1 **Human knockouts and phenotypic analysis**

2 **in a cohort with a high rate of consanguinity**

3

4 Danish Saleheen[1,2]†*, Pradeep Natarajan[3,4]†, Irina Armean[4,5], Wei Zhao[1], Asif Rasheed[2],

5 Sumeet Khetarpal[6], Hong-Hee Won[7], Konrad J. Karczewski[4,5], Anne H. O'Donnell-

6 Luria[4,5,8], Kaitlin E. Samocha[4,5], Namrata Gupta[4], Mozzam Zaidi[2], Maria Samuel[2], Atif

7 Imran[2], Shahid Abbas[9], Faisal Majeed[2], Madiha Ishaq[2], Saba Akhtar[2], Kevin Trindade[6],

8 Megan Mucksavage[6], Nadeem Qamar[10], Khan Shah Zaman[10], Zia Yaqoob[10], Tahir

9 Saghir[10], Syed Nadeem Hasan Rizvi[10], Anis Memon[10], Nadeem Hayyat Mallick[11],

10 Mohammad Ishaq[12], Syed Zahed Rasheed[12], Fazal-ur-Rehman Memon[13], Khalid

11 Mahmood[14], Naveeduddin Ahmed[15], Ron Do[16,17], Ronald M. Krauss[18], Daniel G.

12 MacArthur[4,5], Stacey Gabriel[4], Eric S. Lander[4], Mark J. Daly[4,5], Philippe Frossard[2]†, John

13 Danesh[19,20]†, Daniel J. Rader[6,21]†, Sekar Kathiresan[3,4]†*

14

15 †Contributed equally

16

17 [1] Department of Biostatistics and Epidemiology, Perelman School of Medicine at the

18 University of Pennsylvania, Philadelphia, PA, USA

19 [2] Center for Non-Communicable Diseases, Karachi, Pakistan

20 [3] Center for Human Genetic Research and Cardiovascular Research Center,

21 Massachusetts General Hospital, Boston, MA, USA

22 [4] Broad Institute of Harvard and MIT, Cambridge, MA, USA

23  [5] Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts

24  General Hospital and Harvard Medical School, Boston, MA

25  [6] Division of Translational Medicine and Human Genetics, Department of Medicine,

26  Perelman School of Medicine at the University of Pennsylvania, Philadelphia, PA, USA

27  [7] Samsung Advanced Institute for Health Sciences and Technology (SAIHST),

28  Sungkyunkwan University, Samsung Medical Center, Seoul, Korea

29  [8] Division of Genetics and Genomics, Boston Children's Hospital, Boston, MA, USA

30  [9] Faisalabad Institute of Cardiology, Faisalabad, Pakistan

31  [10] National Institute of Cardiovascular Disorders, Karachi, Pakistan

32  [11] Punjab Institute of Cardiology, Lahore, Pakistan

33  [12] Karachi Institute of Heart Diseases, Karachi, Pakistan

34  [13] Red Crescent Institute of Cardiology, Hyderabad, Pakistan

35  [14] The Civil Hospital, Karachi, Pakistan

36  [15] Liaquat National Hospital, Karachi, Pakistan

37  [16] Department of Genetics and Genomic Sciences, Mount Sinai Medical Center, Icahn

38  School of Medicine at Mount Sinai, New York, NY, USA

39  [17] The Charles Bronfman Institute of Personalized Medicine, Icahn School of Medicine at

40  Mount Sinai, New York, NY, USA

41  [18] Children's Hospital Oakland Research Institute, Oakland, CA, USA

42  [19] Department of Public Health and Primary Care, University of Cambridge, UK

43  [20] Wellcome Trust Sanger Institute, Hinxton, Cambridge, UK

44  [21] Department of Human Genetics, University of Pennsylvania, USA

45

46    *Corresponding authors:

47    Danish Saleheen, MBBS, PhD

48    Department of Biostatistics and Epidemiology

49    University of Pennsylvania

50    11-134 Translational Research Center

51    3400 Civic Center Boulevard

52    Philadelphia, PA 19104

53    Tel: 215-573-6323

54    Fax: 215-573-2094

55    Email: saleheen@mail.med.upenn.edu

56

57    Sekar Kathiresan, MD

58    Broad Institute and Massachusetts General Hospital

59    CPZN 5.830

60    185 Cambridge Street

61    Boston, MA 02114

62    Tel: 617-643-6120

63    Email: sekar@broadinstitute.org

64

65    Word Count, Summary Paragraph: 313

66    Word Count, Main Text: 3,123

67 **Summary Paragraph**

68 A major goal of biomedicine is to understand the function of every gene in the human

69 genome.[1] Loss-of-function (LoF) mutations can disrupt both copies of a given gene in

70 humans and phenotypic analysis of such 'human knockouts' can provide insight into gene

71 function. To date, comprehensive analysis of genes knocked out in humans has been

72 limited by the fact that LoF mutations are infrequent in the general population and so,

73 observing an individual homozygous LoF for a given gene is exceedingly rare.[2,3]

74 However, consanguineous unions are more likely to result in offspring who carry LoF

75 mutations in a homozygous state. In Pakistan, consanguinity rates are notably high.[4]

76 Here, in order to understand consequences of complete gene disruption in humans, we

77 sequenced the protein-coding regions of 10,503 adult participants living in Pakistan,

78 identified individuals carrying predicted LoF (pLoF) mutations in the homozygous state,

79 and performed phenotypic analysis involving >200 traits. We enumerated 49,138 rare (<1

80 % minor allele frequency) pLoF mutations. These pLoF mutations are predicted to knock

81 out 1,317 genes in at least one participant. Homozygosity for pLoF mutations at *PLAG27*

82 was associated with absent enzymatic activity of soluble lipoprotein-associated

83 phospholipase A2; at *CYP2F1*, with higher plasma interleukin-8 concentrations; at

84 *TREH*, with lower concentrations of apoB-containing lipoprotein subfractions; at either

85 *A3GALT2* or *NRG4*, with markedly reduced plasma insulin C-peptide concentrations; and

86 at *SLC9A3R1*, with mediators of calcium and phosphate signaling. Finally, *APOC3* is a

87 gene which regulates metabolism of plasma triglyceride-rich lipoproteins and where

88 heterozygous deficiency confers resistance to coronary heart disease.[5,6] In Pakistan, we

89 now observe *APOC3* homozygous pLoF carriers; we recalled these knockout humans and

90    challenged with an oral fat load. Compared with wild-type family members, *APOC3*

91    knockouts displayed marked blunting of the usual post-prandial rise in plasma

92    triglycerides. Overall, these observations provide a roadmap for a 'human knockout

93    project', a systematic effort to understand the phenotypic consequences of complete

94    disruption of genes in humans.

**Main Text**

We studied adult participants in the Pakistan Risk of Myocardial Infarction Study (PROMIS) designed to understand the determinants of cardiometabolic diseases in South Asians.[7] Consanguineous marriages have been common in this region of South Asia for many generations.[8] In PROMIS, 39.0% of participants reported that their parents were cousins and 39.8% reported themselves being married to a cousin. An expectation from consanguinity is long regions of autozygosity, defined as homozygous loci identical by descent.[9] Using genome-wide genotyping data available in 18,541 PROMIS participants, we quantified the length of runs of homozygosity, defined as homozygous segments at least 1.5 megabases long. We compared the lengths of runs of homozygosity among PROMIS participants with those seen in other populations from the International HapMap3 Project. Median length of genome-wide homozygosity among PROMIS participants was 6-7 times higher than participants of European (CEU, TSI) ($P = 3.6$ x $10^{-37}$), East Asian (CHB, JPT, CHD) ($P = 5.4$ x $10^{-48}$) and African ancestries (YRI, MKK) ($P = 1.3$ x $10^{-40}$), respectively (**Supplementary Figure 1**).

In order to identify individuals who are homozygous for predicted loss-of-function (pLoF) mutations (i.e., nonsense, frameshift, or canonical splice-site mutations predicted to inactivate a gene), we performed whole exome sequencing in 10,503 PROMIS participants (**Table 1**) with genetic ancestry similar to the overall cohort. Across all participants, 1,639,223 exonic and splice-site sequence variants in 19,026 autosomal genes passed quality control metrics. Of these, 57,137 mutations across 14,345 autosomal genes were annotated as pLoF.

117     To increase the probability that mutations annotated as pLoF by automated

118     algorithms are *bona fide*, we removed nonsense and frameshift mutations occurring

119     within the last 5% of the transcript and within exons flanked by non-canonical splice

120     sites, splice site mutations at small (<15 bp) introns, at non-canonical splice sites, and

121     where the purported pLoF allele is observed across primates. Common pLoF alleles are

122     less likely to exert strong functional effects as they are less constrained by purifying

123     selection; thus, we define pLoF mutations in the rest of the manuscript as variants with a

124     minor allele frequency (MAF) of < 1% and passing the aforementioned bioinformatic

125     filters. Applying these criteria, we generated a set of 49,138 pLoF mutations across

126     13,074 autosomal genes.[10] The site-frequency spectrum for these pLoF mutations

127     revealed that the majority was seen only in one or a few individuals (**Supplementary**

128     **Figure 2**).

129     Across all 10,503 PROMIS participants, both copies of 1,317 distinct genes were

130     predicted to be inactivated due to pLoF mutations. A full listing of all 1,317 genes

131     knocked out, the number of knockout participants for each gene, and the specific pLoF

132     mutation(s) are provided in **Supplementary Table 1**. 891 (67.7 %) of the genes were

133     knocked out only in one participant (**Fig. 1a**). Nearly 1 in 5 sequenced participants (1,843

134     individuals, 17.5 %) had at least one gene knocked by a homozygous pLoF mutation.

135     1,504 of these 1,843 individuals (81.6 %) were homozygous pLoF carriers for just one

136     gene, but a minority of participants were knockouts for more than one gene and one

137     participant had six genes with homozygous pLoF genotypes.

138     We compared the coefficient of inbreeding (F coefficient) in PROMIS

139     participants with that of 15,249 individuals from outbred populations of European or

140    African American ancestry. The F coefficient estimates the excess homozygosity

141    compared with an estimated outbred ancestor. PROMIS participants had a 4-fold higher

142    median inbreeding coefficient compared to outbred populations (0.016 v 0.0041; $P < 2$ x

143    $10^{-16}$) (**Fig. 1b**). Additionally, those in PROMIS who reported that their parents were

144    closely related had even higher median inbreeding coefficients than those who did not

145    (0.023 v 0.013; $P < 2$ x $10^{-16}$). The F inbreeding coefficient was correlated with the

146    number of homozygous pLoF genes present in each individual. (Spearman r = 0.31; $P = 5$

147    x $10^{-231}$) (**Fig. 1c**). When restricted to individuals with high levels of inbreeding (F

148    inbreeding coefficient > 6.25%, the expected degree of autozygosity from a first-cousin

149    union), 721 of 1,585 individuals (45%) were homozygous for at least one pLoF mutation.

150         We tested the hypothesis that genes observed in the homozygous pLoF state in

151    PROMIS participants are under less evolutionary constraint. We calculated the

152    probability of being LoF intolerant (at >90% threshold) for each gene (see Methods) [11,12]

153    and compared this to 1,317 randomly selected genes. The observed 1,317 homozygous

154    pLoF genes were less likely to be classified as highly constrained (odds ratio 0.14; 95%

155    CI 0.12, 0.16; $P < 1$ x $10^{-10}$). Additionally, the 1,317 homozygous pLoF genes are

156    substantially depleted of genes described to be essential for survival and proliferation in

157    four human cancer cell lines (12 of 870 essential genes observed, 1.4%).[13]

158         A number of genes previously predicted to be required for viability in humans

159    were observed in the homozygous pLoF state in humans (**Supplementary Table 2**). For

160    example, 40 of the 1,317 genes have been associated with embryonic or perinatal

161    lethality as homozygous pLoF in mice.[14] Furthermore, 56 genes predicted to be essential

162    using mouse/human conservation data[15] are tolerated as homozygous pLoF in Pakistani

163 adults. In fact, 9 genes are in both datasets and are also modeled as LoF intolerant.[12] One

164 such gene, *EP400* (also known as *p400*), influences cell cycle regulation via chromatin

165 remodeling[16] and is critical for maintaining the identity of murine embryonic stem cells[17]

166 but we observe an adult human homozygous for disruption of a canonical splice site

167 (intron 3 of 52; c.1435+1G>A) in *EP400*. Conversely, we observed 90 genes where the

168 heterozygous pLoF genotype is of appreciable frequency but the homozygous pLoF

169 genotype is depleted (at *P* value threshold < 0.05) (**Supplementary Table 3**).

170       We compared our results to three recent reports where homozygous pLoF genes

171 have been catalogued: in Pakistanis living in Britain, in Icelanders, and in the Exome

172 Aggregation Consortium (ExAC). 3,223 Pakistanis living in Britain with a high degree of

173 parental relatedness (mean 5.62% autozygosity) were sequenced to find 781 homozygous

174 pLoF genes.[18] The sequencing of 2,636 Icelanders and subsequent imputation into

175 104,220 chip-genotyped Icelanders yielded 1,171 genes in the homozygous pLoF state.[3]

176 Analysis of 52,451 multi-ethnic participants from ExAC (i.e., those not overlapping with

177 current PROMIS study) found 877 genes to be knocked out.[19] Here, we identify a total of

178 734 unique genes in the homozygous pLoF state that were not observed in the other three

179 studies (**Supplementary Figure 3**).

180       Intersection of the four sets of genes from these studies revealed only 25 common

181 to all four studies. For example, at phosphodiesteriase 11A (encoded by *PDE11A*),

182 different mutations across the four populations lead to homozygous pLoF state

183 (PROMIS: c.2424-1G>G, p.Cys554ValfsTer14, p.Arg307Ter; ExAC Latino:

184 p.Arg307Ter; ExAC non-Finnish European: p.Cys554ValfsTer14, p.Arg307Ter;

185 Icelanders: p.Arg7ThrfsTer30, p.Arg307Ter; British Pakistani: p.Arg57Ter). The *Pde11a⁻*

186    [/-] mouse shows behavioral phenotypes and *PDE11A* is implicated in depression and

187    schizophrenia in humans.[20] Whether humans lacking *PDE11A* also display

188    neuropsychiatric phenotypes remains to be determined.

189        In order to understand the phenotypic consequences of complete disruption of the

190    1,317 pLoF genes identified in the PROMIS study, we applied three approaches. First,

191    for 426 genes where two or more participants were homozygous pLoF, we conducted an

192    association screen against a panel of 201 phenotypic traits (**Supplementary Table 4**).

193    Second, in blood samples from each of 84 participants, we measured 1,310 protein

194    biomarkers using a new, multiplexed, aptamer-based proteomics assay. Third, at a single

195    gene, apolipoprotein C-III (encoded by *APOC3*), we recalled participants based on

196    genotype (three classes: 'wild-type', heterozygous pLoF, and homozygous pLoF) and

197    performed provocative physiologic testing.

198        At 426 genes where two or more participants were homozygous pLoF, we

199    performed association analyses to determine whether homozygous pLoF mutation status

200    was associated with variation in any of 201 traits. For quantitative traits, we compared

201    mean trait values in homozygous pLoF carriers with non-carriers. For dichotomous traits,

202    we performed logistic regression with trait status as the outcome variable and

203    homozygous pLoF carrier status as the predictor variable. Details of covariate

204    adjustments are presented in the Methods. Across quantitative and dichotomous traits,

205    this resulted in the analysis of 18,959 gene-trait pairs and thus, we set Bonferroni-

206    adjusted significance threshold at $P = 3 \times 10^{-6}$.

207        The quantile-quantile plot of expected versus observed association results shows

208    an excess of highly significant results without systematic inflation (**Supplementary**

209    **Figure 4**). Association results surpassed the Bonferroni significance threshold for 26

210    gene-trait pairs (**Supplementary Table 5**). Below, we highlight seven results: *PLA2G7*,

211    *CYP2F1*, *TREH*, *A3GALT2*, *NRG4*, *SLC9A3R1*, and *APOC3*.

212            Lipoprotein-associated phospholipase A2 (Lp-PLA2, encoded by *PLA2G7*)

213    hydrolyzes phospholipids to generate lysophosphatidylcholine and oxidized nonesterified

214    fatty acids. In observational epidemiologic studies, higher soluble Lp-PLA2 enzymatic

215    activity has been correlated with increased risk for coronary heart disease; small molecule

216    inhibitors of Lp-PLA2 have been developed for the treatment of coronary heart disease.[21]

217    In PROMIS, we identified participants who are naturally deficient in the Lp-PLA2

218    enzyme. Two participants are homozygous for a splice-site mutation, *PLA2G7*

219    c.663+1G>A, and 106 are heterozygous for this same mutation. We observed a dose-

220    dependent response relationship between genotype and enzymatic activity: when

221    compared with non-carriers, c.663+1G>A homozygotes have markedly lower Lp-PLA2

222    enzymatic activity (-245 nmol/ml/min, $P = 2 \times 10^{-7}$) whereas the 106 heterozygotes had

223    an intermediate effect (-120 nmol/ml/min, $P = 2 \times 10^{-77}$) (**Fig. 2a-b**). If Lp-PLA2 plays a

224    causal role for coronary heart disease, one might expect those naturally deficient for this

225    enzyme to have reduced risk for coronary heart disease. We tested the association of

226    *PLA2G7* c.663+1G>A with myocardial infarction across all participants and found that

227    carriers of the pLoF allele did not have reduced risk (OR 0.97; 95% CI, 0.70 – 1.34; $P =$

228    0.87) (**Fig. 2c**). In contrast, at two positive control genes, we replicated prior observations

229    (**Supplementary Table 6**); at *LDLR*, heterozygous pLoF mutations increased MI risk 20-

230    fold and, at *PCSK9*, heterozygous pLoF mutations reduced risk by 78%. Of note, in two

231    recent randomized controlled trials, pharmacologic Lp-PLA2 inhibition failed to reduce

232     risk for coronary heart disease,[22,23] a result that might have been anticipated by this

233     genetic analysis.[24]

234        Cytochrome P450 2F1 (encoded by *CYP2F1*) is primarily expressed in the lung

235     and metabolizes pulmonary-selective toxins, such as cigarette smoke, and thus,

236     modulates the expression of environment-associated pulmonary diseases.[25] At *CYP2F1*,

237     we identified two participants homozygous for a splice-site mutation, c.1295-2A>G.

238     When compared with non-carriers, c.1295-2A>G homozygotes displayed higher soluble

239     interleukin 8 concentrations (3.7-fold increase, $P = 2 \times 10^{-6}$) (**Supplementary Figure 5**).

240     *CYP2F1* c.1295-2A>G heterozygosity had a more modest effect (2.4-fold increase, $P = 2$

241     $\times 10^{-4}$). Interleukin 8 induces migration of neutrophils in airways and is a mediator of

242     acute pulmonary inflammation and chronic obstructive pulmonary disease (COPD).[26,27]

243     However, neither carrier reports a personal or family history of obstructive pulmonary

244     disease; further studies of these participants are required to assess the roles of CYP2F1

245     and interleukin 8 on pulmonary physiology.

246        Trehalase (encoded by *TREH*) is an intestinal enzyme that splits the naturally-

247     found unabsorbed disaccharide, trehalose, into two glucose molecules.[28] Trehalase

248     deficiency, an autosomal recessive trait, leads to abdominal pain, distention, and

249     flatulence after trehalose ingestion. We identified six participants homozygous for a

250     deletion of a splice acceptor site (c.90-

251     9_106delTCTCTGCAGTGAGATTTACTGCCACG) in exon 2. Homozygotes, unlike

252     heterozygotes or non-carriers, had lower concentrations of several apolipoprotein B-

253     containing lipoprotein subfractions (**Supplementary Table 5**) (**Supplementary Figure

254     6**).

255　　　　　Alpha-1,3-galactosyltransferase 2 (encoded by *A3GALT2*) catalyzes the formation

256　　of the Gal-α1-3Galβ1-4GlcNAc-R (α-gal) epitope; the biological role of this enzyme in

257　　humans is uncertain.[29] At *A3GALT2*, we identified two participants homozygous for a

258　　frameshift mutation, p.Thr106SerfsTer4. Compared with non-carriers, p.Thr106SerfsTer4

259　　homozygotes both had dramatically reduced concentrations of fasting C-peptide (-97.4%;

260　　$P = 6 \times 10^{-12}$) and insulin (-92.3%; $P = 1 \times 10^{-4}$). Such an association was only observed

261　　in the homozygous state (**Supplementary Figure 7**). *A3galt2*$^{-/-}$ mice and pigs have

262　　recently been shown to have glucose intolerance.[30,31]

263　　　　　To understand if the identification of only a single homozygote may still be

264　　informative, we performed a complementary analysis, focusing on those with the most

265　　extreme standard Z scores (|Z score| > 5) and requiring that there be evidence for

266　　association in heterozygotes as well (see Methods). This procedure highlighted neureglin

267　　4 (*NRG4*), a member of the epidermal growth factor family extracellular ligands which is

268　　highly expressed in brown fat, particularly during adipocyte differentiation.[32,33] At *NRG4*,

269　　we identified a single participant homozygous for a frameshift mutation,

270　　p.Ile75AsnfsTer23, who had nearly absent fasting insulin C-peptide concentrations (-99.3

271　　%; $P = 1 \times 10^{-10}$). When compared with non-carriers, heterozygotes for *NRG4*

272　　p.Ile75AsnfsTer23 (n = 8) displayed 48.3 % reduction in insulin C-peptide ($P = 1 \times 10^{-2}$).

273　　Mice deleted for *Nrg4* have recently been shown to have glucose intolerance.[33] The

274　　single *NRG4* pLoF homozygote participant did not have diabetes nor elevated fasting

275　　glucose. Heterozygosity for a *NRG4* pLoF mutation (n=26) was also not associated with

276　　diabetes or fasting glucose. More detailed phenotyping will be required to definitively

277　　assess any relationship of *NRG4* deficiency in humans with glucose intolerance.

278    To further dissect the consequences of a subset of homozygous pLoF genes, we

279    measured 1,310 protein biomarkers in 84 participants through a new, multiplexed,

280    proteomic assay (SOMAscan). Among the 84 participants, there were nine genes with at

281    least two pLoF homozygotes and we associated these genotypes across 1,310 protein

282    biomarkers and observed a number of associations (**Supplementary Table 7**). We

283    highlight two PROMIS participants who are homozygous pLoF at *SLC9A3R1*; these

284    participants have increased circulating concentrations of several proteins involved in

285    parathyroid hormone or osteoclast signaling including calcium / calmodulin-dependent

286    protein kinase II (CAMK2) alpha, beta, and delta subunits, cAMP-regulated

287    phosphoprotein 19, and signal transducer and activator of transcription (STAT) 1, 3, and

288    6 (**Supplementary Table 7**). *SLC9A3R1* (aka *NHERF1*) encodes a Na+/H+ exchanger

289    regulatory cofactor that interacts with and regulates the parathyroid hormone receptor;

290    *Nherf1*[-/-] mice display hyperphosphaturia and disrupted protein kinase A-dependent

291    cAMP-mediated phosphorylation.[34,35] Humans carrying rare missense mutations in

292    *SLC9A3R1* have nephrolithiasis, osteoporosis, and hypophosphatemia.[36]

293    Apolipoprotein C-III (apoC-III, encoded by *APOC3*) is a major protein

294    component of chylomicrons, very low-density lipoprotein cholesterol, and high-density

295    lipoprotein cholesterol.[37] We and others recently reported that *APOC3* pLoF mutations in

296    heterozygous form lower plasma triglycerides and reduce risk for coronary heart

297    disease[5,6,38]; there is now substantial interest in *APOC3* as a therapeutic target.[39-41] In

298    published studies, no *APOC3* pLoF homozygotes have been identified despite study of

299    nearly 200,000 participants from the U.S. and Europe, raising concerns that complete

300    *APOC3* deficiency may be harmful. However, in our study of ~10,000 Pakistanis, we

301    identified four participants homozygous for *APOC3* p.Arg19Ter. When compared with

302    non-carriers, p.Arg19Ter homozygotes displayed near-absent plasma apoC-III protein (-

303    88.9 %, $P = 5 \times 10^{-23}$), lower plasma triglyceride concentrations (-59.6 %, $P = 7 \times 10^{-4}$),

304    higher high-density lipoprotein (HDL) cholesterol  (+26.9 mg/dL, $P = 3 \times 10^{-8}$); and

305    similar levels of low-density lipoprotein (LDL) cholesterol ($P = 0.14$) (**Fig. 3a-d**).

306          ApoC-III functions as a brake on the metabolism of dietary fat and thus, the

307    complete lack of this protein should promote handling of ingested fat. We re-contacted

308    one homozygous pLoF proband, his wife, and 27 of his first-degree relatives for

309    genotyping and physiologic investigation. We found that the proband's wife, a first

310    cousin, was also a pLoF homozygote, leading to all nine children being obligate

311    homozygotes (**Fig. 3e**). In this family, we challenged pLoF homozygotes ($APOC3^{-/-}$; n =

312    6) and non-carriers ($APOC3^{+/+}$; n = 7) with a 50 g/m$^2$ oral fat load followed by serial

313    blood testing for six hours. *APOC3* p.Arg19Ter homozygotes had significantly lower

314    post-prandial triglyceride excursions (triglycerides area under the curve 468.3 mg/dL*6

315    hours vs 1267.7 mg/dL*6 hours; $P = 1 \times 10^{-4}$) (**Fig. 3f**). These data show that complete

316    lack of apoC-III markedly improves clearance of plasma triglycerides after a fatty meal

317    and are consistent with and extend an earlier report of diminished post-prandial lipemia

318    in *APOC3* pLoF heterozygotes.[38]

319          Targeted gene disruption in model organisms followed by phenotypic analysis has

320    been a fruitful approach to understand gene function[42]; here, we extend this concept to

321    the human organism, leveraging naturally-occurring pLoF mutations, consanguinity, and

322    biochemical phenotyping. These results permit several conclusions. First, power to

323    identify human knockouts is improved with the study of multiple populations and

324   particularly those with high degrees of consanguinity. Using the observed median

325   inbreeding coefficient of sequenced participants and genotypes from the first 7,078

326   sequenced Pakistanis, we estimate that the sequencing of 200,000 Pakistanis, may result

327   in up to 8,754 genes (95% CI, 8,669-8,834) completely knocked out in at least one

328   participant (**Fig. 4**).

329   Second, a panel of phenotypes measured in a blood sample can yield hypotheses

330   regarding phenotypic consequences of gene disruption as observed for *PLA2G7*,

331   *CYP2F1*, *TREH*, *A3GALT2, NRG4*, *SLC9A3R1*, and *APOC3*. Finally, recall by genotype

332   followed by provocative testing may provide physiologic insights. We used this approach

333   to demonstrate that complete lack of apoC-III is tolerated and results in both lowered

334   fasting triglyceride concentrations as well as substantially blunted post-prandial lipemia.

335   Several limitations deserve mention. First and most importantly, any given

336   mutation annotated as pLoF may not truly lead to loss of protein function. In addition to

337   bioinformatics filtration, we manually curated all homozygous pLoF variants (n=1,580)

338   to assess confidence in variant fidelity and predicted biochemical impact

339   (**Supplementary Table 1** and **Supplementary Table 8**). We found 56 variants with

340   genotypes with a low number of supportive reads, 55 with poorly mapped reads

341   (**Supplementary Table 9**), and an additional 66 where there were potential mechanisms

342   of protein-truncation rescue (**Supplementary Figure 8**) or occurred within exons or

343   splice sites where conservation was low. Thus, we found the majority of pLoF calls (1403

344   out of 1580; 89%) to be free of mapping or annotation error. However, for any given

345   pLoF, experimental validation will be required to prove loss of gene function (e.g.,

346   targeted assays such as RT-PCR of transcript and/or Western blot of protein to confirm

347  its absence). Second, statistical power for genotype-phenotype correlation is low if a gene

348  is knocked-out in only 1 or 2 participants. However, this situation should improve with

349  larger sample sizes (**Supplementary Figure 9**). Third, statistical power in the proteomics

350  analysis may be low because of the limited number of samples assayed and the impact of

351  non-genetic factors on plasma concentrations.[43] Finally, our analysis was limited to

352  available phenotypes and in only one instance did we recall participants for deeper

353  phenotyping; rather, a standardized clinical phenotyping protocol is desirable for each

354  participant where a gene is observed to be knocked out.

355      To date, most human genetic studies have pursued a phenotype-first ("forward"

356  genetics) approach, beginning with traits of interest followed by genetic mapping. It is

357  now feasible to pursue a systematic genotype-first ("reverse" genetics) approach, starting

358  with homozygous pLoF humans followed by methodical examination of a diverse set of

359  traits.

360      These observations set the stage for a 'human knockout project,' a systematic

361  effort to understand the phenotypic consequences of complete disruption of every gene in

362  the human genome. Key elements for a human knockout project include: 1) identification

363  of populations where homozygous genotypes may be enriched[18,44]; 2) deep-coverage

364  sequencing of the protein-coding regions of the genome[3]; 3) availability of a broad array

365  biochemical as well as clinical phenotypes across the population; 4) ability to re-contact

366  knockout humans as well as family members; 5) a thorough clinical evaluation in each

367  participant where a gene is observed to be knocked out; and 6) hypothesis-driven

368  provocative phenotyping in selected participants.

369 **Methods**

370 **General overview of the Pakistan Risk for Myocardial Infarction Study (PROMIS).**

371 The PROMIS study was designed to investigate determinants of cardiometabolic diseases

372 in Pakistan. Since 2005, the study has enrolled close to 38,000 participants; the present

373 investigation sequenced 10,503 participants selected as 4,793 cases with myocardial

374 infarction and 5,710 controls free of myocardial infarction. Participants aged 30-80 years

375 were enrolled from nine recruitment centers based in five major urban cities in Pakistan.

376 Type 2 diabetes in the study was defined based on self-report or fasting glucose levels

377 >125 mg/dL or HbA1c > 6.5 % or use of glucose lowering medications. The institutional

378 review board at the Center for Non-Communicable Diseases (IRB: 00007048,

379 IORG0005843, FWAS00014490) approved the study and all participants gave informed

380 consent.

381

382 **Phenotype descriptions.**

383 Non-fasting blood samples (with the time since last meal recorded) were drawn and

384 centrifuged within 45 minutes of venipuncture. Serum, plasma and whole blood samples

385 were stored at -70°C within 45 minutes of venipuncture. All samples were transported on

386 dry ice to the central laboratory at the Center for Non-Communicable Diseases (CNCD),

387 Pakistan, where serum and plasma samples were aliquoted across 10 different storage

388 vials. Samples were stored at -70°C for any subsequent laboratory analyses. All

389 biochemical assays were conducted in automated auto-analyzers. At CNCD Pakistan,

390 measurements for total-cholesterol, HDL cholesterol, LDL cholesterol, triglycerides, and

391 creatinine were made in serum samples using enzymatic assays; whereas levels of HbA1c

392    were measured using a turbidemetric assay in whole-blood samples (Roche Diagnostics,

393    USA). For further measurements, aliquots of serum and plasma samples were transported

394    on dry ice to the Smilow Research Center, University of Pennsylvania, USA, where

395    following biochemical assays were conducted: apolipoproteins (apoA-I, apoA-II, apoB,

396    apoC-III, apoE) and non-esterified fatty acids were measured through

397    immunoturbidometric assays using kits by Roche Diagnostics or Kamiya; lipoprotein (a)

398    levels were determined through a turbidimetric assay using reagents and calibrators from

399    Denka Seiken (Niigata, Japan); LpPLA2 mass and activity levels were determined using

400    immunoassays manufactured by diaDexus (San Francisco, CA, USA); measurements for

401    insulin, leptin and adiponectin were made using radio-immunoassays by LINCO (MO,

402    USA); levels of adhesion molecules (ICAM-1, VCAM-1, P- and E-Selectin) were

403    determined through enzymatic assays by R&D (Minneapolis, MN, USA); and

404    measurements for C-reactive protein, alanine transaminase, aspartate transaminase,

405    cystatin-C, ferritin, ceruloplasmin, thyroid stimulating hormone, alkaline phosphatase,

406    sodium, potassium, choloride, phosphate, sex-harmone binding globulin were made using

407    enzymatic assays manufactured by Abbott Diagnostics (NJ, USA).  Glomerular filtration

408    rate (eGFR) was estimated from serum creatinine levels using the MDRD equation.

409    ApoC-III levels were determined in an autoanalyzer using a commercially available

410    ELISA by Sekisui Diagnostics (Lexington, USA). We also measured the following 52

411    protein biomarkers by multiplex immunoassay using a customised panel on the Luminex

412    100/200 instrument by RBM (Myriad Rules Based Medicine, Austin, TX, USA): fatty

413    acid binding protein, granuloctye monocyte colony stimulating factor, granulocyte colony

414    stimulating factor, interferon gamma, interleukin-1 beta, interleukin 1 receptor,

415    interleukin 2, interleukin 3, interleukin 4, interleukin 5, interleukin 6, interleukin 7,

416    interleukin 8, interleukin 10, interleukin 18, interleukin p40, interleukin p70, interleukin

417    15, interleukin 17, interleukin 23, macrophage inflammatory protein 1 alpha, macrophage

418    inflammatory protein 1 beta, malondialdehyde-modified LDL, matrix metalloproteinase

419    2, matrix metalloproteinase 3, matrix metalloproteinase 9, nerve growth factor beta,

420    tumor necrosis factor alpha, tumor necrosis factor beta, brain-derived neurotrophic factor,

421    CD40, CD40 ligand, eotaxin, factor VII, insulin-like growth factor 1, lecithin-type

422    oxidized LDL receptor 1, monocyte chemoattractant protein 1, myeloperoxidase, N-

423    terminal prohormone of brain natriuretic peptide, neuronal cell adhesion molecule,

424    pregnancy-associated plasma protein A, soluble receptor for advanced glycation end-

425    products, sortilin, stem cell factor, stromal cell-derived factor 1, thrombomodulin, S100

426    calcium binding protein B, and vascular endothelial growth factor.

427

428    **Laboratory methods for array-based genotyping.**

429    As previously described, a genomewide association scan was performed using the

430    Illumina 660 Quad array at the Wellcome Trust Sanger Institute (Hinxton, UK) and using

431    the Illumina HumanOmniExpress at Cambridge Genome Services, UK.[45] Initial quality

432    control (QC) criteria included removal of participants or single nucleotide

433    polymorphisms (SNPs) that had a missing rate >5%. SNPs with a MAF <1% and a P-

434    value of $<10^{-7}$ for the Hardy-Weinberg equilibrium test were also excluded from the

435    analyses. In PROMIS, further QC included removal of participants with discrepancy

436    between their reported sex and genetic sex determined from the X chromosome. To

437    identify sample duplications, unintentional use of related samples (cryptic relatedness)

438    and sample contamination (individuals who seem to be related to nearly everyone in the

439    sample), identity-by-descent (IBD) analyses were conducted in PLINK.[46]

440

441    **Laboratory methods for exome sequencing.**

442    **Exome sequencing.** Exome sequencing was performed at the Broad Institute.

443    Sequencing and exome capture methods have been previously described.[47,48] A brief

444    description of the methods is provided below.

445    **Receipt/quality control of sample DNA**. Samples were shipped to the Biological

446    Samples Platform laboratory at the Broad Institute of MIT and Harvard (Cambridge, MA,

447    USA). DNA concentration was determined by PicoGreen (Invitrogen; Carlsbad, CA,

448    USA) prior to storage in 2D-barcoded 0.75 ml Matrix tubes at  -20 $^{\circ}$C in the SmaRTStore

449    (RTS, Manchester, UK) automated sample handling system. Initial quality control (QC)

450    on all samples involving sample quantification (PicoGreen), confirmation of high-

451    molecular weight DNA and fingerprint genotyping and gender determination (Illumina

452    iSelect; Illumina; San Diego, CA, USA). Samples were excluded if the total mass,

453    concentration, integrity of DNA or quality of preliminary genotyping data was too low.

454    **Library construction.** Library construction was performed as previously described[49],

455    with the following modifications: initial genomic DNA input into shearing was reduced

456    from 3µg to 10-100ng in 50µL of solution. For adapter ligation, Illumina paired end

457    adapters were replaced with palindromic forked adapters, purchased from Integrated

458    DNA Technologies, with unique 8 base molecular barcode sequences included in the

459    adapter sequence to facilitate downstream pooling. With the exception of the palindromic

460    forked adapters, the reagents used for end repair, A-base addition, adapter ligation, and

461     library enrichment PCR were purchased from KAPA Biosciences (Wilmington, MA,

462     USA) in 96-reaction kits. In addition, during the post-enrichment SPRI cleanup, elution

463     volume was reduced to 20 µL to maximize library concentration, and a vortexing step

464     was added to maximize the amount of template eluted.

465     **In-solution hybrid selection.** 1,970 samples underwent in-solution hybrid selection as

466     previously described[49], with the following exception: prior to hybridization, two

467     normalized libraries were pooled together, yielding the same total volume and

468     concentration specified in the publication. 8,808 samples underwent hybridization and

469     capture using the relevant components of Illumina's Rapid Capture Exome Kit and

470     following the manufacturer's suggested protocol, with the following exceptions: first, all

471     libraries within a library construction plate were pooled prior to hybridization, and

472     second, the Midi plate from Illumina's Rapid Capture Exome Kit was replaced with a

473     skirted PCR plate to facilitate automation. All hybridization and capture steps were

474     automated on the Agilent Bravo liquid handling system.

475     **Preparation of libraries for cluster amplification and sequencing.** Following post-

476     capture enrichment, libraries were quantified using quantitative PCR (KAPA Biosystems)

477     with probes specific to the ends of the adapters. This assay was automated using

478     Agilent's Bravo liquid handling platform. Based on qPCR quantification, libraries were

479     normalized to 2nM and pooled by equal volume using the Hamilton Starlet. Pools were

480     then denatured using 0.1 N NaOH. Finally, denatured samples were diluted into strip

481     tubes using the Hamilton Starlet.

482     **Cluster amplification and sequencing.** Cluster amplification of denatured templates

483     was performed according to the manufacturer's protocol (Illumina) using HiSeq v3

484    cluster chemistry and HiSeq 2000 or 2500 flowcells. Flowcells were sequenced on HiSeq

485    2000 or 2500 using v3 Sequencing-by-Synthesis chemistry, then analyzed using RTA

486    v.1.12.4.2. Each pool of whole exome libraries was run on paired 76bp runs, with and 8

487    base index sequencing read was performed to read molecular indices, across the number

488    of lanes needed to meet coverage for all libraries in the pool.

489    **Read mapping and variant discovery**. Samples were processed from real-time base-

490    calls (RTA v.1.12.4.2 software [Bustard], converted to qseq.txt files, and aligned to a

491    human reference (hg19) using Burrows–Wheeler Aligner (BWA).[50] Aligned reads

492    duplicating the start position of another read were flagged as duplicates and not analysed.

493    Data was processed using the Genome Analysis ToolKit (GATK v3).[51-53] Reads were

494    locally realigned around insertions-deletions (indels) and their base qualities were

495    recalibrated. Variant calling was performed on both exomes and flanking 50 base pairs of

496    intronic sequence across all samples using the HaplotypeCaller (HC) tool from the

497    GATK to generate a gVCF. Joint genotyping was subsequently performed and 'raw'

498    variant data for each sample was formatted (variant call format (VCF)). Single nucleotide

499    polymorphisms (SNVs) and indel sites were initially filtered after variant calibration

500    marked sites of low quality that were likely false positives.

501    **Data analysis QC**. Fingerprint concordance between sequence data and fingerprint

502    genotypes was evaluated. Variant calls were evaluated on both bulk and per- sample

503    properties: novel and known variant counts, transition–transversion (TS–TV) ratio,

504    heterozygous–homozygous non-reference ratio, and deletion/insertion ratio. Both bulk

505    and sample metrics were compared to historical values for exome sequencing projects at

506    the Broad Institute. No significant deviation of from historical values was noted.

507

**Data processing and quality control of exome sequencing.**

**Variant annotation.** Variants were annotated using Variant Effect Predictor[54] and the LOFTEE[10] plugin to identify protein-truncating variants predicted to disrupt the respective gene's function with "high confidence." Each allele at polyallelic sites was separately annotated.

**Sample level quality control.** We performed quality control of samples using the following steps. For quality control of samples, we used bi-allelic SNVs that passed the GATK VQSR filter and were on genomic regions targeted by both ICE and Agilent exome captures. We removed samples with discordance rate > 10% between genotypes from exome sequencing with genotypes from array-based genotyping and samples with sex mismatch between inbreeding coefficient on chromosome X and fingerprinting. We tested for sample contamination using the verifyBamID software, which examines the proportion of non-reference bases at reference sites, and excluded samples with high estimated contamination (FREEMIX scores > 0.2).[55] After removing monozygotic twins or duplicate samples using the KING software[56], we removed outlier samples with too many or too few SNVs (>17,000 or <12,000 total variants; >400 singletons; and >300 doubletons). We removed those with extreme overall transition-to-transversion ratios (>3.8 or <3.3) and heterozygosity (heterozygote:non-reference homozygote ratio >6 or <2). Finally, we removed samples with high missingness (>0.05).

**Variant level quality control.** Variant score quality recalibration was performed separately for SNVs and indels use the GATK VariantRecalibrator and ApplyRecalibration to filter out variants with lower accuracy scores. Additionally, we

24

530    removed sites with an excess of heterozygosity calls (InbreedingCoeff <-0.3). To further

531    reduce the rate of inaccurate variant calls, we further filtered out SNVs with low average

532    quality (quality per depth of coverage (QD) < 2) and a high degree of missingness (> 20

533    %), and indels also with low average quality (quality per depth of coverage (QD) < 3)

534    and a high degree of missingness (> 20 %).

535

536    **Laboratory methods for proteomics.**

537    **Protein capture.** For 91 participants, enriched for homozygous pLoF mutations, we

538    measured 1,310 protein analytes in plasma using the SOMAscan assay (SomaLogic,

539    Boulder, CO, USA). Protein-capture was performed using modified aptamer technology

540    as previously described.[57] Briefly, modified nucleotides, analogous to antibodies, on a

541    custom DNA microarray recognize intact tertiary protein structures. After washing,

542    complexes are released from beads by photocleavage of the linker with UV light and the

543    resultant relative fluorescent unit is proportional to target protein.

544    **Quality control.** Samples (n = 7) were excluded if they showed evidence of systematic

545    inflation of association, or >5 % of traits in the top or bottom 1[st] percentile of the analytic

546    distribution.

547

548    **Methods for manual curation of a subset of** pLoF **variants.**

549    Manual curation was performed collaboratively by three geneticists: 25 pLoF variant

550    calls were reviewed independently by two reviewers and compared to ensure similar

551    review criteria before the remainder was divided and separately assessed by each of the

552    two reviewers separately. A third reviewer resolved discrepancies. Read and genotype

25

553   support was confirmed by review of reads in Integrative Genomics Viewer. We flagged

554   pLoF variants for any of the following six reasons:  1) read-mapping flags; 2) genotyping

555   flags; 3) presence of an additional polymorphism which rescues protein truncation; 4)

556   presence of an additional polymorphism which rescues splice site; 5) if affecting a

557   minority of transcripts; and 6) polymorphism occurs at exon or splice site with low

558   conservation.  Criteria for these reasons are provided in **Supplementary Table 8**.

559

560   **Methods for inbreeding analyses.**

561   **Array-derived runs of homozygosity.** Analyses were conducted in PLINK[46] using

562   genome-wide association (GWAS) data in PROMIS and HapMap 3 populations.

563   Segments of the genome that were at-least 1.5 Mb long, had a SNP density of 1 SNP per

564   20 kb and had 25 consecutive homozygous SNPs (1 heterozygous and/or 5 missing SNPs

565   were permitted within a segment) were defined to be in a homozygous state (or referred

566   as "runs of homozygosity" (ROH)), as described previously.[58] Homozygosity was

567   expressed as the percentage of the autosomal genome found in a homozygous state, and

568   was calculated by dividing the sum of ROH length within each individual by the total

569   length of the autosome in PROMIS and HapMap 3 populations respectively. To

570   investigate variability in homozygosity explained by parental consanguinity, the

571   difference in $R^2$ is reported for a linear regression model of homozygosity including and

572   excluding parental consanguinity on top of age, sex and the first 10 principal components

573   derived from the typed autosomal GWAS data.

574   **Sequencing-derived coefficient of inbreeding.** We compared the coefficient of

575   inbreeding distributions of 10,503 exome sequenced PROMIS participants with 15,248

26

576    participants (European ancestry = 12,849, and African ancestry = 2,399) who were

577    exome sequenced at the Broad Institute (Cambridge, MA) from the Myocardial Infarction

578    Genetics consortium.[48] We extracted approximately 5,000 high-quality polymorphic

579    SNVs in linkage equilibrium present on both target intervals that passed variant quality

580    control metrics based on HapMap 3 data.[59] Using PLINK, we estimated the coefficient of

581    inbreeding separately within each ethnicity group.[46] The coefficient of inbreeding was

582    estimated as the observed degree of homozygosity compared with the anticipated

583    homozygosity derived from an estimated common ancestor.[60] The Wilcoxon-Mann-

584    Whitney test was used to test whether PROMIS participants had different median

585    coefficients of inbreeding compared to other similarly sequenced outbred individuals and

586    whether the median coefficient of inbreeding was different between PROMIS participants

587    who reported parental relatedness versus not. A two-sided $P$ of 0.05 was the pre-specified

588    threshold for statistical significance.

589

590    **Methods for sequencing projection analysis.**

591    To compare the burden of unique completely inactivated genes in the PROMIS cohort

592    with outbred cohorts of diverse ethnicities, we extracted the minor allele frequencies

593    (maf) of "high confidence" loss-of-function mutations observed in the first 7,078

594    sequenced PROMIS participants, and in European, African, and East Asian ancestry

595    participants from the Exome Aggregation Consortium (ExAC r0.3;

596    exac.broadinstitute.org). For each gene and for each ethnicity, the combined minor allele

597    frequency (cmaf) of rare (maf < 0.1%) "high confidence" loss-of-function mutations was

598    calculated. We then simulated the number of unique completely inactivated genes across

599    a range of sample sizes per ethnicity and PROMIS. The expected probability of observing

600    complete inactivation (two pLoF copies in an individual) of a gene was calculated as

601    $(1 - F) * cmaf^2 + F * cmaf$, which accounts for allozygous and autozygous,

602    respectively, mechanisms for complete genie knockout. F, the inbreeding coefficient, is

603    defined

604    as $F = 1 - (expected\ heterozygosity\ rate\ /\ observed\ heterozygosity\ rate)$. For

605    PROMIS, the median F inbreeding coefficient (0.016) was used for estimation. Down-

606    sampling within the observed sample size for both high-confidence pLoF mutations and

607    synonymous variants did not deviate significantly from the expected trajectory

608    (**Supplementary Figure 11**).  For a range of sample sizes (0-200,000), each gene was

609    randomly sampled under a binomial distribution ($X \sim B(n, cmaf)$) and it was

610    determined if the gene was successfully sampled at least once. To refine the estimated

611    count of unique genes per sample size, each sampling was replicated ten times.

612

613    **Methods for constraint score analysis.**

614    We sought to determine whether the observed homozygous pLoF genes were under less

615    evolutionary constraint by first obtaining constraint loss of function constraint scores

616    derived from the Exome Aggregation Consortium (Lek M et al, in preparation).[11,12]

617    Briefly, we used the number of observed and expected rare (MAF < 0.1%) loss of

618    function variants per gene to determine to which of three classes it was likely to belong:

619    pLoF (observed variation matches expectation), recessive (observed variation is ~50%

620    expectation), or haploinsufficient (observed variation is <10% of expectation). The

621    probability of being loss of function intolerant (pLI) of each transcript was defined as the

622    probability of that transcript falling into the haploinsufficient category. Transcripts with a

623    pLI $\geq$ 0.9 are considered very likely to be loss of function intolerant; those with pLI $\leq$ 0.1

624    are not likely to be loss of function intolerant. A list of 1,317 genes were randomly

625    sampled from a list of sequenced genes 1,000 times and the proportion of loss of function

626    intolerant genes compared to the proportion of the observed homozygous pLoF genes

627    was compared using the chi square test. The likelihood that the distribution of the test

628    statistics deviated from the pLoF was ascertained.

629

630    Additionally, we sought to determine whether there were genes with appreciate pLoF

631    allele frequencies yet relative depletion of homozygous pLoF genotypes. We computed

632    estimated genotype frequencies based on Hardy-Weinberg equilibrium and the F

633    inbreeding coefficient and compared the frequencies to the observed genotype counts

634    with the chi square goodness-of-fit test. A nominal $P < 0.05$ is used to demonstrate at

635    least nominal association.

636

637    **Methods for rare variant association analysis.**

638    **Recessive model association discovery**. We sought to determine whether complete loss-

639    of-function of a gene was associated with a dense array of phenotypes. We extracted a list

640    of individuals per gene who were homozygous for a high confidence pLoF allele that was

641    rare (minor allele frequency < 1 %) in the cohort. From a list of 1,317 genes where there

642    was at least one participant homozygous pLoF and a list of 201 traits, we initially

643    considered 264,717 gene-trait pairings. To reduce the likelihood of false positives, we

644    only considered gene-trait pairs where there were at least two homozygous pLoF alleles

645    per gene phenotyped for a given trait yielding 18,959 gene-trait pairs for analysis.

646    For all analyses, we constructed generalized linear models to test whether complete loss

647    of function versus non-carriers was associated with trait variation. A logit link was used

648    for binomial outcomes. Right-skewed continuous traits were natural log transformed.

649    Age, sex, and myocardial infarction status were used as covariates in all analyses. We

650    extracted principal components of ancestry using EIGENSTRAT to control for

651    population stratification in all analyses.[61] For lipoprotein-related traits, the use of lipid-

652    lowering therapy was used as a covariate. For glycemic biomarkers, only non-diabetics

653    were used in the analysis. The $P$ threshold for statistical significance was 0.05 / 18,959 =

654    $3 \times 10^{-6}$.

655    **Heterozygote association replication**. We hypothesized that some of the associations

656    for homozygous pLoF alleles will display a more modest effect for heterozygous pLoF

657    alleles. Thus, the aforementioned analyses were performed comparing heterozygous

658    pLoF carriers to non-carriers for the 26 homozygous pLoF-trait associations that

659    surpassed prespecified statistical significance. A $P$ of 0.05 / 26 = 0.002 was set for

660    statistical significance for these restricted analyses.

661    **Association for single genic homozygotes**. We performed an exploratory analysis of

662    gene-trait pairs where there was only one phenotyped homozygous pLoF. We performed

663    the above association analyses for genes where there was only one homozygous pLoF

664    phenotyped for a given trait and we focused on those with the most extreme standard Z

665    score statistics (|Z score| > 5) from the primary association analysis and required that

666    there to also be nominal evidence for association ($P < 0.05$) in heterozygotes as well to

667    maximize confidence in an observed single homozygous pLoF-trait association.

668    **Recessive model association discovery for proteomics**. Among the 84 participants with

669    proteomic analyses of 1,310 protein analytes, 9 genes were observed in the homozygous

670    pLoF state at least twice. We log transformed each analyte and associated with

671    homozygous pLoF genotype status, adjusting for proteomic plate, age, sex, myocardial

672    infarction status, and principal components. Gene-analyte associations were considered

673    significant if P values were less than $0.05 / (1{,}310 \times 9) = 4.3 \times 10^{-6}$.

674

675    **Methods for recruitment and phenotyping of an *APOC3* p.Arg19Ter proband and**

676    **relatives.**

677    **Methods for Sanger sequencing**. We collected blood samples from a total of 28

678    subjects, including one of the four *APOC3* p.Arg19Ter homozygous participants along

679    with 27 of his family and community members for DNA extraction and separated into

680    plasma for lipid and apolipoprotein measurements. All subjects were consented prior to

681    initiation of the studies (IRB: 00007048 at the Center for Non-Communicable Diseases,

682    Paksitan). DNA was isolated from whole blood using a reference phenol-chloroform

683    protocol.[62] Genotypes for the p.Arg19Ter variant were determined in all 28 participants

684    by Sanger sequencing. A 685 bp region of the *APOC3* gene including the base position

685    for this variant was amplified by PCR (Expand HF PCR Kit, Roche) using the following

686    primer sequences: Forward primer CTCCTTCTGGCAGACCCAGCTAAGG, Reverse

687    primer CCTAGGACTGCTCCGGGGAGAAAG. PCR products were purified with Exo-

688    SAP-IT (Affymetrix) and sequenced via Sanger sequencing using the same primers.

689 **Oral fat tolerance test.** Six non-carriers and seven homozygotes also participated in an

690 oral fat tolerance test. Participants fasted overnight and then blood was drawn for

691 measurement of baseline fasted lipids. Following this, participants were administered an

692 oral load of heavy cream (50 g fat per square meter of body surface area as calculated by

693 the method of Mosteller[63]). Participants consumed this oral load within a time span of 20

694 minutes and afterwards consumed 200 mL of water. Blood was drawn at 2, 4, and 6 hours

695 after oral fat consumption as done previously.[38,64] All lipid and apolipoprotein

696 measurements from these plasma samples were determined by immunoturbidimetric

697 assays on an ACE Axcel Chemistry analyzer (Alfa Wasserman). A comparisons of area-

698 under-the curve triglycerides was performed between *APOC3* p.Arg19Ter homozygotes

699 and non-carriers using a two independent sample Student's t test; $P < 0.05$ was

700 considered statistically significant.

701 **Tables**

702 **Table 1. Baseline characteristics of exome sequenced study participants.**

| Characteristic | Value (n = 10,503) |
|---|---|
| Age (yrs) – mean (sd) | 52.0 (9.0) |
| Women – no. (%) | 1,802 (17.2 %) |
| Parents closely related – no. (%) | 4,101 (39.0 %) |
| Spouse closely related – no. (%) | 4,182 (39.8 %) |
| Ethnicity – no. (%) | |
| Urdu | 3,846 (36.6 %) |
| Punjabi | 3,668 (34.9 %) |
| Sindhi | 1,128 (10.7 %) |
| Pathan | 589 (5.6 %) |
| Memon | 141 (1.3 %) |
| Gujarati | 109 (1.0 %) |
| Balochi | 123 (1.2 %) |
| Other | 891 (8.5 %) |
| Hypertension – no. (%)[*] | 4,744 (45.2 %) |
| Hypercholesterolemia – no. (%)[†] | 2,924 (27.8 %) |
| Diabetes mellitus – no. (%)[‡] | 4,264 (40.6 %) |
| Coronary heart disease – no. (%)[§] | 4,793 (45.6 %) |
| Smoking – no. (%)[‖] | 4,201 (40.0 %) |

| | |
|---|---|
| **BMI (m/kg$^2$) – mean (sd)** | 25.9 (4.2) |

703   *Hypertension defined as systolic blood pressure ≥140 mmHg, diastolic blood pressure

704   ≥90 mmHg, or antihypertensive treatment.

705   †Hypercholesterolemia defined as serum total cholesterol >240 mg/dL, lipid lowering

706   therapy or self-report.

707   ‡Diabetes defined as fasting blood glucose ≥126 mg/dL, or HbA1c >6.5 %, oral

708   hypoglycemics, insulin treatment, or self-report.

709   §Coronary heart disease defined as acute myocardial infarction as determined by clinical

710   symptoms with typical EKG findings or elevated serum troponin I.

711   ‖Smoking defined as active current or prior tobacco smoking.

712

**Figure Legends**

**Fig 1. a**, Most genes are observed in the homozygous pLoF state in only single

individuals. **b.** The distribution of F inbreeding coefficient of PROMIS participants is

compared to those of outbred samples of African (AFR) and European (EUR) ancestry. **c**,

The burden of homozygous pLoF genes per individual is correlated with coefficient of

inbreeding.

**Fig 2. a.-b.** Carriage of a splice-site mutation, c.663+1G>A, in *PLA2G7* leads to a dose-

dependent reduction of both lipoprotein-associated phospholipase A2 (Lp-PLA2) mass

and activity, with homozygotes having no circulating Lp-PLA2. **c.** Despite substantial

reductions of Lp-PLA2 activity, *PLA2G7* c.663+1G>A heterozygotes and homozygotes

have similar coronary heart disease risk when compared with non-carriers.

**Fig 3. a.-d.** *APOC3* pLoF genotype status, apolipoprotein C-III, triglycerides, HDL

cholesterol and LDL cholesterol distributions among all sequenced participants.

Apolipoprotein C-III concentration is displayed on a logarithmic base 10 scale. **e.** A

proband with *APOC3* pLoF homozygote genotype as well as several family members

were recalled for provocative phenotyping. Surprisingly, the spouse of the proband was

also a pLoF homozygote, leading to nine obligate homozygote children. Given the

extensive first-degree unions, the pedigree is simplified for clarity. **f.** *APOC3* p.Arg19Ter

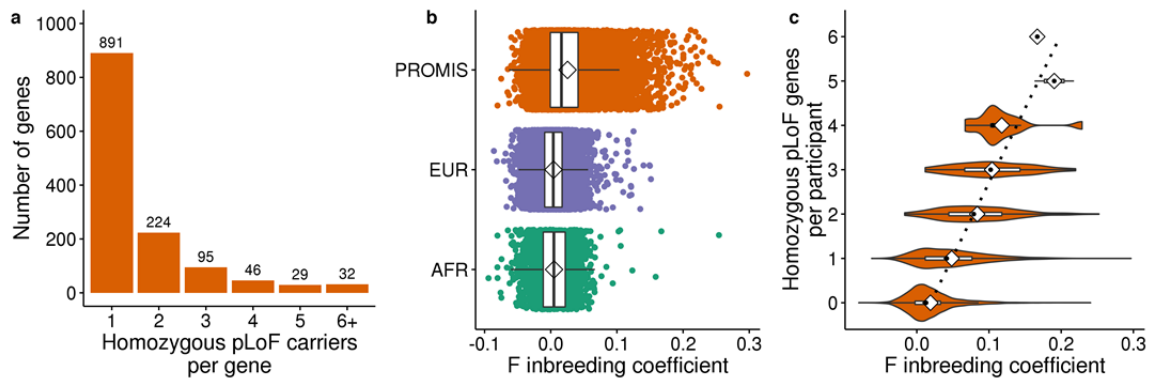homozygotes and non-carriers within the same family were challenged with a 50 g/m$^2$ fat

735     feeding. Homozygotes had lower baseline triglyceride concentrations and displayed

736     marked blunting of post-prandial rise in plasma triglycerides.

737

738     **Fig 4.** Number of unique homozygous pLoF genes anticipated with increasing sample

739     sizes sequenced in PROMIS compared with similar African (AFR) and European (EUR)

740     sample sizes. Estimates derived using observed allele frequencies and degree of
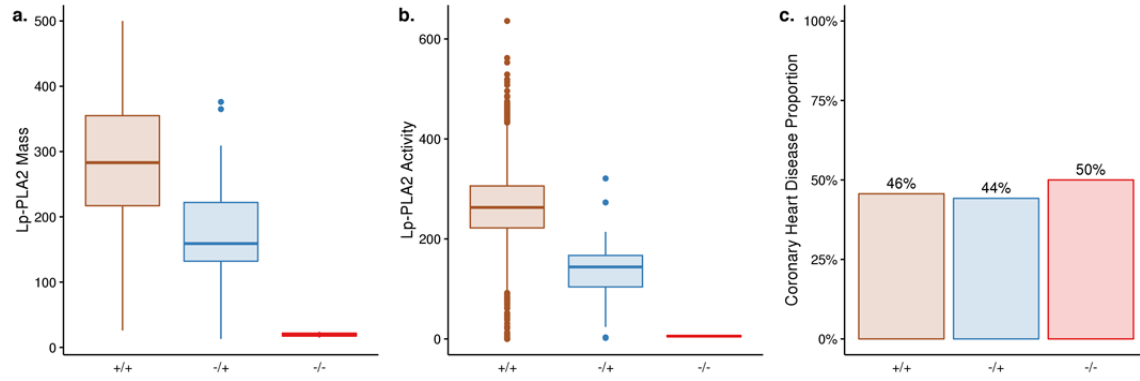
741     inbreeding.

**<u>Figures</u>**

744

745 **Fig. 1. Homozygous pLoF burden in PROMIS is driven by excess autozygosity.**

746

747

**Fig. 2. Carriers of *PLA2G7* splice mutation have diminished Lp-PLA2 mass ($P = 6$ x $10^{-5}$) and activity ($P = 2$ x $10^{-7}$) but similar risk for coronary heart disease risk when compared to non-carriers ($P = 0.87$).**
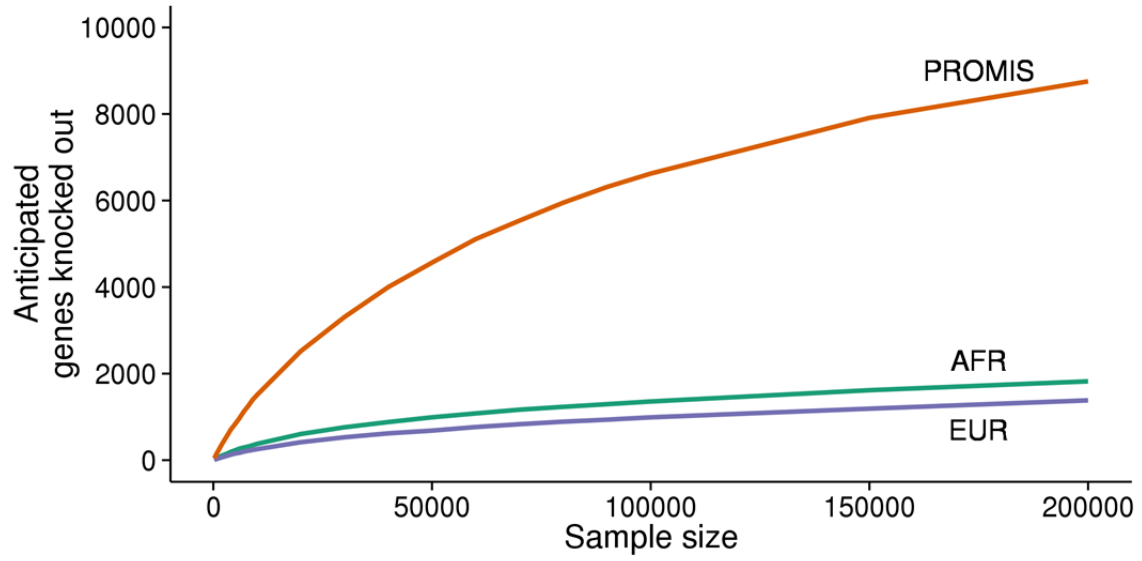
748

749

750

751

752

**Fig 3.** *APOC3* **pLoF homozygotes have diminished fasting triglycerides and blunted**

**post-prandial lipemia.**

755

**Fig 4. Simulations anticipate many more homozygous pLoF genes in the PROMIS**

757 **cohort.**

758 **References**

759 1      Eisenberg, D., Marcotte, E. M., Xenarios, I. & Yeates, T. O. Protein Function in
760        the Post-Genomic Era. *Nature* **405**, 823-826, doi:10.1038/35015694 (2000).
761 2      MacArthur, D. G. *et al.* A Systematic Survey of Loss-of-Function Variants in
762        Human Protein-Coding Genes. *Science* **335**, 823-828,
763        doi:10.1126/science.1215040 (2012).
764 3      Sulem, P. *et al.* Identification of a Large Set of Rare Complete Human
765        Knockouts. *Nat Genet*, doi:10.1038/ng.3243 (2015).
766 4      Bittles, A. H., Mason, W. M., Greene, J. & Rao, N. A. Reproductive Behavior and
767        Health in Consanguineous Marriages. *Science* **252**, 789-794 (1991).
768 5      Crosby, J. *et al.* Loss-of-Function Mutations in Apoc3, Triglycerides, and
769        Coronary Disease. *N Engl J Med* **371**, 22-31, doi:10.1056/NEJMoa1307095
770        (2014).
771 6      Jorgensen, A. B., Frikke-Schmidt, R., Nordestgaard, B. G. & Tybjaerg-Hansen,
772        A. Loss-of-Function Mutations in Apoc3 and Risk of Ischemic Vascular
773        Disease. *N Engl J Med* **371**, 32-41, doi:10.1056/NEJMoa1308027 (2014).
774 7      Saleheen, D. *et al.* The Pakistan Risk of Myocardial Infarction Study: A
775        Resource for the Study of Genetic, Lifestyle and Other Determinants of
776        Myocardial Infarction in South Asia. *Eur J Epidemiol* **24**, 329-338,
777        doi:10.1007/s10654-009-9334-y (2009).
778 8      Modell, B. & Darr, A. Science and Society: Genetic Counselling and Customary
779        Consanguineous Marriage. *Nat Rev Genet* **3**, 225-229, doi:10.1038/nrg754
780        (2002).
781 9      Lander, E. S. & Botstein, D. Homozygosity Mapping: A Way to Map Human
782        Recessive Traits with the DNA of Inbred Children. *Science* **236**, 1567-1570
783        (1987).
784 10     Karczewski, K. J. *Loftee (Loss-of-Function Transcript Effect Estimator)*,
785        <https://github.com/konradjk/loftee> (2015).
786 11     De Rubeis, S. *et al.* Synaptic, Transcriptional and Chromatin Genes Disrupted
787        in Autism. *Nature* **515**, 209-215, doi:10.1038/nature13772 (2014).
788 12     Samocha, K. E. *et al.* A Framework for the Interpretation of De Novo Mutation
789        in Human Disease. *Nat Genet* **46**, 944-950, doi:10.1038/ng.3050 (2014).
790 13     Wang, T. *et al.* Identification and Characterization of Essential Genes in the
791        Human Genome. *Science* **350**, 1096-1101, doi:10.1126/science.aac7041
792        (2015).
793 14     Eppig, J. T. *et al.* The Mouse Genome Database (Mgd): Facilitating Mouse as a
794        Model for Human Biology and Disease. *Nucleic Acids Res* **43**, D726-736,
795        doi:10.1093/nar/gku967 (2015).
796 15     Georgi, B., Voight, B. F. & Bucan, M. From Mouse to Human: Evolutionary
797        Genomics Analysis of Human Orthologs of Essential Genes. *PLoS Genet* **9**,
798        e1003484, doi:10.1371/journal.pgen.1003484 (2013).
799 16     Fuchs, M. *et al.* The P400 Complex Is an Essential E1a Transformation Target.
800        *Cell* **106**, 297-307 (2001).

801    17    Fazzio, T. G., Huff, J. T. & Panning, B. An Rnai Screen of Chromatin Proteins
802          Identifies Tip60-P400 as a Regulator of Embryonic Stem Cell Identity. *Cell*
803          **134**, 162-174, doi:10.1016/j.cell.2008.05.031 (2008).
804    18    Narasimhan, V. M. *et al.* Health and Population Effects of Rare Gene
805          Knockouts in Adult Humans with Related Parents. *Science*,
806          doi:10.1126/science.aac8624 (2016).
807    19    Lek, M. *et al.* Analysis of Protein-Coding Genetic Variation in 60,706 Humans.
808          *Nature* **536**, 285-291, doi:10.1038/nature19057 (2016).
809    20    Kelly, M. P. *et al.* Phosphodiesterase 11a in Brain Is Enriched in Ventral
810          Hippocampus and Deletion Causes Psychiatric Disease-Related Phenotypes.
811          *Proc Natl Acad Sci U S A* **107**, 8457-8462, doi:10.1073/pnas.1000730107
812          (2010).
813    21    Di Angelantonio, E. *et al.* Lipid-Related Markers and Cardiovascular Disease
814          Prediction. *Jama* **307**, 2499-2506, doi:10.1001/jama.2012.6571 (2012).
815    22    White, H. D. *et al.* Darapladib for Preventing Ischemic Events in Stable
816          Coronary Heart Disease. *N Engl J Med* **370**, 1702-1711,
817          doi:10.1056/NEJMoa1315878 (2014).
818    23    O'Donoghue, M. L. *et al.* Effect of Darapladib on Major Coronary Events after
819          an Acute Coronary Syndrome: The Solid-Timi 52 Randomized Clinical Trial.
820          *Jama* **312**, 1006-1015, doi:10.1001/jama.2014.11061 (2014).
821    24    Polfus, L. M., Gibbs, R. A. & Boerwinkle, E. Coronary Heart Disease and
822          Genetic Variants with Low Phospholipase A2 Activity. *N Engl J Med* **372**, 295-
823          296, doi:10.1056/NEJMc1409673 (2015).
824    25    Carr, B. A., Wan, J., Hines, R. N. & Yost, G. S. Characterization of the Human
825          Lung Cyp2f1 Gene and Identification of a Novel Lung-Specific Binding Motif. *J
826          Biol Chem* **278**, 15473-15483, doi:10.1074/jbc.M300319200 (2003).
827    26    Standiford, T. J. *et al.* Interleukin-8 Gene Expression by a Pulmonary
828          Epithelial Cell Line. A Model for Cytokine Networks in the Lung. *J Clin Invest*
829          **86**, 1945-1953, doi:10.1172/JCI114928 (1990).
830    27    Barnes, P. J. Chronic Obstructive Pulmonary Disease. *N Engl J Med* **343**, 269-
831          280, doi:10.1056/NEJM200007273430407 (2000).
832    28    Murray, I. A., Coupland, K., Smith, J. A., Ansell, I. D. & Long, R. G. Intestinal
833          Trehalase Activity in a Uk Population: Establishing a Normal Range and the
834          Effect of Disease. *Br J Nutr* **83**, 241-245 (2000).
835    29    Christiansen, D. *et al.* Humans Lack Igb3 Due to the Absence of Functional
836          Igb3-Synthase: Implications for Nkt Cell Development and Transplantation.
837          *PLoS Biol* **6**, e172, doi:10.1371/journal.pbio.0060172 (2008).
838    30    Dahl, K., Buschard, K., Gram, D. X., d'Apice, A. J. & Hansen, A. K. Glucose
839          Intolerance in a Xenotransplantation Model: Studies in Alpha-Gal Knockout
840          Mice. *APMIS* **114**, 805-811, doi:10.1111/j.1600-0463.2006.apm_393.x
841          (2006).
842    31    Casu, A. *et al.* Insulin Secretion and Glucose Metabolism in Alpha 1,3-
843          Galactosyltransferase Knock-out Pigs Compared to Wild-Type Pigs.
844          *Xenotransplantation* **17**, 131-139, doi:10.1111/j.1399-3089.2010.00572.x
845          (2010).

846   32   Schneider, M. R. & Wolf, E. The Epidermal Growth Factor Receptor Ligands at
847        a Glance. *J Cell Physiol* **218**, 460-466, doi:10.1002/jcp.21635 (2009).
848   33   Wang, G. X. *et al.* The Brown Fat-Enriched Secreted Factor Nrg4 Preserves
849        Metabolic Homeostasis through Attenuation of Hepatic Lipogenesis. *Nat Med*
850        **20**, 1436-1443, doi:10.1038/nm.3713 (2014).
851   34   Murtazina, R. *et al.* Tissue-Specific Regulation of Sodium/Proton Exchanger
852        Isoform 3 Activity in Na(+)/H(+) Exchanger Regulatory Factor 1 (Nherf1)
853        Null Mice. Camp Inhibition Is Differentially Dependent on Nherf1 and
854        Exchange Protein Directly Activated by Camp in Ileum Versus Proximal
855        Tubule. *J Biol Chem* **282**, 25141-25151, doi:10.1074/jbc.M701910200
856        (2007).
857   35   Wang, B., Yang, Y. & Friedman, P. A. Na/H Exchange Regulatory Factor 1, a
858        Novel Akt-Associating Protein, Regulates Extracellular Signal-Regulated
859        Kinase Signaling through a B-Raf-Mediated Pathway. *Mol Biol Cell* **19**, 1637-
860        1645, doi:10.1091/mbc.E07-11-1114 (2008).
861   36   Karim, Z. *et al.* Nherf1 Mutations and Responsiveness of Renal Parathyroid
862        Hormone. *N Engl J Med* **359**, 1128-1135, doi:10.1056/NEJMoa0802836
863        (2008).
864   37   Huff, M. W. & Hegele, R. A. Apolipoprotein C-Iii: Going Back to the Future for a
865        Lipid Drug Target. *Circ Res* **112**, 1405-1408,
866        doi:10.1161/CIRCRESAHA.113.301464 (2013).
867   38   Pollin, T. I. *et al.* A Null Mutation in Human Apoc3 Confers a Favorable Plasma
868        Lipid Profile and Apparent Cardioprotection. *Science* **322**, 1702-1705,
869        doi:10.1126/science.1161524 (2008).
870   39   Gaudet, D. *et al.* Antisense Inhibition of Apolipoprotein C-Iii in Patients with
871        Hypertriglyceridemia. *N Engl J Med* **373**, 438-447,
872        doi:10.1056/NEJMoa1400283 (2015).
873   40   Gaudet, D. *et al.* Targeting Apoc3 in the Familial Chylomicronemia Syndrome.
874        *N Engl J Med* **371**, 2200-2206, doi:10.1056/NEJMoa1400284 (2014).
875   41   Graham, M. J. *et al.* Antisense Oligonucleotide Inhibition of Apolipoprotein C-
876        Iii Reduces Plasma Triglycerides in Rodents, Nonhuman Primates, and
877        Humans. *Circ Res* **112**, 1479-1490, doi:10.1161/circresaha.111.300367
878        (2013).
879   42   Brown, S. D. & Moore, M. W. Towards an Encyclopaedia of Mammalian Gene
880        Function: The International Mouse Phenotyping Consortium. *Dis Model Mech*
881        **5**, 289-292, doi:10.1242/dmm.009878 (2012).
882   43   Liu, Y. *et al.* Quantitative Variability of 342 Plasma Proteins in a Human Twin
883        Population. *Mol Syst Biol* **11**, 786, doi:10.15252/msb.20145728 (2015).
884   44   Scott, E. M. *et al.* Characterization of Greater Middle Eastern Genetic Variation
885        for Enhanced Disease Gene Discovery. *Nat Genet*, doi:10.1038/ng.3592
886        (2016).
887   45   Kooner, J. S. *et al.* Genome-Wide Association Study in Individuals of South
888        Asian Ancestry Identifies Six New Type 2 Diabetes Susceptibility Loci. *Nat*
889        *Genet* **43**, 984-989, doi:10.1038/ng.921 (2011).

890    46    Purcell, S. *et al.* Plink: A Tool Set for Whole-Genome Association and
891          Population-Based Linkage Analyses. *Am J Hum Genet* **81**, 559-575,
892          doi:10.1086/519795 (2007).
893    47    Tennessen, J. A. *et al.* Evolution and Functional Impact of Rare Coding
894          Variation from Deep Sequencing of Human Exomes. *Science* **337**, 64-69,
895          doi:10.1126/science.1219240 (2012).
896    48    Do, R. *et al.* Exome Sequencing Identifies Rare Ldlr and Apoa5 Alleles
897          Conferring Risk for Myocardial Infarction. *Nature* **518**, 102-106,
898          doi:10.1038/nature13917 (2015).
899    49    Fisher, S. *et al.* A Scalable, Fully Automated Process for Construction of
900          Sequence-Ready Human Exome Targeted Capture Libraries. *Genome Biol* **12**,
901          R1, doi:10.1186/gb-2011-12-1-r1 (2011).
902    50    Li, H. & Durbin, R. Fast and Accurate Short Read Alignment with Burrows-
903          Wheeler Transform. *Bioinformatics* **25**, 1754-1760,
904          doi:10.1093/bioinformatics/btp324 (2009).
905    51    McKenna, A. *et al.* The Genome Analysis Toolkit: A Mapreduce Framework for
906          Analyzing Next-Generation DNA Sequencing Data. *Genome Res* **20**, 1297-
907          1303, doi:10.1101/gr.107524.110 (2010).
908    52    DePristo, M. A. *et al.* A Framework for Variation Discovery and Genotyping
909          Using Next-Generation DNA Sequencing Data. *Nat Genet* **43**, 491-498,
910          doi:10.1038/ng.806 (2011).
911    53    Van der Auwera, G. A. *et al.* From Fastq Data to High Confidence Variant Calls:
912          The Genome Analysis Toolkit Best Practices Pipeline. *Curr Protoc*
913          *Bioinformatics* **11**, 11 10 11-11 10 33, doi:10.1002/0471250953.bi1110s43
914          (2013).
915    54    McLaren, W. *et al.* Deriving the Consequences of Genomic Variants with the
916          Ensembl Api and Snp Effect Predictor. *Bioinformatics* **26**, 2069-2070,
917          doi:10.1093/bioinformatics/btq330 (2010).
918    55    Jun, G. *et al.* Detecting and Estimating Contamination of Human DNA Samples
919          in Sequencing and Array-Based Genotype Data. *Am J Hum Genet* **91**, 839-848,
920          doi:10.1016/j.ajhg.2012.09.004 (2012).
921    56    Manichaikul, A. *et al.* Robust Relationship Inference in Genome-Wide
922          Association Studies. *Bioinformatics* **26**, 2867-2873,
923          doi:10.1093/bioinformatics/btq559 (2010).
924    57    Gold, L. *et al.* Aptamer-Based Multiplexed Proteomic Technology for
925          Biomarker Discovery. *PLoS One* **5**, e15004,
926          doi:10.1371/journal.pone.0015004 (2010).
927    58    Hunter-Zinck, H. *et al.* Population Genetic Structure of the People of Qatar.
928          *Am J Hum Genet* **87**, 17-25, doi:10.1016/j.ajhg.2010.05.018 (2010).
929    59    Purcell, S. M. *et al.* A Polygenic Burden of Rare Disruptive Mutations in
930          Schizophrenia. *Nature* **506**, 185-190, doi:10.1038/nature12975 (2014).
931    60    Wright, S. Coefficients of Inbreeding and Relationship. *Am Nat* **56**, 330-338
932          (1922).
933    61    Price, A. L. *et al.* Principal Components Analysis Corrects for Stratification in
934          Genome-Wide Association Studies. *Nat Genet* **38**, 904-909,
935          doi:10.1038/ng1847 (2006).

936    62    Sambrook, J. & Russell, D. W. Purification of Nucleic Acids by Extraction with
937        Phenol:Chloroform. *CSH Protoc* **2006**, doi:10.1101/pdb.prot4455 (2006).
938    63    Mosteller, R. D. Simplified Calculation of Body-Surface Area. *N Engl J Med*
939        **317**, 1098, doi:10.1056/NEJM198710223171717 (1987).
940    64    Maraki, M. *et al.* Validity of Abbreviated Oral Fat Tolerance Tests for
941        Assessing Postprandial Lipemia. *Clin Nutr* **30**, 852-857,
942        doi:10.1016/j.clnu.2011.05.003 (2011).
943

944    **Supplementary Information** is linked to the online version of the paper at

945    www.nature.com/nature.

946

963 Muhammad Wajid, Irfan Ali, Muhammad Ikhlaq, Danish Sheikh, Muhammad Imran,

964 Matthew Walker, Nadeem Sarwar, Sarah Venorman, Robin Young, Adam Butterworth,

965 Hannah Lombardi, Binder Kaur and Nasir Sheikh. Fieldwork in the PROMIS study has

966 been supported through funds available to investigators at the Center for Non-

967 Communicable Diseases, Pakistan and the University of Cambridge, UK.

968

969

970 **Author Contributions** Sample recruitment and phenotyping was performed by D.S.,

971 P.F., J.D., A.R., M.Z., M.S., M.F., A.I., N.K.S., S.A., F.M., M.I., S.A., K.T., N.H.M.,

972 K.S.Z., N.Q., M.I., S.Z.R., F.M., K.M., N.A., and R.M.K.. D.S., P.F., J.D., and W.Z.

973 performed array-based genotyping and runs-of-homozygosity analyses. Exome

974 sequencing was coordinated by D.S., N.G., S.G., E.S.L., D.J.R., and S.K.. P.N., W.Z.,

975 H.H.W., and R.D. performed exome sequencing quality control and association analyses.

976 P.N., I.A., K.J.K., A.H.O., and D.G.M. performed variant annotation. D.S., S.K. and

977 D.J.R. performed confirmatory genotyping and lipoprotein biomarker assays. D.S. and

978 A.R. conducted recall based studies for the *APOC3* knockouts. P.N. and M.J.D.

979 performed bioinformatics simulations. P.N. and K.E.S. performed constraint score

980 analyses. D.S., P.N., and S.K. designed the study and wrote the paper. D.S. and P.N.

981 contributed equally. All authors discussed the results and commented on the manuscript.

982

986    (exac.broadinstitute.org). Reprints and permissions information is available at

987    www.nature.com/reprints. The authors do not declare competing financial interests.

988    Correspondence and requests for materials should be addressed to D.S. or S.K.

989    (saleheen@mail.med.upenn.edu or sekar@broadinstitute.org).