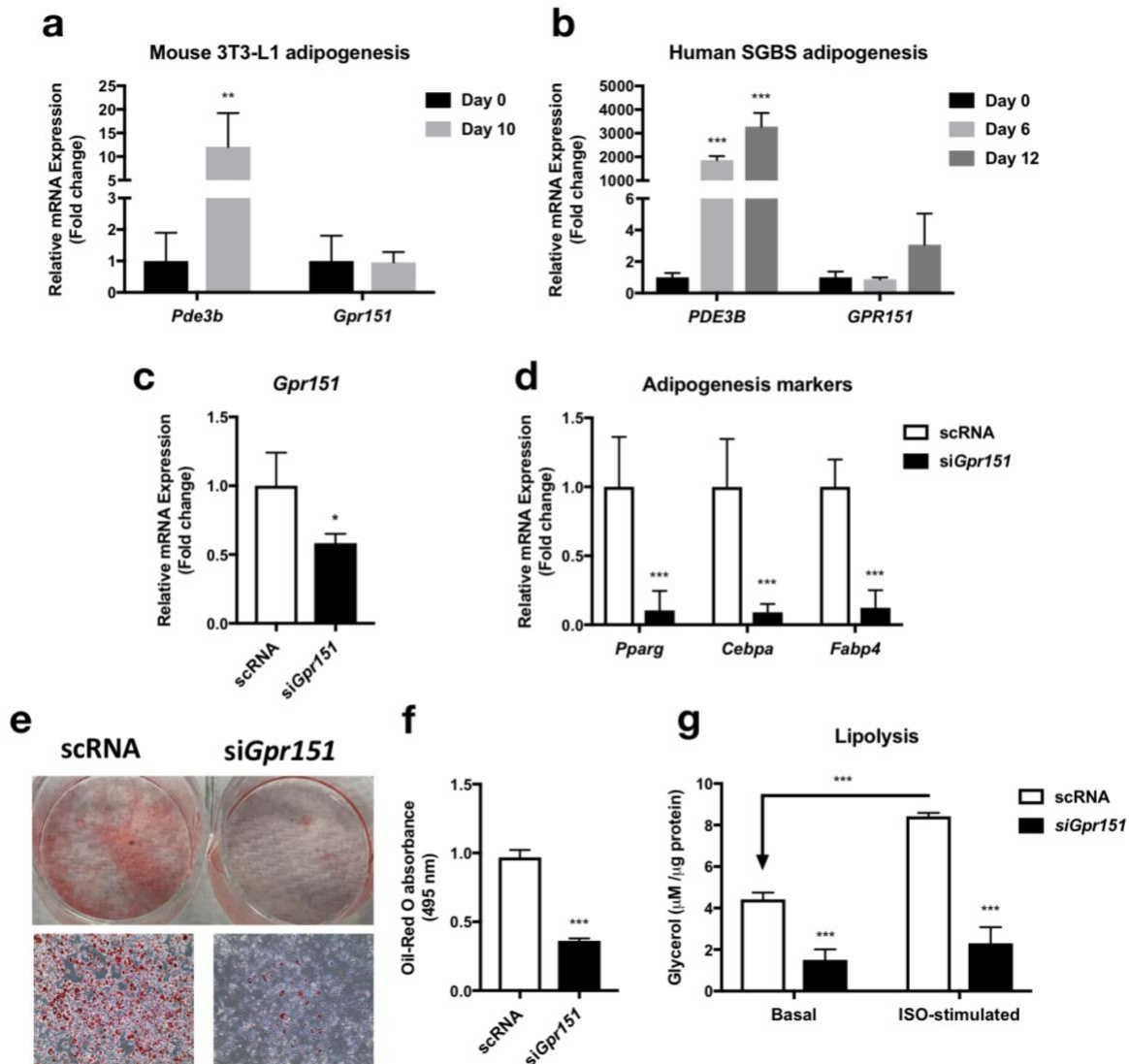
967 represents a phenotype with at least 0.5% of phenotype contribution scores for each of the component and the rest of the
968 phenotypes are aggregated as others and shown as the gray bar on the top. For BMI, additional phenotype grouping is applied
969 (Methods, Supplementary Table S3). **b** Biological characterization of driving non-coding and coding variants of the key components
970 for BMI with GREAT. The key components are shown proportional to their squared cosine score along with significantly enriched
971 terms in mouse MGI phenotype ontology. The radius represents binomial fold change and the color gradient represents p-value
972 from GREAT ontology enrichment analysis. pred.: predicted, #: number, %: percentage, mass/% mass and percentage, BP: blood
973 pressure, AR: automated reading, L: left, R: right, WA: weighted average. †: Corneal resistance factor (right), ‡: Birth weight of first
974 child, §: Age started wearing glasses or contact lenses, ||: Average weekly beer plus cider intake, ¶: Median z-statistic (in group-
975 defined mask) for shapes activation, ♣: Weighted-mean MD in tract uncinata fasciculus (right).

976 **Figure 4**



977
 978 **Fig. 4** DeGAs applied to the protein-truncating variants (PTVs) dataset. **a-b** Phenotype (**a**) and gene (**b**) contribution scores for the
 979 top key components associated with BMI based on phenotype grouping (Methods, Supplementary Table S3). **c** Variant PCA plot
 980 with biplot annotations for the top two components (Methods). The identified targets for functional follow-up (main text) are marked
 981 as (I) rs114285050 (a stop-gain variant on *GPR151*) and (II) rs150090666 (*PDE3B*). **d** Phenome-wide association analysis for
 982 *GPR151* rs114285050. The p-values (left) and log odds ratio (OR) (binary phenotypes, shown as red) or beta (quantitative
 983 phenotypes, shown as blue) (right) along with 95% confidence interval are shown for the phenotypes with minimum case count of
 984 1,000 (binary phenotypes) or 1,000 individuals with non-missing values (quantitative phenotypes) and strong association ($p < 0.001$)
 985 and with this variants among all the phenotypes used in the study. L: left, R: right, %: percentage, pred: predicted.

986 Figure 5



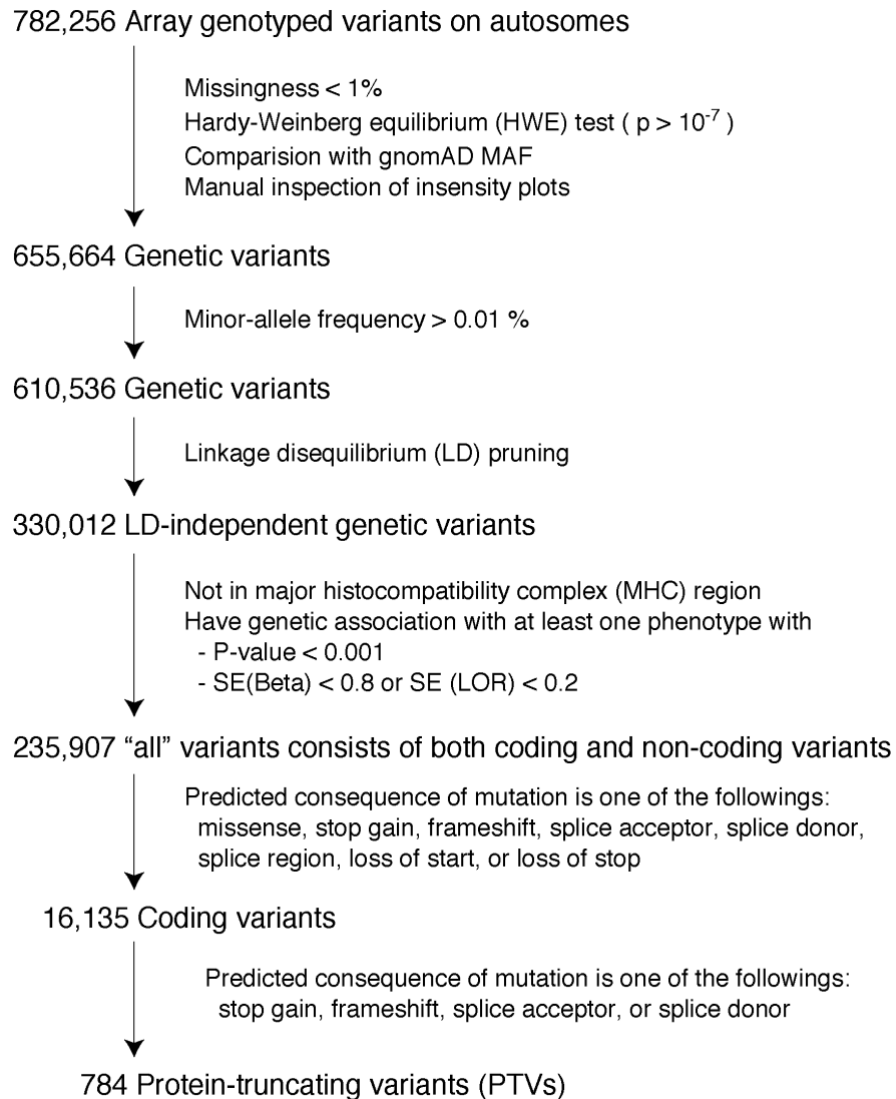
987
 988 **Fig. 5** Experimental validation of *GPR151* and *PDE3B* function in cellular models of adipogenesis. **a-b** qPCR analysis of gene
 989 expression patterns of *PDE3B* and *GPR151* during **(a)** mouse 3T3-L1 adipogenesis and **(b)** human SGBS adipogenesis. **c** qPCR
 990 analysis of *Gpr151* mRNA knockdown in 3T3-L1 preadipocytes. **d** qPCR analysis of the effect of siGpr151 knockdown on
 991 adipogenesis markers, *Pparg*, *Cebpa* and *Fabp4*. **e-g** Oil-Red O staining **(e)**, quantification of lipid droplets **(f)**, and lipolysis **(g)**
 992 in scRNA- or siGpr151-transfected adipocytes. Means \pm SEM are shown (***)p-value<0.001, **p-value<0.01, *p-value <0.05). scRNA:
 993 scrambled siRNA. ISO: isoproterenol.

994 Supplementary Materials

995 List of supplementary materials

- 996 ● Fig. S1: Variant filtering workflow
- 997 ● Fig. S2: Scree plot for TSVD of the GWAS summary statistics
- 998 ● Fig. S3: Squared cosine score (all variants dataset)
- 999 ● Fig. S4: Scree plot for TSVD of the phenotype data
- 1000 ● Fig. S5 TSVD of the phenotype data
- 1001 ● Fig. S6: GWAS analysis of the decomposed phenotypes
- 1002 ● Fig. S7: Genetic correlation of phenotype components
- 1003 ● Fig. S8: Intensity of genetic correlation of phenotype components
- 1004 ● Fig. S9: Gene contribution score (all variants dataset)
- 1005 ● Fig. S10: Variant PCA plot for MI
- 1006 ● Fig. S11: Variant PCA plot for gallstones
- 1007 ● Fig. S12: Robustness analysis – Biplots
- 1008 ● Fig. S13: Robustness analysis – Top 5 PCs
- 1009 ● Fig. S14: Robustness analysis – BMI
- 1010 ● Fig. S15: Robustness analysis – MI
- 1011 ● Fig. S16: Robustness analysis – Gallstones
- 1012 ● Fig. S17: GREAT enrichment analysis for MI
- 1013 ● Fig. S18: GREAT enrichment analysis for gallstones
- 1014 ● Fig. S19: Similarity of the top enriched terms for each DeGAs component
- 1015 ● Fig. S20: Squared cosine score for coding dataset
- 1016 ● Fig. S21: Phenotype contribution scores for coding dataset
- 1017 ● Fig. S22: Gene contribution scores for coding dataset
- 1018 ● Fig. S23: Squared cosine score of BMI (PTVs dataset)
- 1019 ● Fig. S24: PheWAS analysis for *PDB3B*
- 1020 ● Fig. S25: Univariate regression analysis for *GPR151*
- 1021 ● Fig. S26: Univariate regression analysis for *PDE3B*
- 1022 ● Fig. S27: *GPR151* overexpression
- 1023 ● Fig. S28: Effects of *Pde3b* knockdown in 3T3-L1 adipogenesis
- 1024 ● Table S1: List of phenotype categories
- 1025 ● Table S2: List of phenotypes
- 1026 ● Table S3: Phenotype groupings for visualization
- 1027 ● Table S4: Summary of contribution scores for the key components
- 1028 ● Table S5: GREAT enrichment analysis for BMI
- 1029 ● Table S6: GREAT enrichment analysis for MI
- 1030 ● Table S7: GREAT enrichment analysis for gallstones
- 1031 ● Table S8: PheWAS analysis for rs114285050 (*GPR151*)
- 1032 ● Table S9: PheWAS analysis for rs150090666 (*PDE3B*)
- 1033 ● Table S10: Genetic correlation of summary statistics for 10 traits with different GWAS
- 1034 covariates

1035 **Fig. S1: Variant filtering workflow**



1036

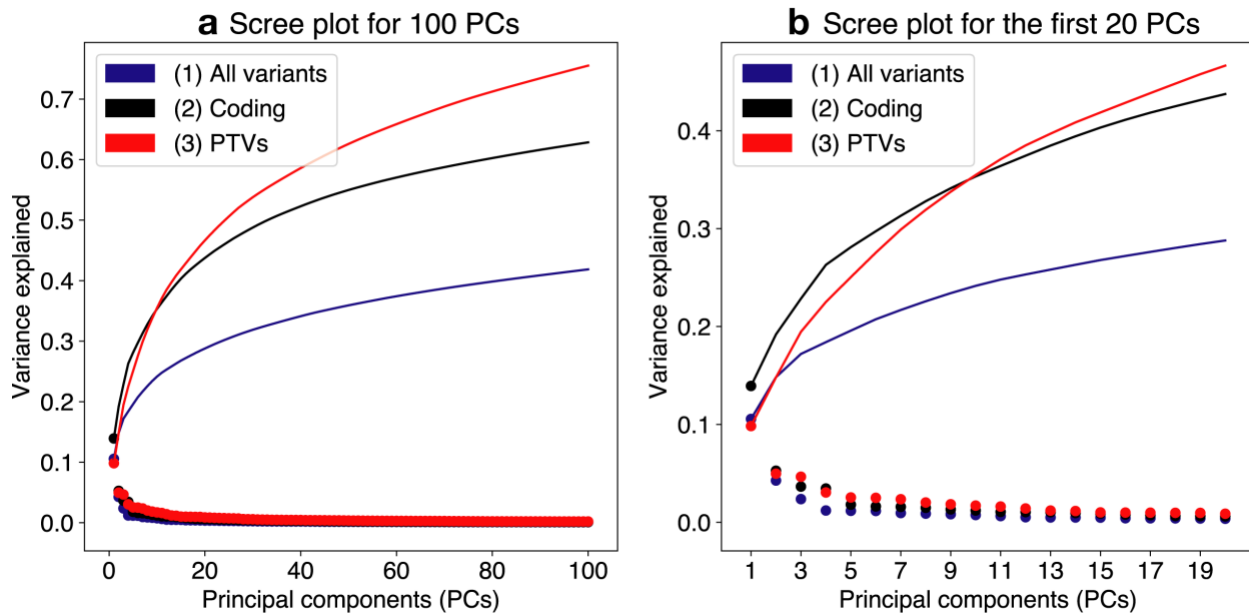
1037 **Fig. S1** Illustrative summary of the variant filters used in the study. The last three variant sets

1038 ("all" variants, coding variants, and PTVs) are used in the study. SE: standard error. LOR: log

1039 odds ratio.

1040 Fig. S2: Scree plot for TSVD of the GWAS summary statistics

1041

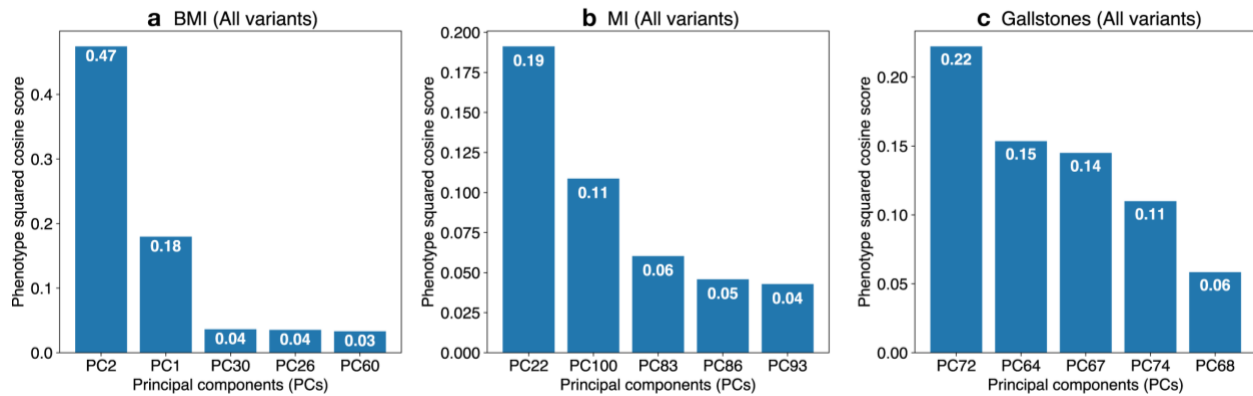


1042

1043

1044 **Fig. S2** Scree plot summarizes variance explained in each of the top 100 (a) and 20 (b)
1045 components. The scree plots are shown for three datasets consists of LD-pruned and QC-
1046 filtered sets of array-genotyped variants outside of MHC region: (1) all array-genotyped variants,
1047 which includes coding and non-coding variants (blue), (2) coding variants (black), and (3)
1048 protein-truncating variants (PTVs, red). For each component, we calculate the variance
1049 explained defined as squared eigenvalues divided by the total variance in the original matrix
1050 (Methods). We plotted those values as dots and cumulative values as lines.

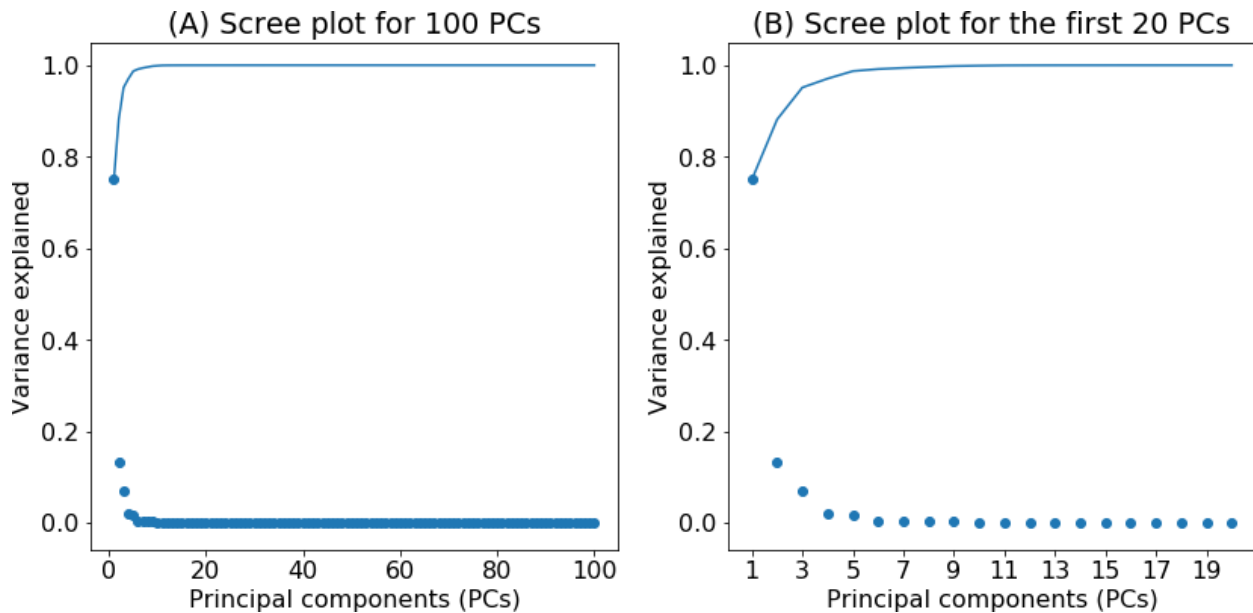
1051 **Fig. S3: Squared cosine score (all variants dataset)**



1052
1053

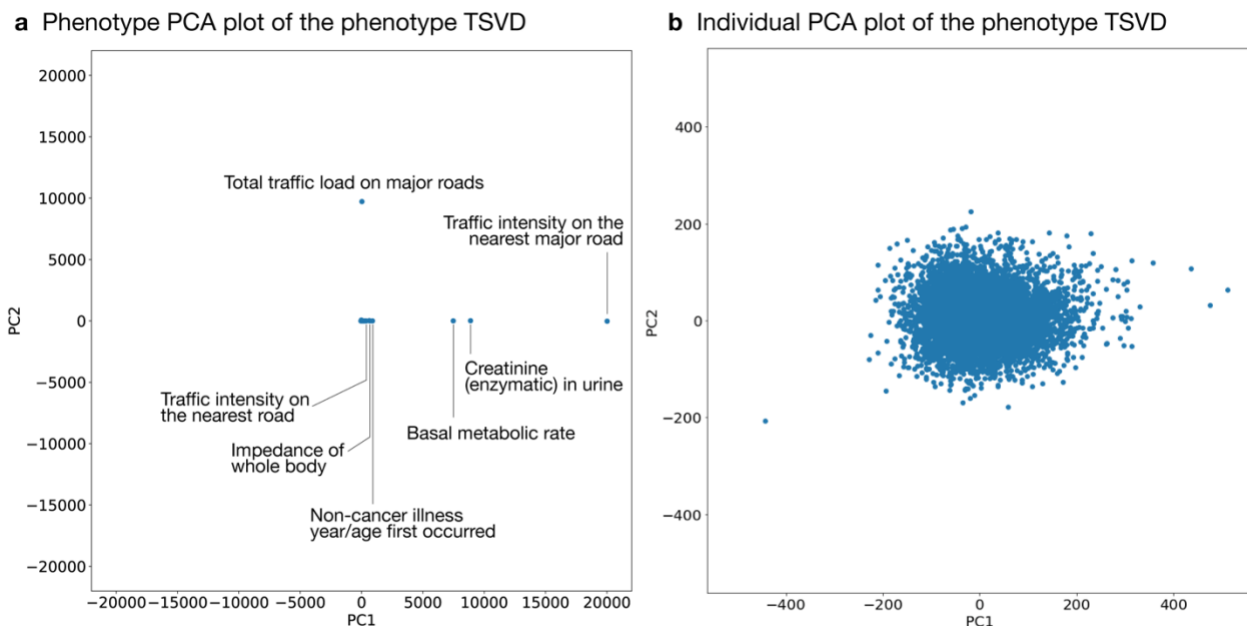
1054 **Fig. S3** Identification of the key components with phenotype squared cosine scores. Squared
1055 cosine score quantifies relative importance of the key components for a given phenotype. The
1056 top five key components are identified for all variant dataset that includes both coding and non-
1057 coding variants for three phenotypes: **a** body mass index (BMI), **b** myocardial infarction (MI),
1058 and **c** gallstones. The top five key components are shown on the horizontal axis and the
1059 corresponding squared cosine scores are shown on the vertical axis.

1060 Fig. S4: Scree plot for TSVD of phenotype data



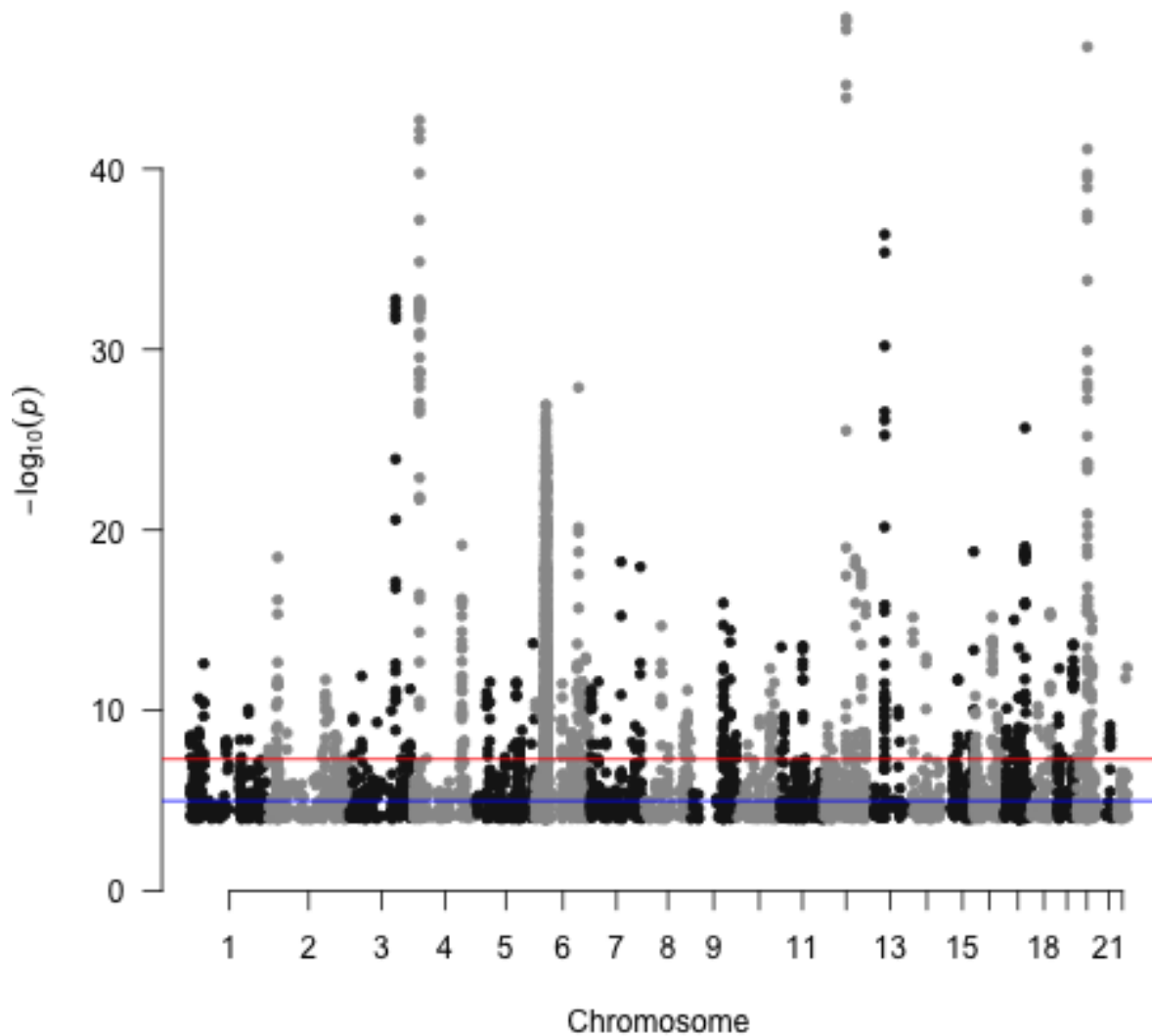
1061
1062 **Fig. S4** Scree plot summarizes variance explained in the top 100 (a) and 20 (b) components
1063 characterized from the imputed and normalized phenotype data. We calculate the variance
1064 explained defined as squared eigenvalues divided by the total variance in the original matrix
1065 (Methods). We plotted those values as dots and cumulative values as lines.

1066 Fig. S5 TSVD of the phenotype data



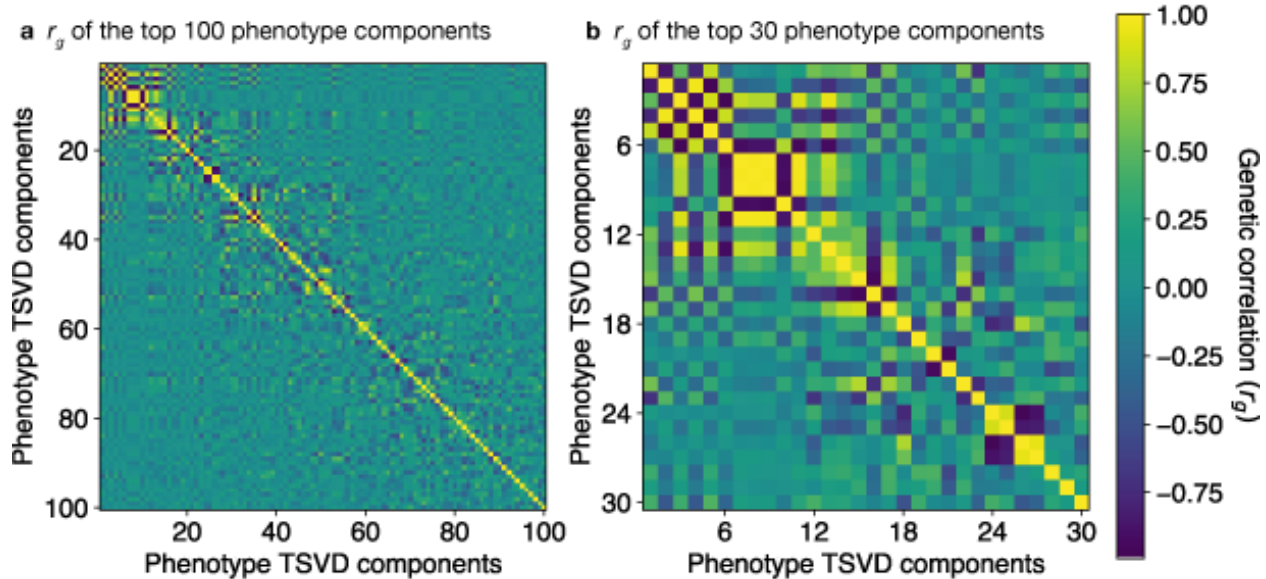
1067
1068 **Fig. S5** Characterization of latent structures of phenotypic data characterized by truncated
1069 singular value decomposition (TSVD) of the imputed and normalized phenotype data.
1070 Phenotype (a) and Individual (b) PCA plots summarizes the first two components.

1071 Fig. S6: GWAS analysis of the decomposed phenotypes



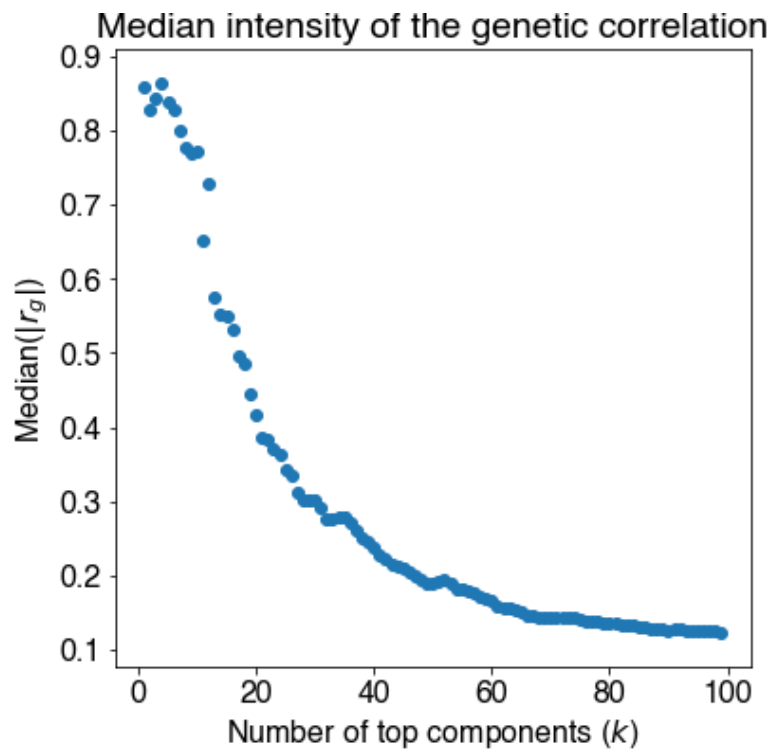
1072
1073 **Fig. S6** Genome-wide association analysis of the phenotype PCs. After characterizing the
1074 phenotype latent space with TSVD on the phenotype data, we performed GWAS analysis. The
1075 statistical significance for the first phenotype component is shown in the plot. The variants with p
1076 $< 1.0 \times 10^{-4}$ are shown. The red and blue lines indicate genome-wide significance (5.0×10^{-8})
1077 and genome-wide suggestive (5.0×10^{-5}) levels, respectively.
1078

1079 Fig. S7: Genetic correlation of phenotype components



1080
1081 **Fig. S7** Genetic correlation (r_g) of phenotype TSVD components shown for the top 100
1082 components (a) and the top 30 components (b), respectively.

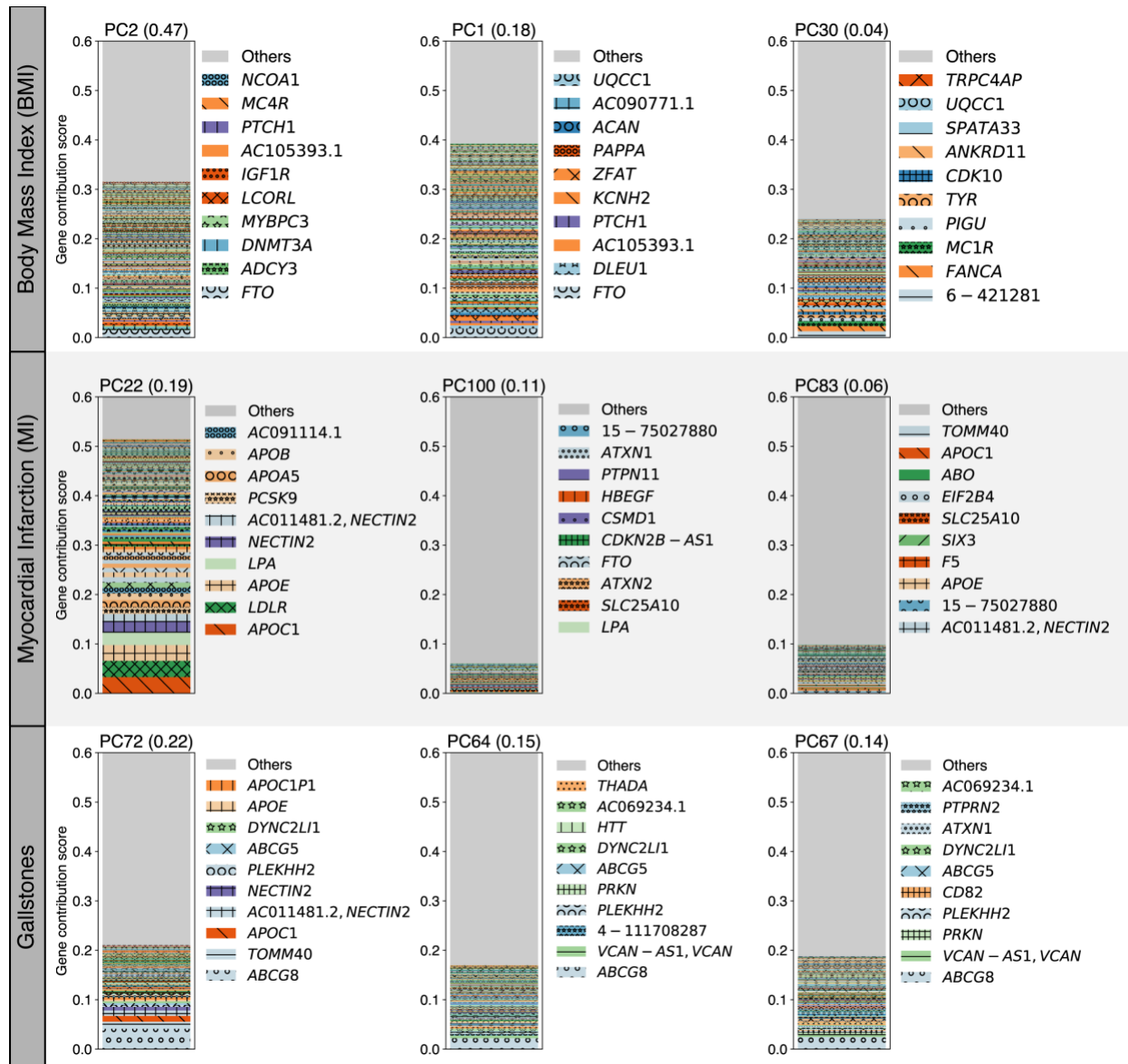
1083 Fig. S8: Intensity of genetic correlation of phenotype components



1084
1085 **Fig. S8** The median of the absolute value of the genetic correlation (r_g) among the top
1086 phenotypic components.

1087

Fig. S9: Gene contribution score (all variants dataset)



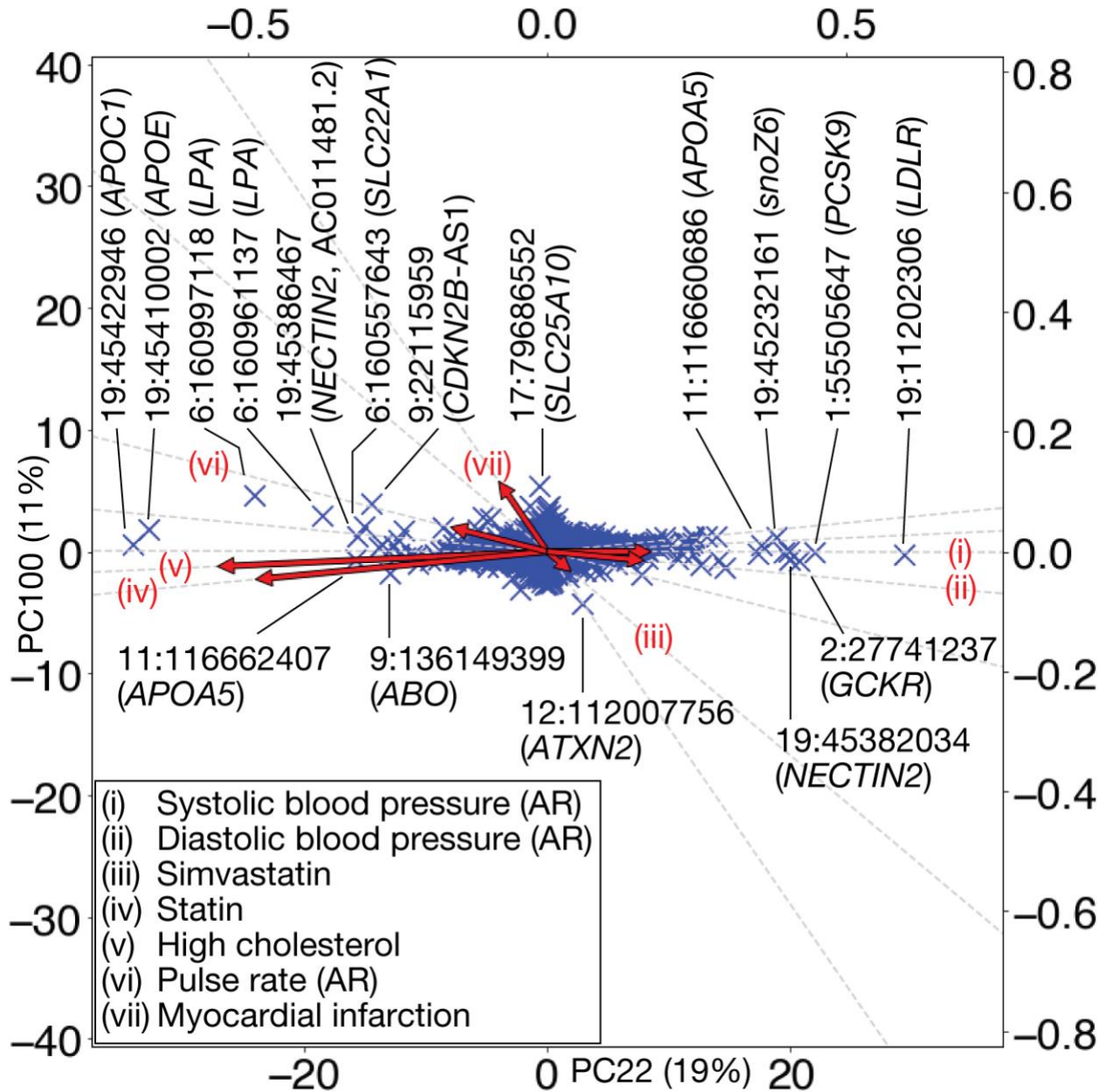
1088

1089

1090 **Fig. S9** Gene contribution scores for the top three key components for body mass index (BMI),
 1091 myocardial infarction (MI), and gallstones using all variant dataset, which includes both coding
 1092 and non-coding variants. For each phenotype, the top three key components with their
 1093 phenotype squared cosine scores are shown on the top of the stacked bar plot and gene
 1094 contribution scores for each of the components are shown as colored segments. Each colored
 1095 segment represents a gene with at least 0.05% of contribution scores and the rest of the genes
 1096 are aggregated as the gray bar at the top. For the visualization, the maximum value of the
 1097 vertical axis is set to be 0.6. For each component, the labels for the top 10 driving genes
 1098 are shown. For non-coding variants, we display their genomic coordinates.

1099 Fig. S10: Variant PCA plot for myocardial infarction.

1100



1101

1102

1103

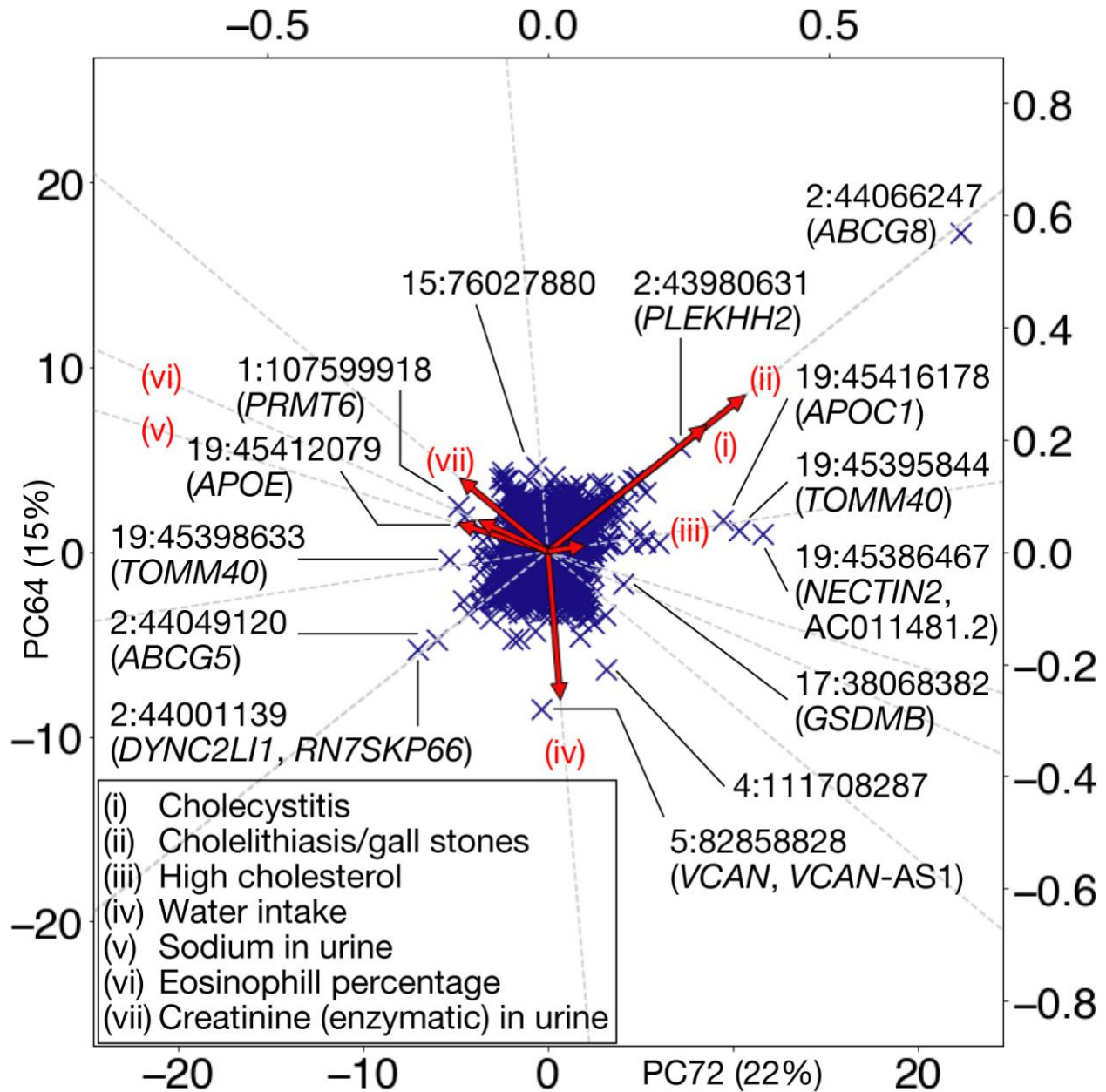
1104

1105

1106

Fig. S10 Variant PCA plot with biplot annotation for the top two key components for myocardial infarction using “all” dataset. Genetic variants projected into the top two key components, PC22 (horizontal axis) and PC100 (vertical axis) are shown as scatter plot. Variants are annotated with gene symbols. Directions of genetic associations for relevant phenotypes are annotated as red arrows using the secondary axes (Methods). Abbreviations. AR: automated reading.

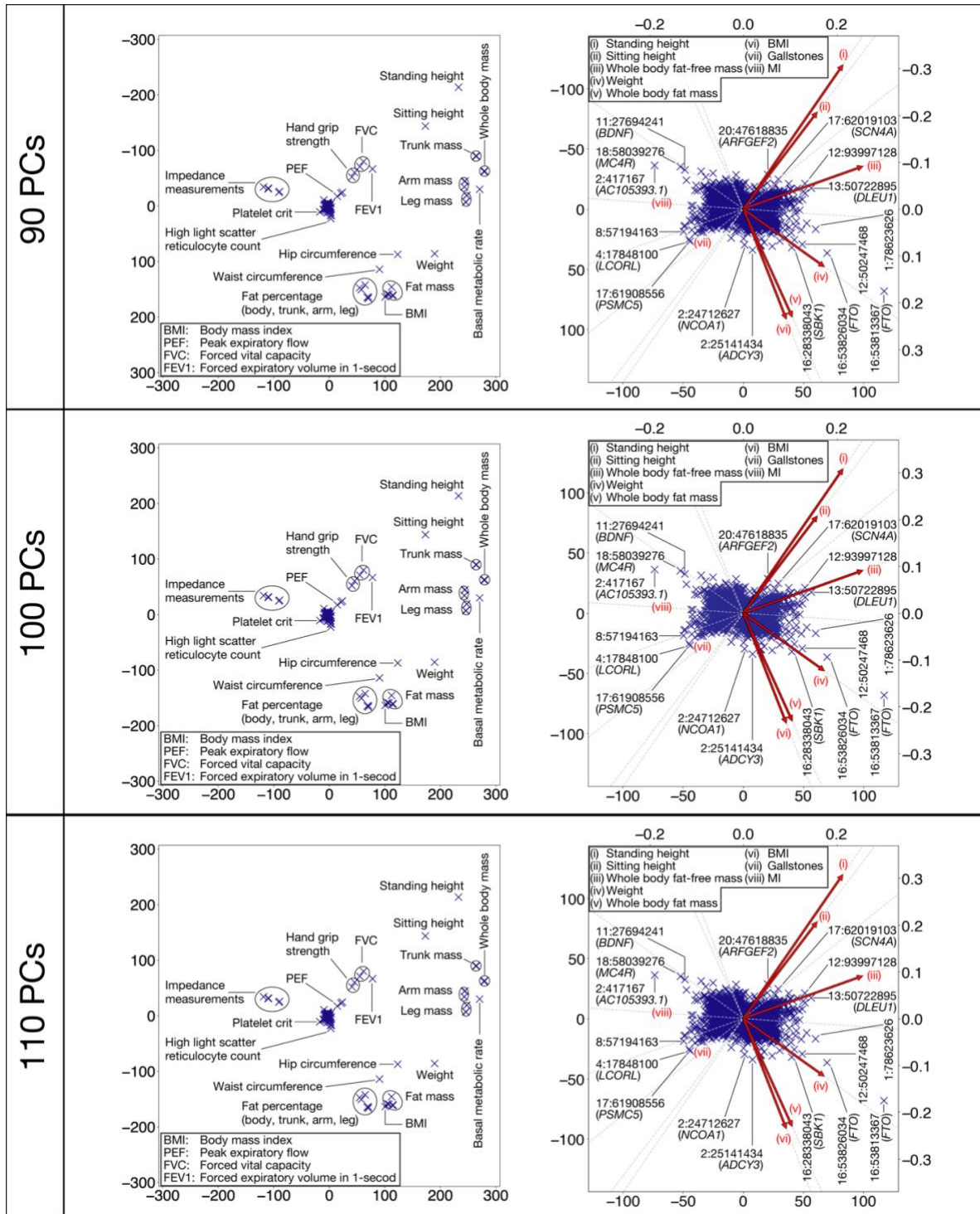
1107 Fig. S11: Variant PCA plot for Gallstones



1108
 1109 **Fig. S11** Variant PCA plot with biplot annotation for the top two key components for gallstones
 1110 using "all" dataset. Genetic variants projected into the top two key components, PC72
 1111 (horizontal axis) and PC64 (vertical axis). Variants are annotated with gene symbols. Directions
 1112 of genetic associations for relevant phenotypes are annotated as red arrows using the
 1113 secondary axes (Methods).

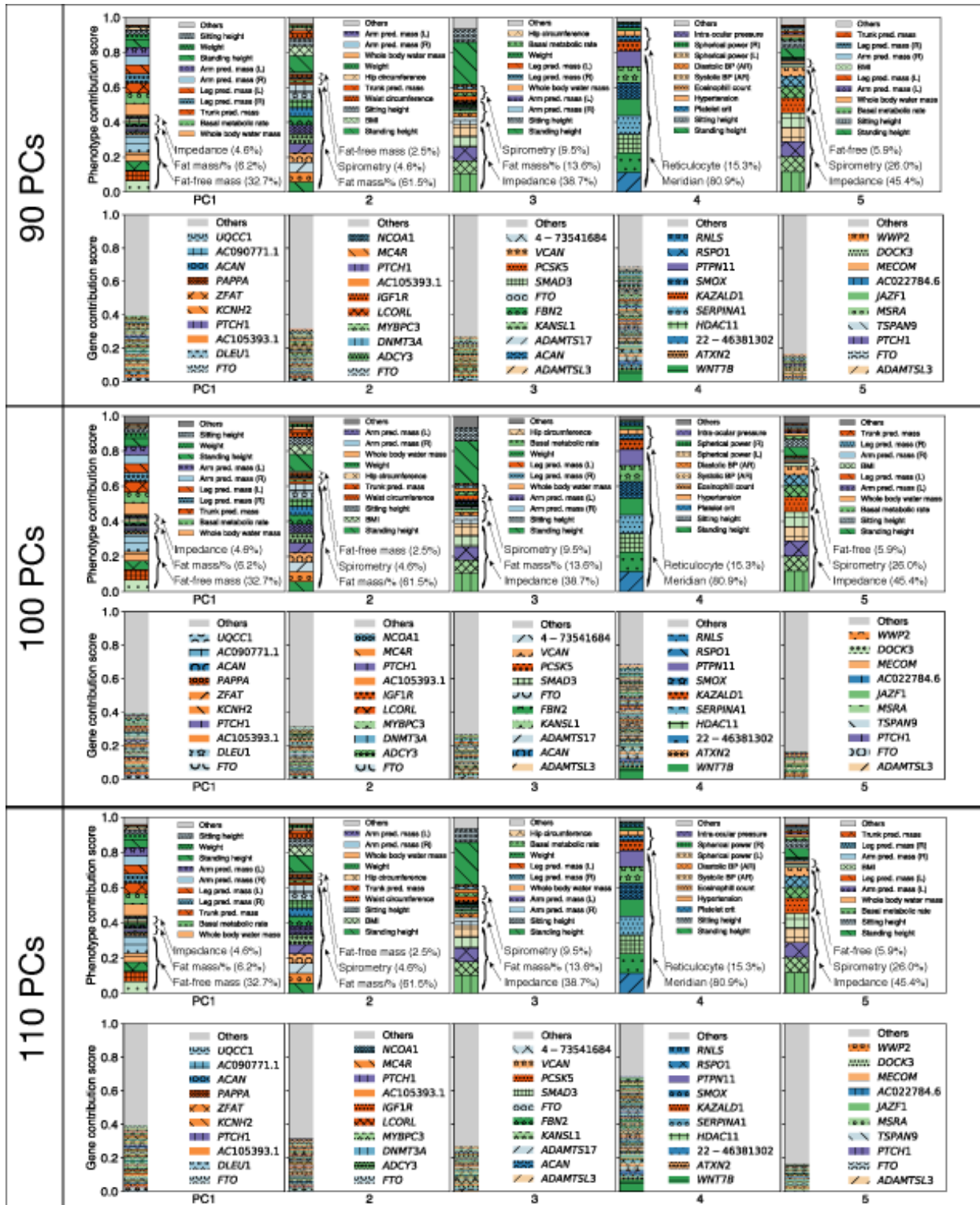
1114

1115 Fig. S12: Robustness analysis – Biplots



1116
 1117 **Fig. S12** Comparison of the top two DeGAs components by robustness analysis with respect to
 1118 the number of latent factors in DeGAs. The phenotype PCA plot (left) and the variant PCA plot
 1119 with the biplot annotations (right) are shown (Methods). To cope with the sign indeterminacy of
 1120 the latent components, the direction of PC2 is reversed in the plots for TSVD with 90 PCs.

1121 Fig. S13: Robustness analysis – Top 5 PCs

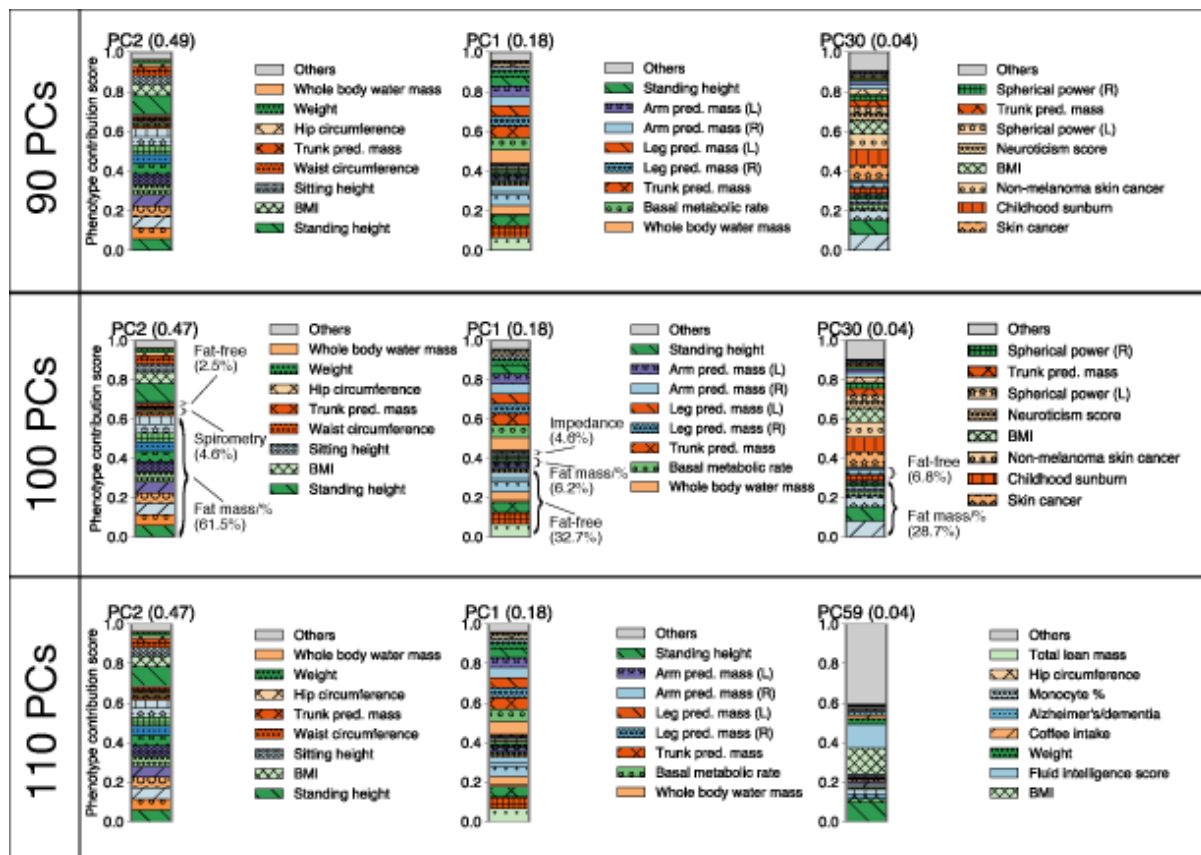


1122
1123
1124

Fig. S13 Comparison of the top five DeGAs components by robustness analysis with respect to the number of latent factors in DeGAs. The phenotype and gene contribution scores are shown.

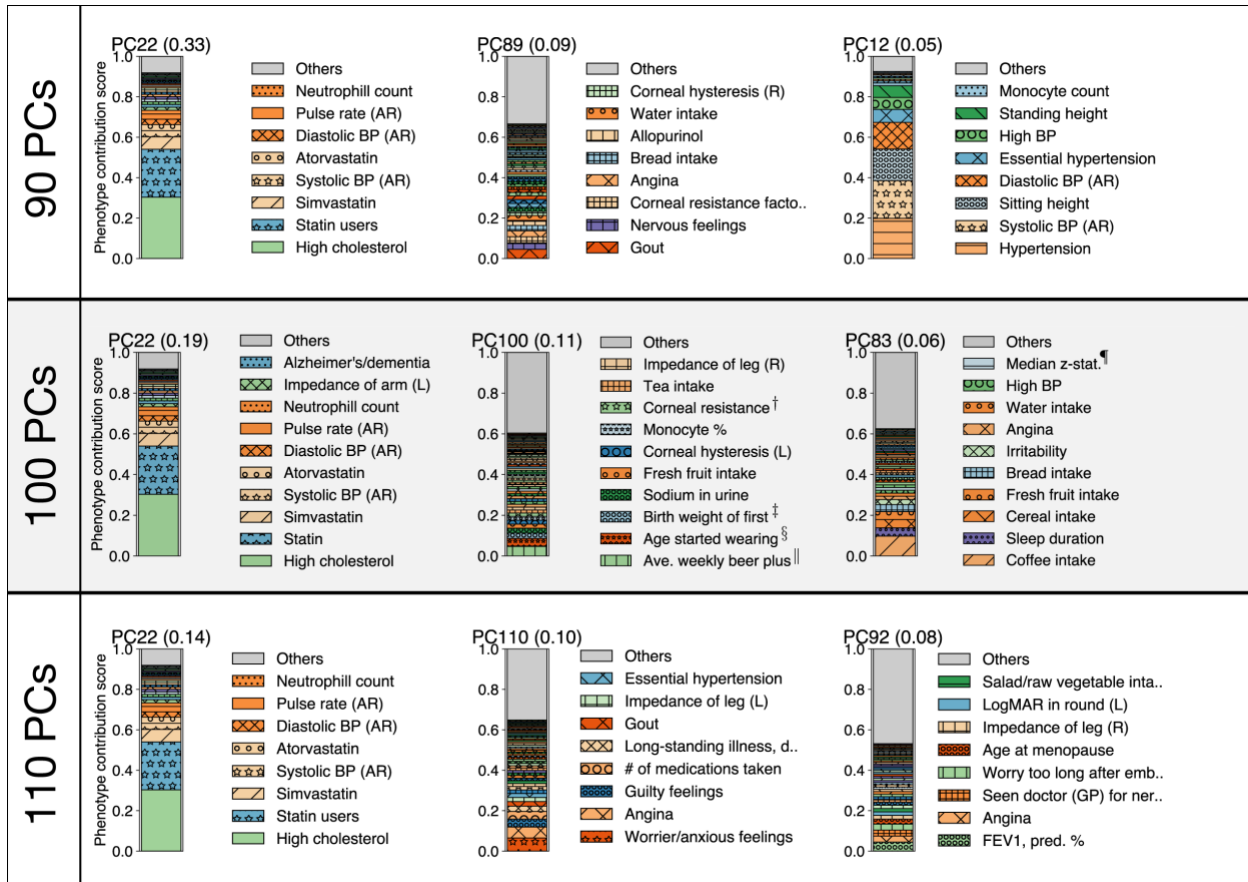
1125 Each colored segment represents a phenotype or gene with at least 0.5% and 0.05% of
 1126 phenotype and gene contribution scores, respectively, and the rest is aggregated as others on
 1127 the top of the stacked bar plots. The major contributing phenotype groups (Methods,
 1128 Supplementary Table S3) and additional top 10 phenotypes and the top 10 genes for each
 1129 component are annotated.

1130 Fig. S14: Robustness analysis – BMI



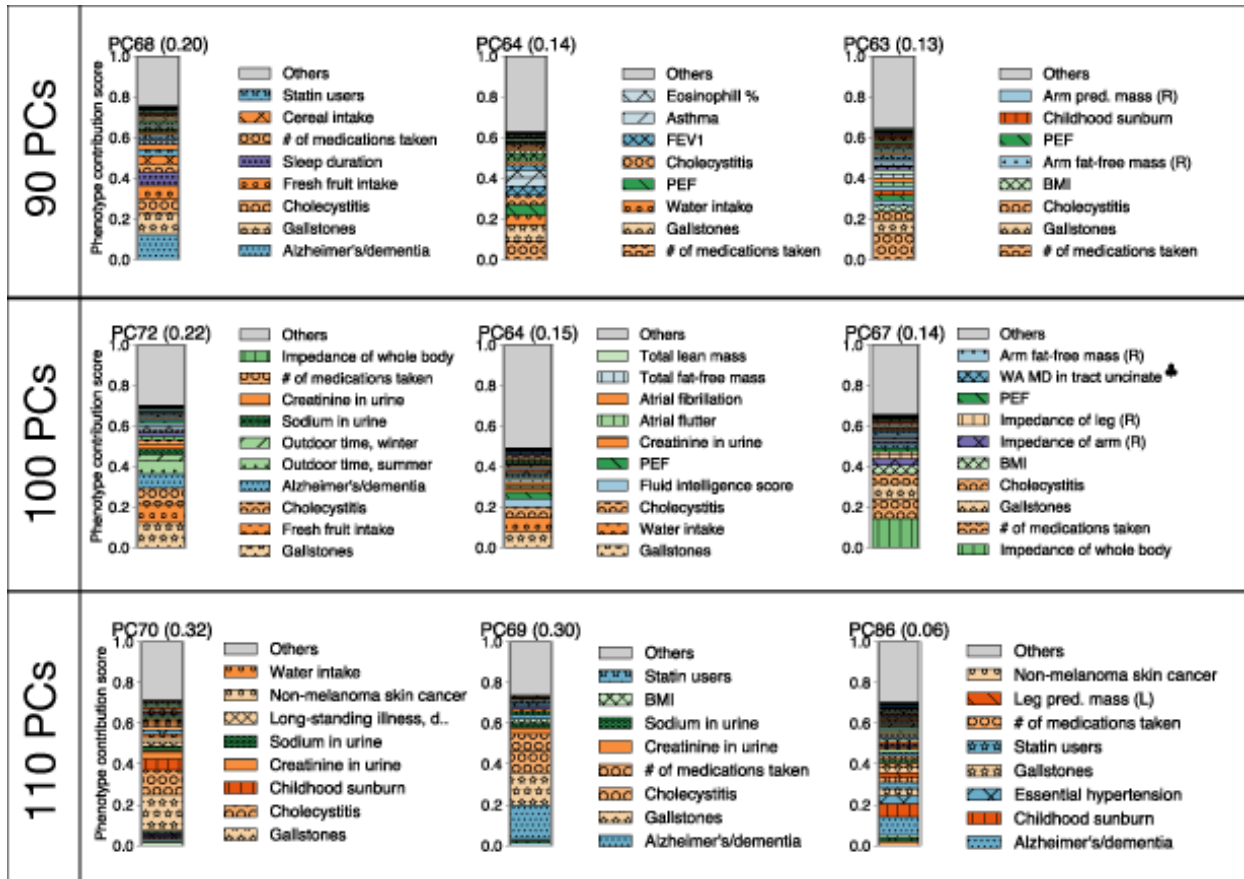
1131
 1132 **Fig. S14** Comparison of the key components for body mass index (BMI) by robustness analysis
 1133 with respect to the number of latent factors in DeGAs. For each condition, the top three key
 1134 components with their phenotype squared cosine scores are shown on the top of the stacked
 1135 bar plot and phenotype contribution scores for each of the components are shown as colored
 1136 segments. Each colored segment represents a gene with at least 0.5% of contribution scores
 1137 and the rest of the phenotypes are aggregated as the gray bar at the top. For each component,
 1138 the labels for the top 6 driving phenotypes are shown.

1139 Fig. S15: Robustness analysis – MI



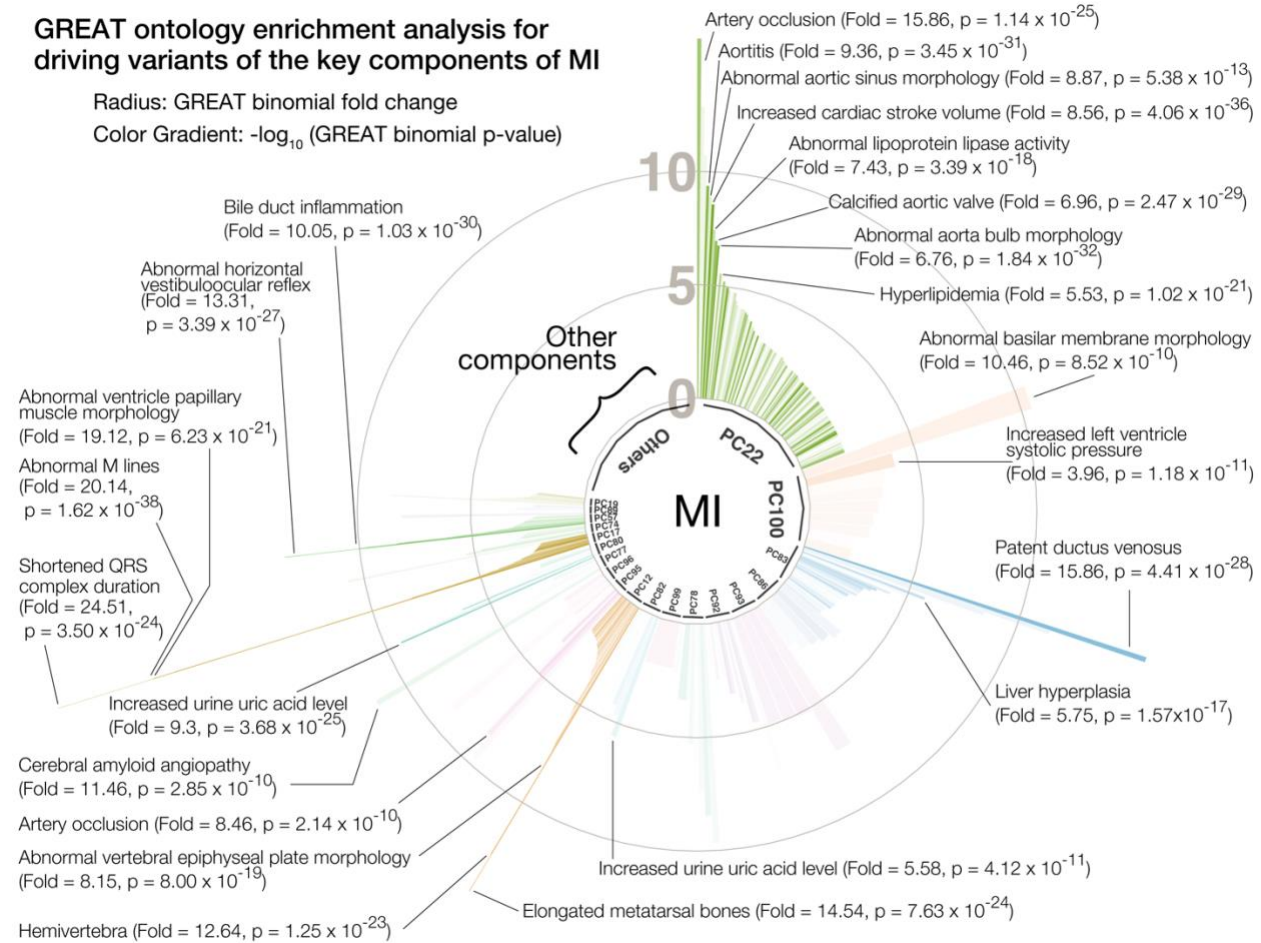
1140
 1141 **Fig. S15** Comparison of the key components for myocardial infarction (MI) by robustness
 1142 analysis with respect to the number of latent factors in DeGAs. For each condition, the top three
 1143 key components with their phenotype squared cosine scores are shown on the top of the
 1144 stacked bar plot and phenotype contribution scores for each of the components are shown as
 1145 colored segments. Each colored segment represents a gene with at least 0.5% of contribution
 1146 scores and the rest of the phenotypes are aggregated as the gray bar at the top. For each
 1147 component, the labels for the top 6 driving phenotypes are shown.

1148 Fig. S16: Robustness analysis – Gallstones



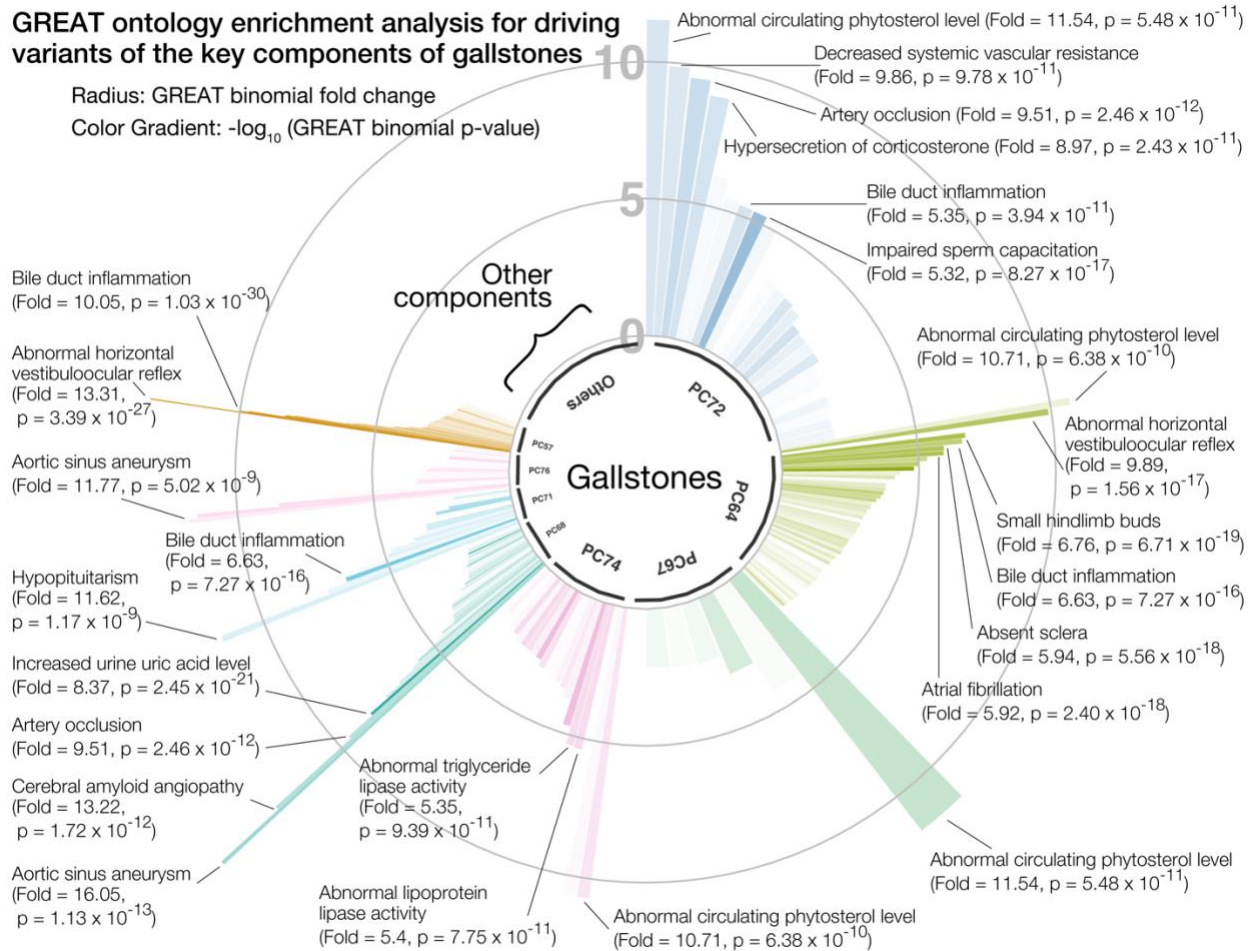
1149
 1150 **Fig. S16** Comparison of the key components for gallstones by robustness analysis with respect
 1151 to the number of latent factors in DeGAs. For each condition, the top three key components with
 1152 their phenotype squared cosine scores are shown on the top of the stacked bar plot and
 1153 phenotype contribution scores for each of the components are shown as colored segments.
 1154 Each colored segment represents a gene with at least 0.5% of contribution scores and the rest
 1155 of the phenotypes are aggregated as the gray bar at the top. For each component, the labels for
 1156 the top 6 driving phenotypes are shown.

1157 Fig. S17: GREAT enrichment analysis for MI



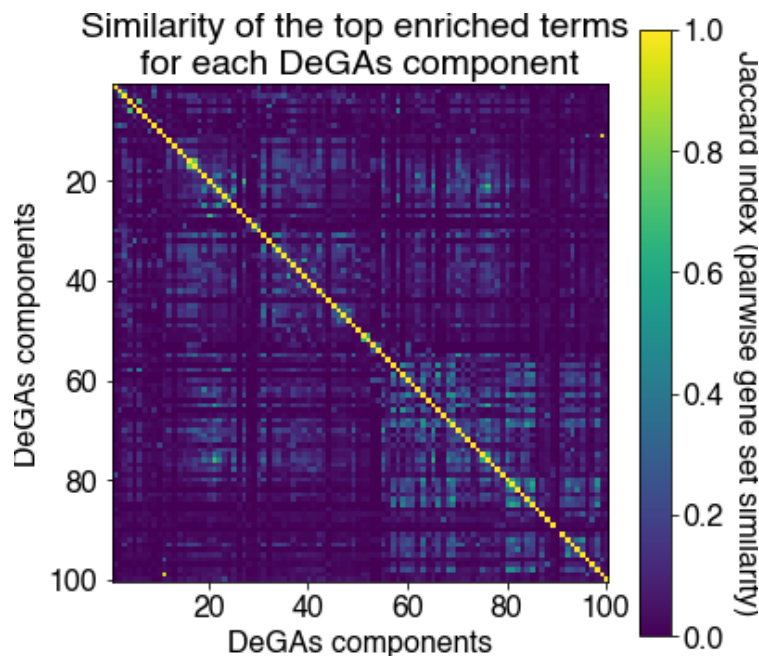
1158
1159 **Fig. S17** Biological characterization of driving non-coding and coding variants of the key
1160 components for myocardial infarction (MI) with the genomic region enrichment analysis tool
1161 (GREAT) using the all variants dataset. The key components are shown proportional to their
1162 squared cosine score along with significantly enriched terms in mouse genome informatics
1163 (MGI) phenotype ontology. The radius represents binomial fold change and the color gradient
1164 represents p-value from GREAT ontology enrichment analysis.

1165 **Fig. S18: GREAT enrichment analysis for gallstones**



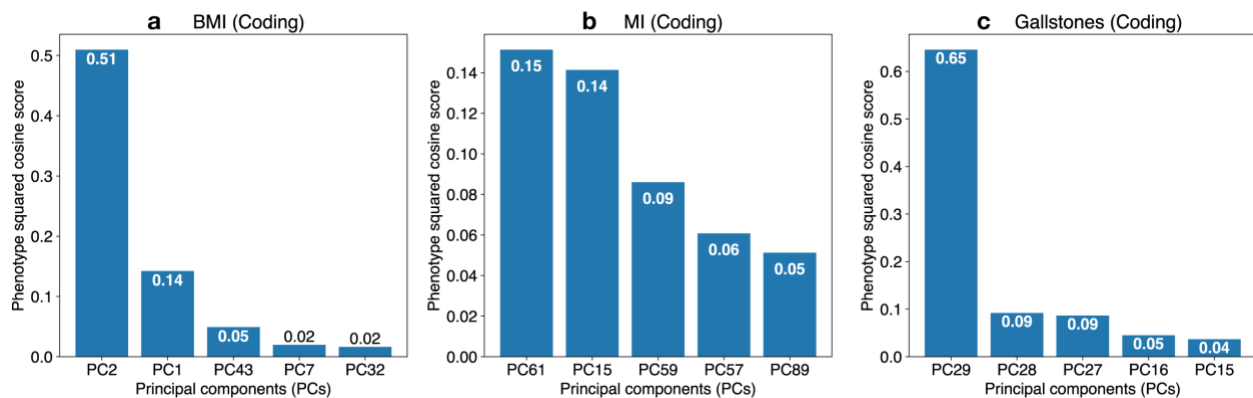
1166
1167 **Fig. S18** Biological characterization of driving non-coding and coding variants of the key
1168 components for gallstones with the genomic region enrichment analysis tool (GREAT) using the
1169 all variants dataset. The key components are shown proportional to their squared cosine score
1170 along with significantly enriched terms in mouse genome informatics (MGI) phenotype ontology.
1171 The radius represents binomial fold change and the color gradient represents p-value from
1172 GREAT ontology enrichment analysis.

1173 Fig. S19: Similarity of the top enriched terms for each DeGAs
1174 component



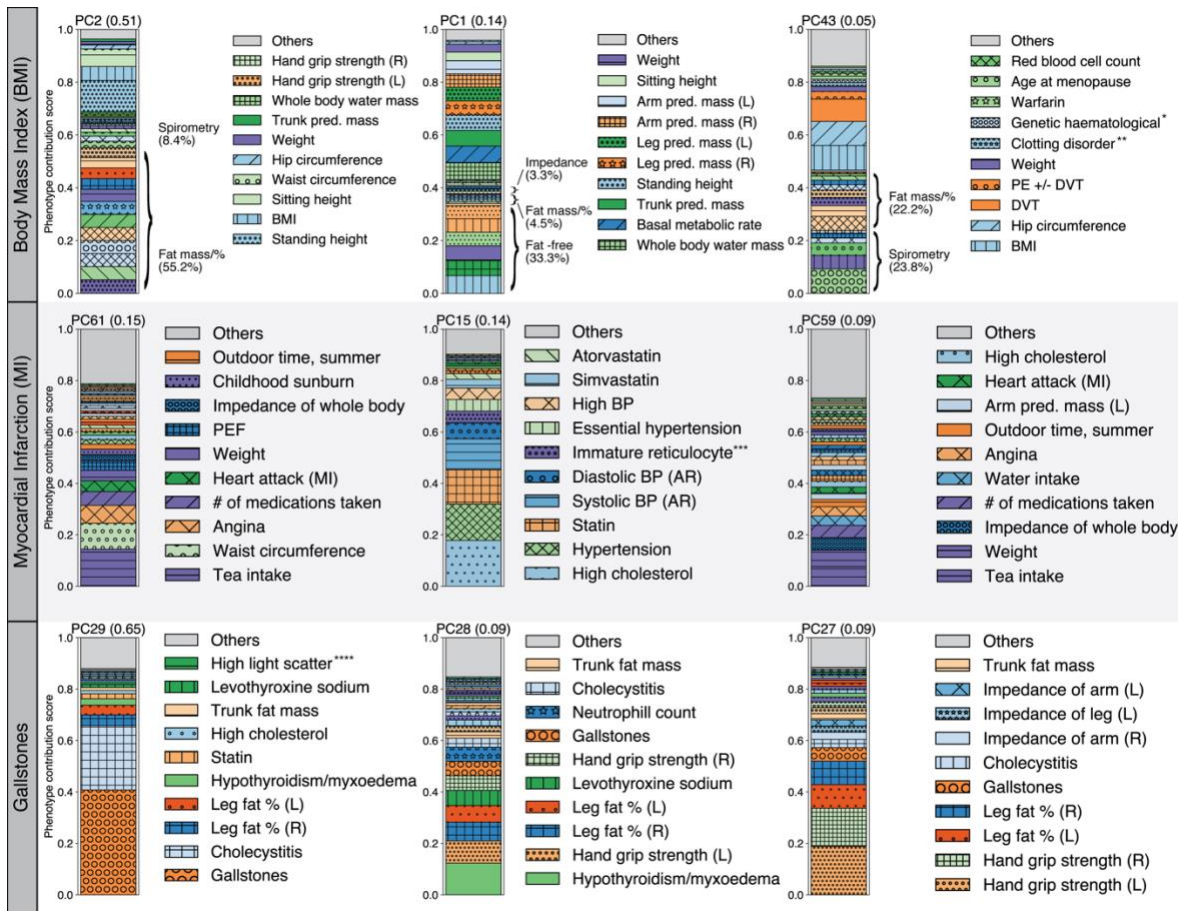
1175
1176 **Fig. S19** Similarity of the top enriched terms for each DeGAs component. For each DeGAs
1177 component, we took the top enriched ontology terms identified by GREAT and obtained the list
1178 of genes annotated with that term. Using these gene sets, we quantified the pairwise gene set
1179 similarity across the 100 DeGAs components using Jaccard Index.

1180 Fig. S20: Squared cosine score for coding dataset



1181
1182 **Fig. S20** Identification of the key components with phenotype squared cosine scores using
1183 coding dataset. Squared cosine score quantifies relative importance of the key components for
1184 a given phenotype. The top five key components are identified for coding dataset for three
1185 phenotypes: **a** body mass index (BMI), **b** myocardial infarction (MI), and **c** gallstones. The top
1186 five key components are shown on the horizontal axis and the corresponding squared cosine
1187 scores are shown on the vertical axis.

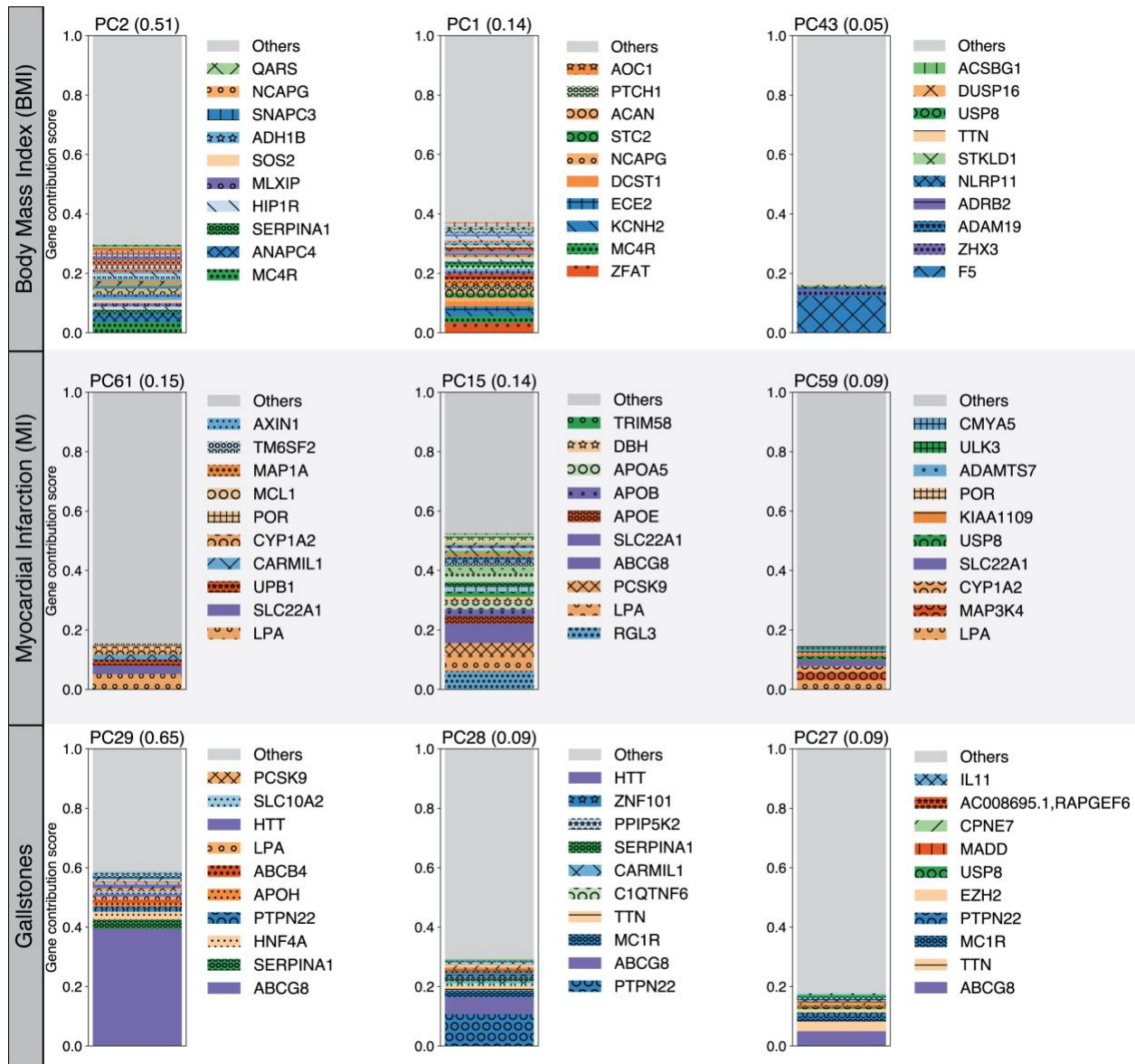
1188 Fig. S21: Phenotype contribution scores for coding dataset



1189
 1190 **Fig. S21** Phenotype contribution scores for the top three key components for body mass index
 1191 (BMI), myocardial infarction (MI), and gallstones using coding dataset. For each phenotype, the
 1192 top three key components with their phenotype squared cosine scores are shown on the top of
 1193 the stacked bar plot and phenotype contribution scores for each of the components are shown
 1194 as colored segments. Each colored segment represents a phenotype with at least 0.5% of
 1195 contribution scores and the rest of the genes are aggregated as the gray bar at the top. For
 1196 BMI, additional phenotype grouping is applied (Methods, Supplementary Table S3). For each
 1197 component, the labels for the top 10 driving genes are shown.
 1198

1199

Fig. S22: Gene contribution scores for coding dataset



1200

1201

1202

1203

1204

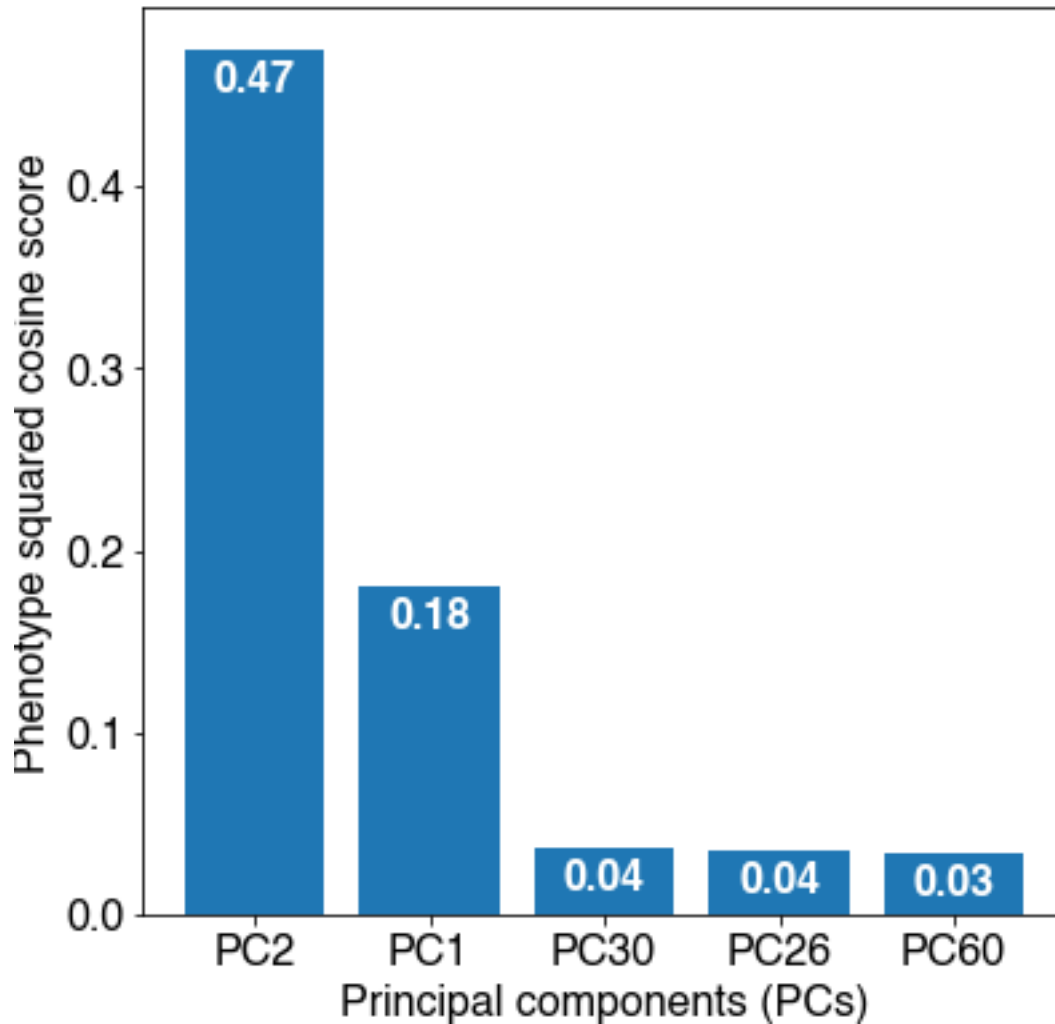
1205

1206

1207

Fig. S22 Gene contribution scores for the top three key components for body mass index (BMI), myocardial infarction (MI), and gallstones using coding dataset. For each phenotype, the top three key components with their phenotype squared cosine scores are shown on the top of the stacked bar plot and gene contribution scores for each of the components are shown as colored segments. Each colored segment represents a gene with at least 0.05% of contribution scores and the rest of the genes are aggregated as the gray bar at the top. For each component, the labels for the top 10 driving genes are shown.

1208 Fig. S23: Squared cosine score of BMI (PTVs dataset)

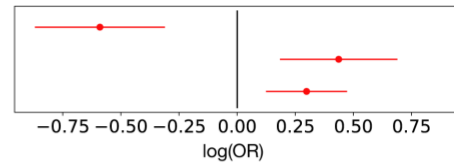
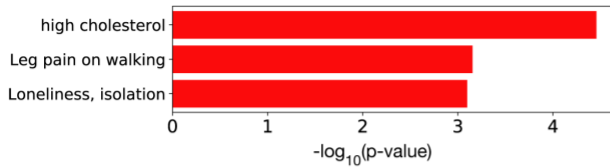


1209

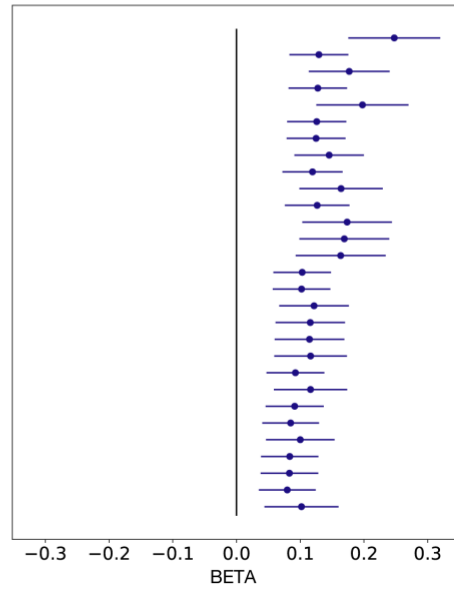
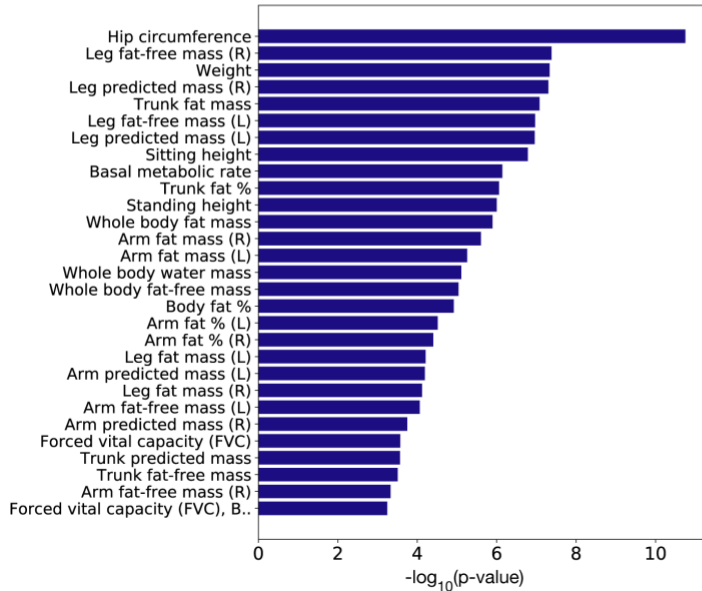
1210 **Fig. S23** Identification of the key components for BMI with phenotype squared cosine scores
1211 using the PTVs dataset. The top five key components are shown on the horizontal axis and the
1212 corresponding squared cosine scores are shown on the vertical axis.

1213 Fig. S24: PheWAS analysis for *PDE3B*

a PheWAS analysis of rs150090666 (*PDE3B*) for binary phenotypes



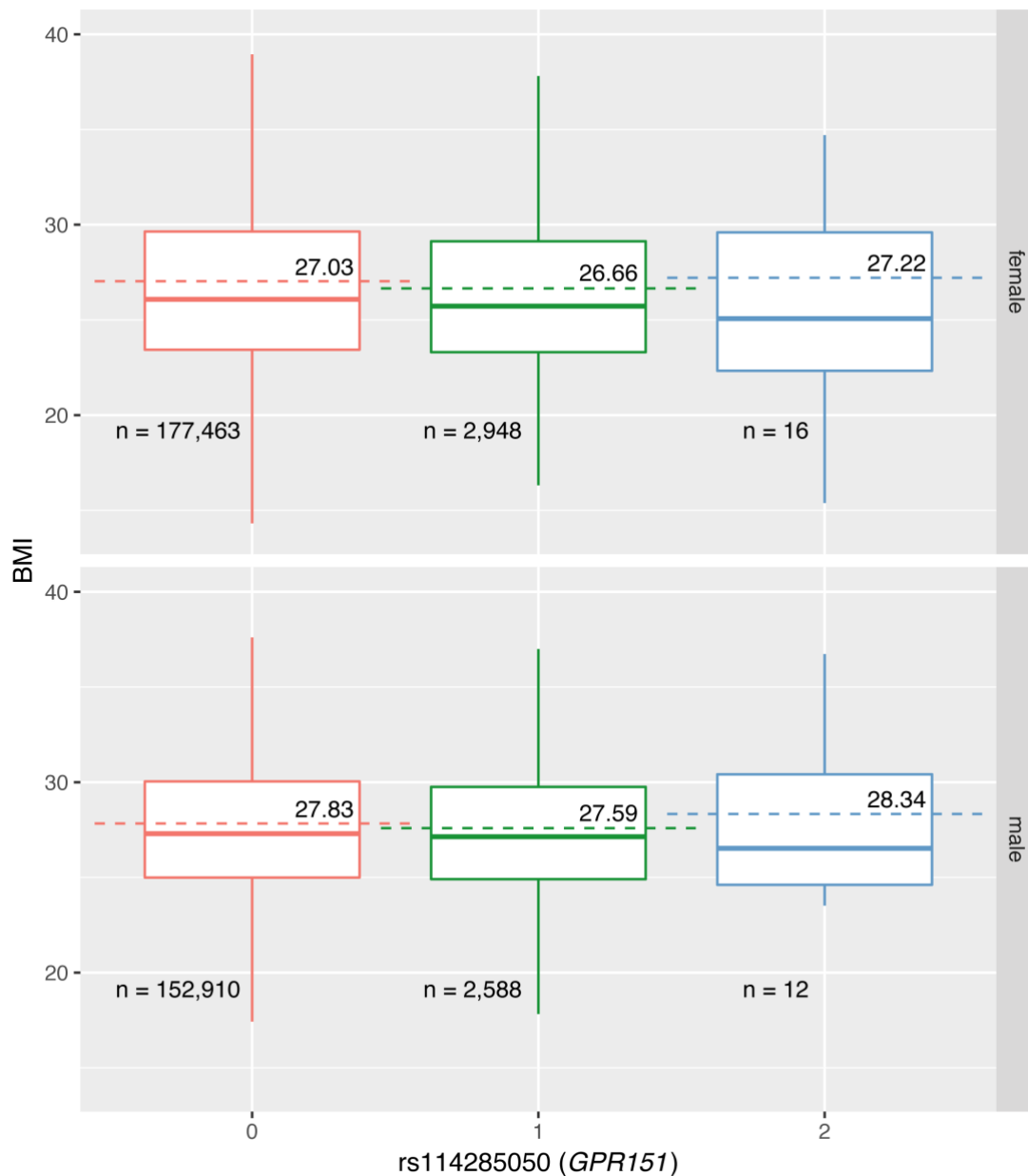
b PheWAS analysis of rs150090666 (*PDE3B*) for quantitative phenotypes



1214

1215 **Fig. S24** Phenome-wide association (PheWAS) analysis for rs150090666, a stop-gain variant in
1216 *PDE3B*. The p-values (left) and log odds ratio (binary phenotypes, shown as red)
1217 (quantitative phenotypes, shown as blue) (right) along with 95% confidence interval are shown
1218 for the phenotypes with minimum case count of 1,000 (binary phenotypes, **a**) or 1,000
1219 individuals with non-missing values (quantitative phenotypes, **b**) and strong association ($p \leq$
1220 0.001) and with this variants among all the phenotypes used in the study.

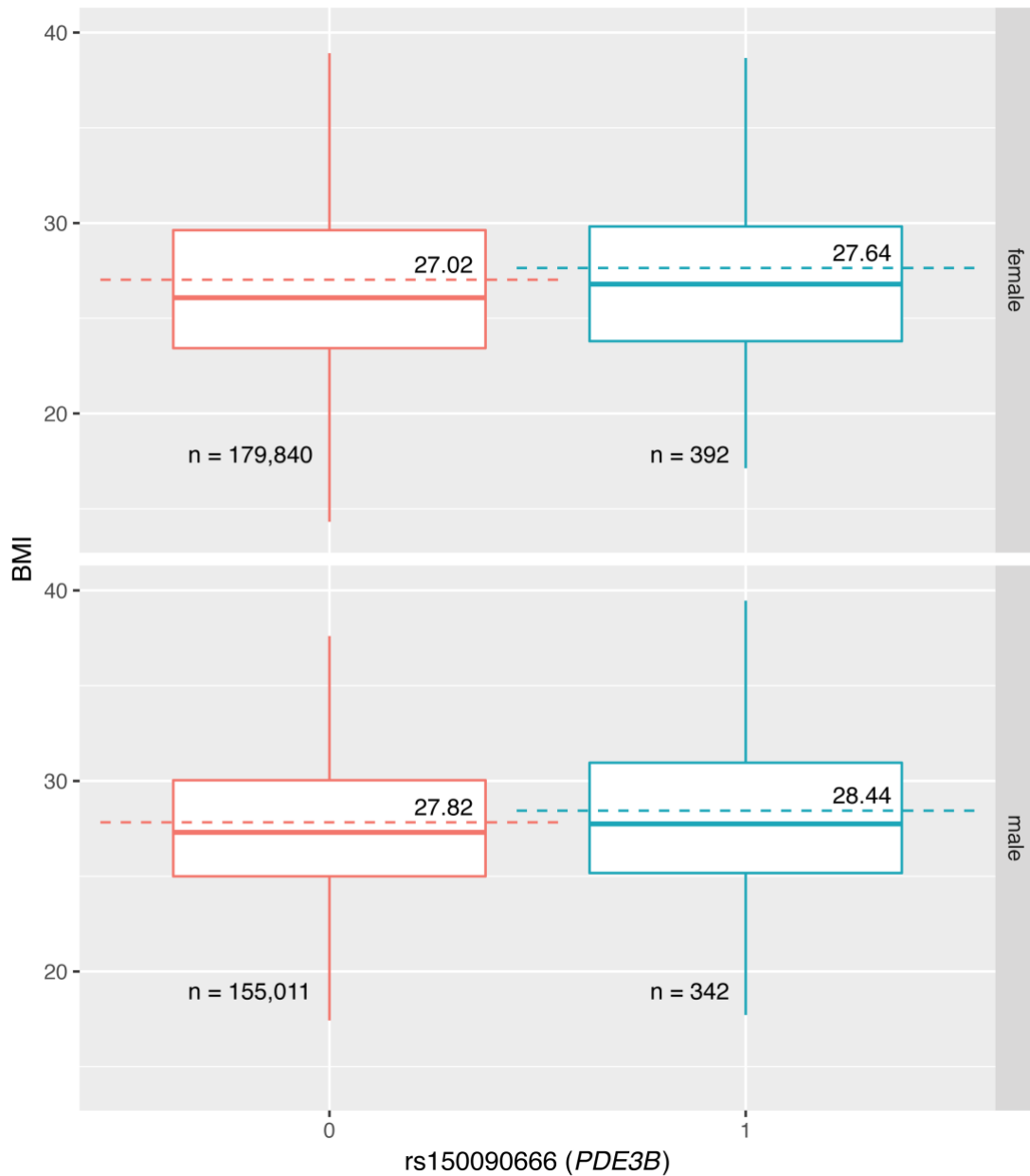
1221 Fig. S25: Univariate regression analysis for *GPR151*



1222
1223

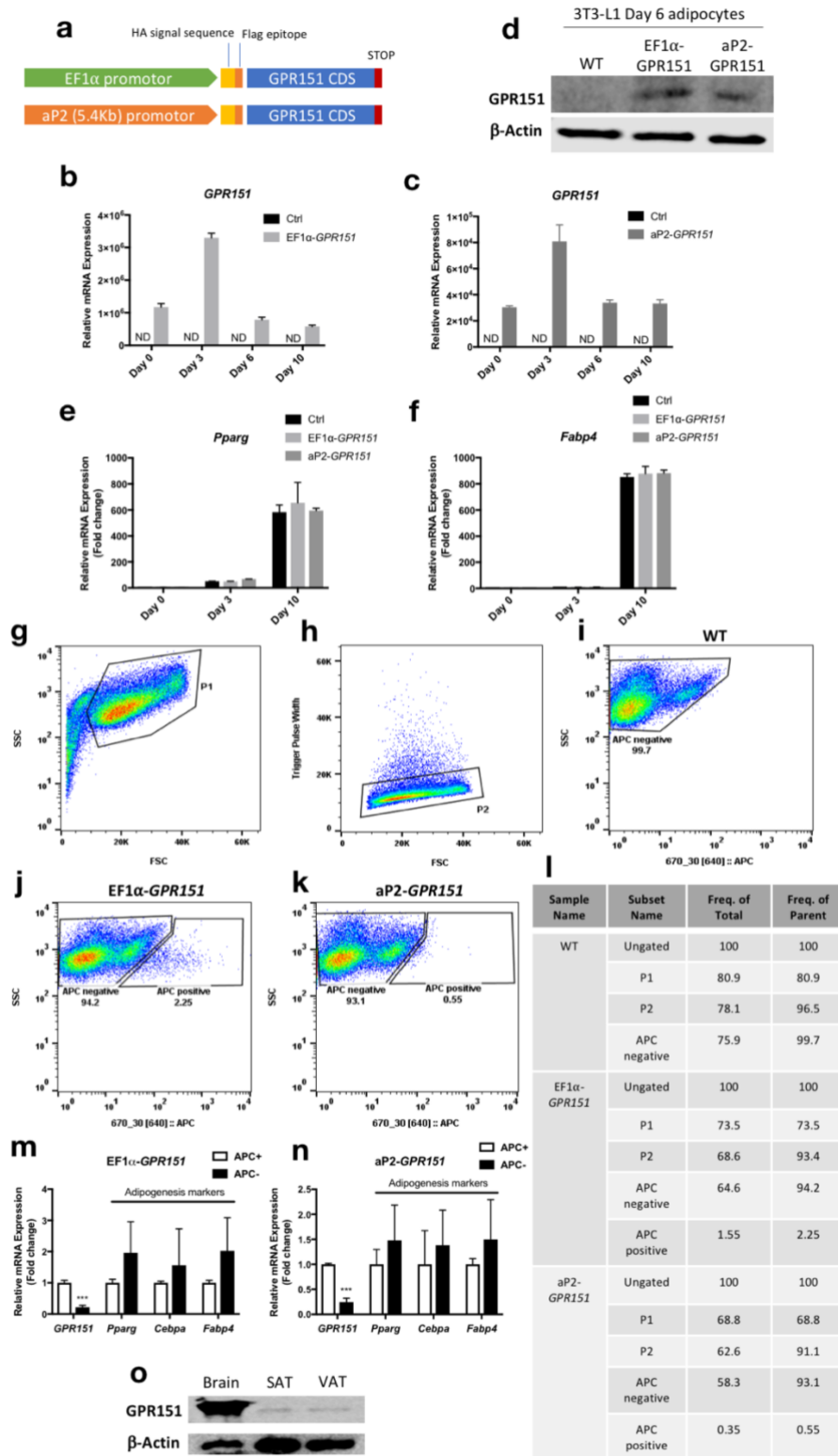
1224 **Fig. S25** Distribution of BMI stratified by sex and genotype of rs114285050, a stop-gain variant
1225 in *GPR151*. The outliers are removed from the plot and the mean values are annotated and
1226 shown as dashed lines. In the box plots, the median, two hinges (the first and the third quartiles)
1227 and two whiskers are shown. The upper whisker extends from the hinge to the largest value no
1228 further than $1.5 \times \text{IQR}$ from the hinge (where IQR is the inter-quartile range, or distance between
1229 the first and third quartiles). The number of carriers of the variants are shown at the bottom.

1230 Fig. S26: Univariate regression analysis for *PDE3B*



1231
1232 **Fig. S26** Distribution of BMI stratified by sex and genotype of rs150090666, a stop-gain variant
1233 in *PDE3B*. The outliers are removed from the plot and the mean values are annotated and
1234 shown as dashed lines. In the box plots, the median, two hinges (the first and the third quartiles)
1235 and two whiskers are shown. The upper whisker extends from the hinge to the largest value no
1236 further than 1.5 * IQR from the hinge (where IQR is the inter-quartile range, or distance between
1237 the first and third quartiles). The number of carriers of the variants are shown at the bottom.

1238 Fig. S27: *GPR151* overexpression

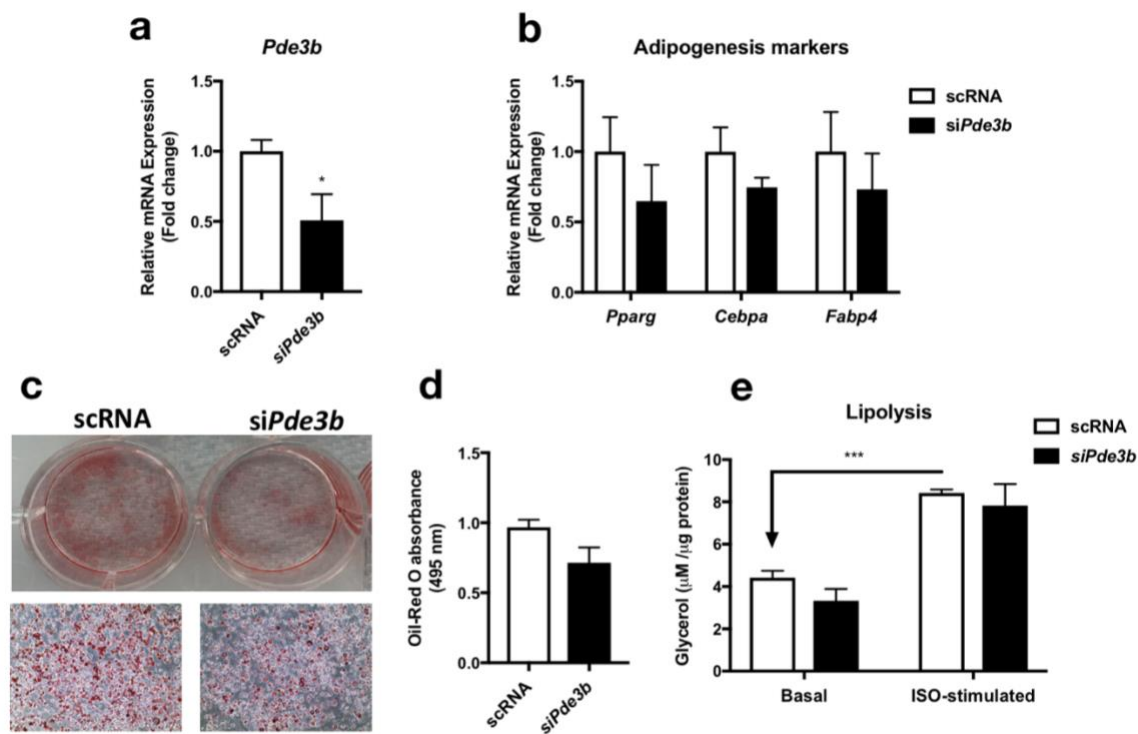


1239
1240
1241

Fig. S27 Effects of *GPR151* overexpression on 3T3-L1 adipogenesis. **a** Structure of *GPR151* overexpression construct driven by either EF1α or aP2 promoter. **b-d** Confirmation of *GPR151*

1242 overexpression at both mRNA (**b-c**) and protein levels (**d**) in 3T3-L1 cells during adipogenesis.
 1243 **e-f** qPCR analysis of the effect of *GPR151* overexpression on adipogenesis markers, *Pparg* (**e**)
 1244 and *Fabp4* (**f**). **g-i** Representative FACS gating strategy used to obtain APC+ and APC-
 1245 adipocytes. Cells were initially selected by size, on the basis of forward scatter (FSC) and side
 1246 scatter (SSC) (**g**). Cells were then gated on both FSC and SSC singlets to ensure that individual
 1247 cells were analyzed (**h**). Non-infected Day 6 3T3-L1 wild-type (WT) adipocytes were used to
 1248 determine background fluorescence levels (**i**). **j-l** Representative FACS collection gates used to
 1249 sort Day 6 3T3-L1 adipocytes infected with either EF1 α -*GPR151* (**j**) or aP2-*GPR151* (**k**) (shown
 1250 as APC positive), in comparison to WT (shown as APC negative). The abundance of the
 1251 relevant cell population in post-sort fractions were listed in **l**. **m-n** Relative mRNA levels of
 1252 *GPR151* and adipogenic markers (*Pparg*, *Cebpa*, *Fabp4*) in purified APC+ and APC- cells from
 1253 Day 6 3T3-L1 adipocytes infected by either EF1 α -*GPR151* (**m**) or aP2-*GPR151* (**n**). **o**
 1254 Comparison of protein levels of GPR151 in mouse brain, subcutaneous adipose tissue (SAT)
 1255 and visceral adipose tissue (VAT). For bar plots, means \pm SEM are shown. ND: not-detectable.

1256 Fig. S28: *Pde3b* knockdown



1257 **Fig. S28** Effects of *Pde3b* knockdown in 3T3-L1 adipogenesis. **a** qPCR analysis of *Pde3b*
 1258 mRNA knockdown in 3T3-L1 preadipocytes. **b** qPCR analysis of the effect of si*Pde3b*
 1259 knockdown on adipogenesis markers, *Pparg*, *Cebpa* and *Fabp4*. **c-d** Oil-Red O staining (**c**) and
 1260 quantification (**d**) of lipid droplets in scRNA- or si*Pde3b*-transfected adipocytes. **e** lipolysis
 1261 assays of scRNA- or si*Pde3b*-transfected adipocytes. Means \pm SEM are shown (**p-
 1262 value<0.001, *p-value<0.05). scRNA: scrambled siRNA. ISO: isoproterenol.
 1263

1264 Table S1 List of phenotype categories

1265 List of phenotype categories used in our study and their data source are shown with one
1266 example phenotype per category. Abbreviation in the type column. B: binary, Q: quantitative, P:
1267 described in previously published literature, F: the UK Biobank data field ID, and C: the UK
1268 Biobank data category ID.

1269 Table S2 List of phenotypes

1270 The list of phenotypes considered in the study. The table is sorted by category, number cases
1271 (for binary phenotypes), and the number of non-missing values (for quantitative phenotypes).
1272 The two columns, "All", "Coding", and "PTVs" indicates whether the phenotype is used in each
1273 of the dataset after imposing the filters on the genome-and phenome-wide summary statistics
1274 matrix. One can browse the summary statistics from genome-wide association studies on the
1275 Global Biobank Engine with the URL in the table.

1276 Table S3: Phenotype groupings for visualization

1277 The list of phenotype groups used in the phenotype contribution score plots are summarized.

1278 Table S4: Summary of contribution scores for the key 1279 components

1280 The list of top 20 driving phenotypes, genes, and variants for the first five principal components
1281 and the top three key components for the phenotypes highlighted in the study are summarized
1282 in the table.

1283 Table S5: GREAT enrichment analysis for BMI

1284 Biological characterization of driving non-coding and coding variants of the key components for
1285 BMI with the genomic region enrichment analysis tool (GREAT) using the all variants dataset.
1286 The results of the enrichment analysis for MGI phenotype ontology, a manually curated
1287 genotype-phenotype relationship knowledgebase for mouse, is summarized by the key
1288 components. The two major summary statistics from GREAT, binomial fold and binomial p-
1289 value, are shown. Abbreviation. BFold: binomial fold, BPval: binomial p-value.

1290 Table S6: GREAT enrichment analysis for MI

1291 Biological characterization of driving non-coding and coding variants of the key components for
1292 MI with the genomic region enrichment analysis tool (GREAT) using the all variants dataset.
1293 The results of the enrichment analysis for MGI phenotype ontology, a manually curated
1294 genotype-phenotype relationship knowledgebase for mouse, is summarized by the key

1295 components. The two major summary statistics from GREAT, binomial fold and binomial p-
1296 value, are shown. Abbreviation. BFold: binomial fold, BPval: binomial p-value.

1297 **Table S7: GREAT enrichment analysis for gallstones**

1298 Biological characterization of driving non-coding and coding variants of the key components for
1299 gallstones with the genomic region enrichment analysis tool (GREAT) using the all variants
1300 dataset. The results of the enrichment analysis for MGI phenotype ontology, a manually curated
1301 genotype-phenotype relationship knowledgebase for mouse, is summarized by the key
1302 components. The two major summary statistics from GREAT, binomial fold and binomial p-
1303 value, are shown. Abbreviation. BFold: binomial fold, BPval: binomial p-value.

1304 **Table S8: PheWAS analysis for rs114285050 (*GPR151*)**

1305 Phenome-wide association (PheWAS) analysis for rs114285050, a stop-gain variant in
1306 *GPR151*.

1307 **Table S9: PheWAS analysis for rs150090666 (*PDE3B*)**

1308 Phenome-wide association (PheWAS) analysis for rs150090666, a stop-gain variant in *PDE3B*.
1309

1310 **Table S10: Genetic correlation of summary statistics for 10 traits 1311 with different GWAS covariates**

1312 For five binary traits and five quantitative traits, genetic correlation is computed for two GWAS
1313 summary statistics computed with four and ten genotype principal components in the covariates.
1314