1    **Genome-wide association study of appendicular lean mass in UK Biobank cohort**

2    Yu-Fang Pei[#1,2], Yao-Zhong Liu[#3], Xiao-Lin Yang[#2,4], Hong Zhang[2,4], Gui-Juan Feng[1,2], Lei Zhang[2,4]

3    [1]Department of Epidemiology and Health Statistics,

4    [2]Jiangsu Key Laboratory of Preventive and Translational Medicine for Geriatric Diseases, School of Public Health,

5    Medical College of Soochow University, Jiangsu, PR China.

6    [3]Department of Biostatistics and Data Science, Tulane University School of Public Health and Tropical Medicine,

7    New Orleans, LA, USA.

8    [4]Center for Genetic Epidemiology and Genomics, School of Public Health, Medical College of Soochow University,

9    Jiangsu, PR China.

10    **Keywords**: appendicular lean mass, GWAS, UK Biobank

11    **Running title**: GWAS of appendicular lean mass

12    #: These authors contributed equally to this work.

13    Corresponding author:

14    Lei Zhang, Ph.D.

15    Associate Professor

16    Center for Genetic Epidemiology and Genomics

17    School of Public Health, Medical College, Soochow University

18    199 Ren-ai Rd., SuZhou City, Jiangsu Province 215123, PR China

19    Tel: (86)0512-65883871        Email: lzhang6@suda.edu.cn

20

21 **Abstract**

22 Lean body mass (LBM), an important physiological measure, has a strong genetic determination.

23 To clarify its genetic basis, a large-scale genome-wide association study (GWAS) of

24 appendicular lean mass (ALM) was conducted in 450,580 UK Biobank subjects. A total of 717

25 variants ($p < 5 \times 10^{-9}$) from 561 loci were identified, which were replicated across genders

26 (achieving $p < 5 \times 10^{-5}$ in both genders). The identified variants explained ~11% phenotypic

27 variance, accounting for one quarter of the total ~40% GWAS-attributable heritability. The

28 identified variants were enriched in gene sets related to musculoskeletal and connective tissue

29 development. Of interest are several genes, including *ADAMTS3*, *PAM*, *SMAD3* and *MEF2C*,

30 that either contain multiple significant variants or serve as the hub genes of the associated gene

31 sets. Polygenic score prediction based on the associated variants was able to distinguish subjects

32 of high from low ALM. Overall, our results offered significant findings on the genetic basis of

33 lean mass through an extraordinarily large sample GWAS. The findings are important to not only

34 lean mass *per se* but also other complex diseases, such as type 2 diabetes and fracture, as our

35 Mendelian randomization analysis showed that ALM is a protective factor for these two diseases.

36

37 **Introduction**

38   Lean body mass (LBM) is an important physiological index. The decline of LBM with aging,

39   also known as sarcopenia, is a critical factor for functional impairment and physical disability

40   and a major modifiable cause of frailty in the elderly [1, 2]. LBM is associated with bone mineral

41   density (BMD), and hence may be also   relevant to risk for osteoporosis [3]. Other LBM-related

42   conditions include dysmobility syndrome [4], sarcopenic obesity [5], and cachexia [6]. Overall,

43   sarcopenia was responsible for an increased risk of mortality, with a hazard ratio of 1.29 to 2.39

44   [7].

45   LBM has a significant genetic component, as evidenced by a high heritability of 50% to 80%

46   as observed in twin studies [8, 9]. However, findings on specific genes for human lean mass

47   variation remain limited even with the powerful genome-wide association study (GWAS)

48   approach. A key reason for the limited findings, as in other human complex traits, is the modest

49   sample size used in most GWAS so far performed in LBM [10-14], resulting in few SNPs (single

50   nucleotide polymorphisms) identified with genome-wide significance.

51   As a notable example, a recent large meta-analysis of GWAS amassed 20 cohorts of European

52   ancestry with a total sample size of >38,000 for whole body lean mass (WBM) and of >28,000

53   for appendicular lean mass (ALM) [15]. However, despite of the large sample used, the percent

54   variance explained by the identified SNPs was still only 0.23% and 0.16% for WBM and ALM,

55   respectively, suggesting that most of the heritability of LBM was still undetected. Therefore,

56   even with such a large GWAS meta-analysis, it is still necessary to boost the sample size further

57   so as to enhance the statistical power for detecting more causal SNPs underlying LBM.

58   Here in this study, with a sample containing ~half-million subjects of European origin from

59   UK Biobank (UKB), we performed a GWAS of appendicular lean mass (ALM). At the stringent

60    genome-wide significance level ($p < 5 \times 10^{-9}$), we identified >700 variants that were replicated

61    across genders. Our findings revealed a large number of genetic variants for LBM and

62    contributed to the characterization of the genetic architecture of this important complex trait.

63    Through this GWAS we demonstrated the power for mapping the genetic landscape of common

64    human complex traits/diseases using extraordinarily large samples.

65

## Results

66

67     Basic characteristics of the studied UKB sample are listed in **Supplementary Table 1**. In this

68     study, we quantified appendicular lean mass (ALM) by appendicular fat-free mass measured by

69     electronic impedance. This measurement of lean mass is reliable based on its strong correlation

70     with ALM measured by DXA in 4,294 UKB subjects (with a Pearson's correlation coefficient of

71     0.96, $p<2.2\times10^{-16}$).

72     *Main association results*

73     Raw ALM was adjusted with appendicular fat mass (AFM) and the adjusted ALM ($ALM_{adj}$)

74     was the phenotype used for the GWAS. Following quality control (QC) of both $ALM_{adj}$ and

75     genome-wide genotypes, data from 19.4 million variants with minor allele frequency

76     (MAF) >0.1% and imputation quality score >0.3 were available in 244,945 female and 205,635

77     male subjects.

78     In each gender group, additive effect of each variant was tested on $ALM_{adj}$ with BOLT-LMM

79     [16], controlling for age, $age^2$, height and $height^2$. The genomic inflation factor showed notable

80     inflation in both gender groups ($\lambda_{female}$=1.92, $\lambda_{male}$=1.77). To examine observed inflation for

81     potential polygenic effects and other biases, linkage disequilibrium score regression (LDSC)

82     analysis was performed [17]. The estimated mean chi-square and intercept were 2.34 and 1.12

83     for females, and 2.53 and 1.15 for males, corresponding to an attenuation ratio (AR) of 0.098 and

84     0.090, respectively. The AR estimates are comparable to those estimated in the subset of 369,968

85     unrelated British white subjects (0.090 and 0.074 for females and males, respectively), who were

86     extracted from the total sample.

87     Using BOLT-REML [18], GWAS-attributable heritability was estimated, which was 0.381

88     (s.e $3.30\times10^{-3}$) and 0.394 (s.e $3.80\times10^{-3}$) in females and males, respectively. LDSC estimated a

89   genetic correlation coefficient as high as 0.90 (s.e 0.01) between the two genders, implying that

90   most GWAS-attributable heritability was shared across genders.

91   Given the shared heritability across genders, across-gender meta-analysis was performed with

92   the inverse variance weighted fixed-effects model to combine the gender-specific GWAS results.

93   The meta-analysis signals have an AR of 0.115 (mean chi-square=3.69, intercept=1.31).

94   Genome-wide significance (GWS) level was set to $\alpha=5\times10^{-9}$, and a suggestive significance level

95   was set to $\alpha=5\times10^{-5}$. An association was declared to be "replicated" if it is 1) significant at the

96   GWS level in the across-gender meta-analysis and 2) significant at the suggestive level within

97   each gender.

98   Based on the above criteria, a total of 589 loci were identified at GWS level in across-gender

99   meta-analysis (p $<5\times10^{-9}$), which were also replicated (p $<5\times10^{-5}$) across genders. To check

100  potential linkage disequilibrium (LD) among these loci, LD analysis was performed on 589 lead

101  variants (each from one of the loci). It was found that 47 lead variants are not in linkage

102  equilibrium with each other (LD $r^2>0.1$) due to long-range LD. After removing 28 loci, the lead

103  variants in the remaining 561 loci were all in linkage equilibrium (LD $r^2<0.1$). Therefore, these

104  561 loci were treated as independent loci for downstream analysis.

105  Approximate conditional association analysis and across-gender meta-analysis were

106  recursively performed, which further identified an additional set of 156 conditionally significant

107  variants (p<$5\times10^{-9}$ in across gender meta-analysis) that were replicated across genders (p<$5\times10^{-5}$

108  within each gender). These additional variants were also in linkage equilibrium (LD $r^2<0.1$) with

109  the lead variants of the 561 loci.

110  In total, 717 (i.e., 561+156) independent variants from 561 distinct loci were associated with

111  ALM$_{adj}$ (**Supplementary Table 2**). Among the 717 lead variants, 172 achieved the strongest

112   significance level ($p<5\times10^{-9}$) in both genders (categorized here as the Tier 1 variants). Also, 144

113   variants achieved p values $<5\times10^{-9}$ in females, and p values $< 5\times10^{-5}$ in males; 62 variants

114   achieved p values $< 5\times10^{-9}$ in males, and p values $< 5\times10^{-5}$ in females (categorized here as the

115   Tier 2 variants). At last, 339 variants achieved p values $< 5\times10^{-5}$ in both genders and p values $<$

116   $5\times10^{-9}$ in across-gender meta-analysis (categorized here as the Tier 3 variants).

117   Of the above identified loci, 17 were reported by GWAS or meta-analysis of DXA-derived

118   lean mass [13, 15, 19]; 104 were reported by a study of electronic impedance measured lean

119   mass in a subset UKB cohort subjects (N=155,961) [20].

120   We also evaluated the overlap of the identified loci with those identified for several obesity

121   traits, including body mass index (BMI), waist circumference (WC), WC adjusted for BMI

122   ($WC_{adj}BMI$), waist-hip ratio (WHR) and WHR adjusted for BMI ($WHR_{adj}BMI$). SNPs in 302

123   loci (defined as the lead SNP + 500 kb flanking at each side) showed association with one or

124   more obesity traits at the conventional significance level $5.0\times10^{-8}$, while no trait was associated

125   with the remaining 259 loci, demonstrating their novelty and possibly, specificity to lean, but not

126   fat mass.

127   *Gender heterogeneity/specificity*

128   In addition to the 561 loci that are replicated across genders, our analysis also identified 152

129   loci that are significant ($p<5\times10^{-9}$) in across-gender meta-analysis but not significant at the

130   suggestive level ($p<5\times10^{-5}$) in either gender group (**Supplementary Table 3**). These loci may

131   represent gender specific signals pending further replication.

132   Of the 717 identified variants (of the 561 loci), 109 (15.2%) have a high level across-gender

133   meta-analysis heterogeneity ($I^2>50\%$), all (except one) of which belong to Tiers 1 or 2 variants.

134   A statistical test on gender difference in allele effect size showed that the difference is significant

135     in only 2 SNPs, rs2972156 ($p_{diff}$=2.49×10$^{-12}$) and rs1933801 ($p_{diff}$=4.65×10$^{-6}$), after accounting

136     for multiple testing (α=0.05/717=6.97×10$^{-5}$), suggesting that almost all of the identified variants

137     may have similar effect sizes across genders. The two SNPs (rs2972156 and rs1933801) with

138     different effect size between genders achieved p values of 1.30×10$^{-46}$ and 2.40×10$^{-26}$,

139     respectively, in males and p values of 4.20×10$^{-7}$ and 1.30×10$^{-6}$, respectively, in females.

140     *Heritability distribution*

141     The 717 identified variants include 654 common variants (MAF>5%), 52 less common

142     variants (5%≥MAF>1%) and 11 rare variants (MAF≤1%). Collectively, these variants explain

143     10.82% phenotypic variance in the total sample, most of which (9.91%) is accounted for by

144     common variants. As expected, variants with a smaller MAF generally have a larger per allele

145     effect size (**Figure 1**). For example, the average per allele effect size in rare variants (mean 0.11,

146     s.d 0.05) is 6-fold larger than common variants (mean 0.02, s.d 0.007).

147     Applying the stratified LDSC analysis, the explained heritability was partitioned into 24

148     functional categories [21]. Statistically significant enrichments were observed for 19 functional

149     categories (p<0.05/24, **Figure 2**). In line with the observations by Finucane et al. [21], regions

150     conserved in mammals showed the strongest enrichment of any category, with 2.6% of SNPs

151     explaining an estimated 34.5% of SNP heritability (enrichment ratio (EA)=13.2, P=3.39×10$^{-19}$).

152     Other categories with significant enrichment included coding regions (EA=8.8, P=1.76×10$^{-7}$), 3'

153     UTR (EA=5.7, P=3.73×10$^{-4}$), transcription starting site (EA=5.1, P=1.71×10$^{-5}$), and H3K9ac

154     histone marks (EA=5.1, P=2.07×10$^{-15}$). Neither promoter nor 5'-UTR region showed significant

155     enrichment, though 5'-UTR region had a high estimate of EA (15.5, p=0.03).

156     A new function of the stratified LDSC method was used to assess focal tissues for heritability

157     enrichment [22], using two gene expression datasets [23, 24]. A total of 19 tissues/cells are

158    enriched at a false discovery rate (FDR) <5% (**Figure 2**). About half (9) of them belong to

159    musculoskeletal and connective system, including cartilage, chondrocytes, osteoblasts,

160    fibroblasts, smooth muscle, myometrium, cervical vertebrae, synovial membrane and stromal

161    cells.

162    *Candidate genes prioritization*

163    To prioritize candidate genes at the associated loci, we used multiple analytical strategies. A

164    set of credible risk variants (CRVs) at each locus were defined as variants with high LD with the

165    lead variant ($r^2$>0.8). A total of 17,968 CRVs were defined (**Supplementary Table 4**). Based on

166    the CRVs, 6 types of supporting evidence were used to prioritize 1,337 candidate genes.

167    (**Supplementary Tables 5-10**).

168    A number of genes have multiple lines of supporting evidence. Peptidylglycine

169    Alpha-Amidating Monooxygenase (*PAM*) at 5q21.1, in particular, has all lines of supporting

170    evidence. This locus contains two independent signals. The first is a mis-sense rare SNP

171    rs78408340 (MAF=0.01%, p=6.10×10$^{-10}$) inside PAM, and the second is a common SNP

172    rs400596 located between *PAM* (129.5 kb from *PAM*) and *SLC06A1* genes (237.2 kb from

173    *SLC06A1*). Polymorphisms at rs400596 are associated with the *PAM* expression level in whole

174    blood (p=2.51×10$^{-21}$) [23] and associated with its protein level in peripheral blood (p=

175    p=2.51×10$^{-30}$) [25]. *PAM* is also prioritized by both SMR [26] and DEPICT [27], strengthening

176    its functional relevance.

177    *Comparison between imputation and sequencing-based association signal*

178    Of the 717 identified variants, 42 are mis-sense coding ones. Forty of them, including 7 rare

179    ones, are available in the recently released UKB exome sequencing data that contain a subset of

180    ~50,000 subjects from the whole UKB cohort. Using a set of 45,554 unrelated European subjects

181  who were both genotyped/imputed and sequenced, we compared the imputation-based

182  association results with exome sequencing based results. The 7 rare variants appeared to have

183  limited imputed dosage variation hence their imputation association p-values were not able to

184  derive. In the sequencing data, 3 of these 7 variants were nominally significant (p<0.05,

185  **Supplementary Table 11**), suggesting limited power in imputation-based association analysis

186  (compared with sequence-based analysis) for rare variants. This limited power may be alleviated

187  by increased sample size since in the whole UKB cohort these 7 rare variants achieved

188  significant p values in imputation-based association analysis.

189  Of the remaining 33 variants, the imputation-based and sequencing-based p-values were

190  highly concordant. For example, the imputation-based p-values are within 2-fold difference of

191  the sequencing-based p-values for up to 29 variants. Overall, these observations support that

192  imputation-based association signals are close to the real sequencing-based association signals in

193  a large sample. Therefore, imputation based GWAS may be able to identify true associations,

194  even for those of rare variants.

195  ***Mis-sense variants and the associated genes***

196  As mentioned above, of the 717 identified variants, 42 are mis-sense coding ones. Majority of

197  these 42 mis-sense mutations are predicted to be deleterious according to more than one

198  bioinformatics tool including PolyPhen2 [28], SIFT [29], PROVEAN [30] and Fathmm [31]

199  (**Supplementary Table 12**), supporting their functional relevance.

200  Mis-sense mutations are enriched among rare variants. Eight of the 11 rare variants are

201  mis-sense mutations, which is in clear contrast to 34 mis-sense mutations among the remaining

202  706 variants (odds ratio (OR)=55.37, Fisher's exact test p=$7.11 \times 10^{-9}$). Evidence of the

203  enrichment became stronger by comparing 21 mis-sense mutations from 63 rare or less common

204    variants vs. 22 mis-sense mutations from 654 common variants (OR=14.36, Fisher's exact test

205    p=$5.75\times10^{-13}$), suggesting that low frequency mutations are more likely to play a direct role in

206    changing protein function.

207    Genes containing mis-sense variants are listed in **Table 1**. In particular, the *ADAMTS3* gene

208    contains 3 rare or less common mis-sense variants (rs141374503 MAF=0.4% p=$2.02\times10^{-27}$;

209    rs150270324 MAF=1.3%, p=$2.36\times10^{-14}$, and rs139921635 MAF=2.4%, p=$4.06\times10^{-15}$). In

210    addition, it also contains multiple non-mis-sense variants, including 3 conditionally significant

211    variants in its intron region: rs72653979 (MAF=7.8%, p=$9.51\times10^{-11}$), rs78862524 (MAF=5.5%,

212    p =$3.24\times10^{-23}$) and rs769821342 (MAF=3.2%, p =$1.27\times10^{-14}$), and 2 in its flanking inter-genic

213    region: chr4:73496010 (MAF=47%, p=$3.87\times10^{-21}$) and rs10518106 (MAF=6%, p=$1.16\times10^{-83}$).

214    Though these SNPs are 367.0 kb apart at most, they are in linkage equilibrium with each other

215    (LD $r^2$<0.1). Together, the 8 variants from the *ADAMTS3* gene explain 0.18% of phenotypic

216    variance, making this region the most contributive locus.

217    ***Gene-based and gene set enrichment analyses***

218    A total of 3,101 genes were significant at the gene-based genome-wide significance level

219    (α=0.05/19,098=$2.62\times10^{-6}$, **Supplementary Table 13**), and 85 gene sets were significant at the

220    gene set significance level (α=0.05/10,655=$4.69\times10^{-6}$, **Supplementary Table 14**).

221    The most significant gene set is GO:0001501 'skeletal system development' (p=$8.88\times10^{-24}$),

222    followed by GO:0031012 'proteinaceous extracellular matrix' (p=$5.01\times10^{-14}$), GO:0061448

223    'connective tissue development' (p=$1.10\times10^{-13}$), GO:0048705 'skeletal system morphogenesis'

224    (p=$9.33\times10^{-13}$) and GO:0031012 'extracellular matrix' (p=$1.05\times10^{-12}$). Additional gene sets with

225    known function related to musculoskeletal and connective system, such as GO:0051216:

226     'cartilage development' ($p=3.30\times10^{-12}$) and GO:0042692 'muscle cell differentiation'

227     ($p=1.45\times10^{-6}$), were also identified.

228       Genes involved in multiple gene sets are likely to act as hub genes and may play a central

229     regulatory role. From the list of significant gene sets, the most frequently involved gene is

230     *SMAD3* (gene-based association $p=8.11\times10^{-42}$), which was involved in 46 out of the 85

231     significant gene sets. It was followed by *SOX9* (p=0.05, in 44 gene sets), *MEF2C* ($p=2.83\times10^{-9}$,

232     in 42 gene sets) and *BMP4* (p=0.15, in 42 genes). All these 4 genes were reported by previous

233     studies as important candidate genes for muscle development [32-35]. However, *SOX9* is only

234     nominally significant and *BMP4* is not significant at single gene level, indicating that the

235     significant pathway signals may not be contributed by the two genes. Altogether, there are 34

236     genes, each of which was involved in more than 30 of the 85 significant gene sets.

237       Protein-protein interaction (PPI) analysis using these 34 hub genes connects them into a tight

238     interactional network (**Figure 3**). This network contains multiple genes that are important for

239     skeletal muscle development, such as TGF signaling pathway genes (*TGFB1*, *TGFB2* and

240     *TGFBR2*), BMP signaling pathway genes (*BMP2* and *BMP4*) and SMAD family genes (*SMAD1*,

241     *SMAD2*, *SMAD3* and *SMAD4*).

242     ***Polygenic risk score profiling***

243       To assess the ability of the GWAS findings to predict ALM, a polygenic scoring analysis was

244     performed in the subset of 369,968 unrelated British white subjects from the UKB cohort. Three

245     quarters of the subjects (277,762 participants, including 149,329 females) were randomly

246     selected as the training sample, with the remaining subjects (92,206 participants, including

247     49,660 females) as the validation sample.

248    The training sample identifies 72,456 variants that achieved a p-value $<1\times10^{-5}$ for association

249    with $ALM_{adj}$. Using these variants as predictor, the predicted genome-wide polygenic score (GPS)

250    and the real phenotype residual in the validation sample are significantly correlated (Pearson's

251    correlation coefficient 0.22, 95% CI (0.21, 0.22), $p<2.2\times10^{-16}$). Mean phenotype residuals in the

252    top tail are significantly higher than that in the bottom tail of the GPS distribution (**Figure 4**).

253    For example, the predicted top 1% subjects have an increased average residual of 1.16 than the

254    predicted bottom 1% participants (0.57 (s.d 0.96) vs. -0.59 (s.d 0.94)), corresponding to an 1.69

255    kilo-gram (kg) increase of raw ALM (24.61 kg (s.d 5.89 kg) vs. 22.92 kg (s.d 5.27 kg)). In the

256    female group, the predicted top 1% participants have on average 1.39 kg increase of raw ALM

257    than the predicted bottom 1% participants (20.26 kg (s.d 2.75 kg) vs. 18.87 kg (s.d 2.45 kg)). In

258    males, the increase is 2.29 kg (29.82 kg (s.d 4.18 kg) vs. 27.53 kg (s.d 3.56 kg)). These results

259    demonstrate that the GPS prediction based on the current GWAS finding is capable of

260    identifying subjects of high or low levels of ALM.

261    *Genetic correlations with other traits*

262    To test whether lean mass has a shared genetic etiology with other diseases and relevant traits,

263    a genetic correlation analysis was performed with the LDSC method [17]. Here, ALM studied in

264    our study is strongly genetically correlated with DXA-derived whole body lean mass and the

265    ALM, which were studied by a previous GWAS meta-analysis [15] ($r_g$=0.87 and 0.78) (**Figure

266    5**). Furthermore, ALM is modestly correlated with BMI ($r_g$=0.31). However, the correlation with

267    BMD is low ($r_g$=0.05). ALM is most negatively correlated with BMI adjusted leptin ($r_g$=-0.41). It

268    is also negatively correlated with body fat ($r_g$=-0.17), suggesting a reverse developmental

269    direction towards lean and fat mass.

270    *Mendelian randomization analysis*

271       To investigate whether ALM is causally linked with other complex diseases, a Mendelian

272       randomization analysis was performed with GSMR [26]. Ten diseases from a variety of

273       categories were chose for evaluation, including fracture [36], type 2 diabetes (T2D) [37], asthma

274       [38], insomnia [39], inflammatory bowel disease (IBD) [40], smoking addiction [41], coronary

275       artery disease (CAD) [42], amyotrophic lateral sclerosis (ALS) [43], bipolar disorder [44] and

276       autistic spectrum disorder (ASD) [45]. At the corrected significance level $5\times10^{-3}$ (0.05/10), ALM

277       is causally associated with type 2 diabetes (T2D, $p=4.38\times10^{-8}$) and fracture ($p=1.18\times10^{-3}$), but

278       not with any other disease (**Supplementary Table 15**). Specifically, a negative association is

279       observed between ALM and both diseases, indicating that ALM is a protective factor for both

280       diseases. For T2D, an increase in the unit of one s.d of ALM residual corresponds to a decreased

281       OR of 0.91 (95% CI [0.88, 0.94]). For fracture, an increase in the unit of one s.d. of ALM

282       residual corresponds to a decreased OR of 0.95 (95% CI [0.92, 0.98]).

283

### Discussion

284

285     The incapacity in GWAS to detect and replicate specific genetic variants for human complex

286     traits, contradicting to a trait's established high heritability, e.g., height, was formally recognized

287     as the "missing heritability" problem a decade ago [46, 47].   An explanation is the so called

288     "polygenic model", where hundreds or even thousands of common SNP variants act additively,

289     with each contributing only a "tiny" fraction of the trait variation. The effect of each individual

290     variant is so small that a GWAS with a limited sample size (n<20,000) may be extremely

291     difficult, if not impossible, to detect (let alone replicate) a variant at the genome-wide

292     significance threshold.

293     The polygenic model was supported by the genome-wide complex trait analysis (GCTA),

294     where trait similarity among unrelated subjects was correlated with and explained to a large

295     fraction by similarity of common SNPs at genome-wide scale [48]. Furthermore, with sample

296     sizes at the scale of hundreds of thousands, two GWAS indeed identified at genome-wide

297     significance ~700 variants for adult height [49] and >100 loci for schizophrenia [50]. The

298     successful stories offer a promising prospect for a GWAS with an extraordinarily large sample

299     size to ultimately unravel the puzzling genetic architecture for human complex traits and

300     common diseases.

301     In this study of lean mass with around half million subjects, the largest sample used for

302     GWAS of lean mass so far, a successful endeavor was accomplished again. More than 700

303     variants were identified at the significance of genome-wide scale ($p<5\times10^{-9}$). In particular, more

304     than half of these variants achieved genome-wide significance ($p<5\times10^{-9}$) in one gender and

305     were replicated also in another gender ($p<5\times10^{-5}$). Overall, these >700 variants contributed ~11%

306 of ALM variation, again, the largest explainable fraction of variation for lean mass reported so

307 far in a GWAS.

308 Our findings of >700 variants are expected for a complex trait with a high heritability,

309 particularly considering another trait with comparable heritability, height, which detected also

310 ~700 variants [49]. Interestingly, the majority loci in previous smaller GWAS [13] or

311 meta-analysis [15] of lean mass are also significant in the present study, providing replication

312 evidence from independent samples.

313 The functional relevance of our identified variants is supported by the gene set enrichment

314 analysis, where GO terms, including GO:0001501 'skeletal system development', GO:0061448

315 'connective tissue development', GO:0051216 'cartilage development' and GO:0042692 'muscle

316 cell differentiation', are among the top gene sets of significance. Specifically, the "hub genes"

317 involved in these terms are tightly connected into a network that contains TGF pathway genes,

318 BMP pathway genes and SMAD family genes, which are all important musculoskeletal

319 development genes/pathways. This finding is concordant with developmental biology since cells

320 from bone, cartilage, muscle and fat share the same progenitor, the mesenchymal stem cells, and

321 pleiotropy of muscle and bone is well recognized in both humans [51] and animal models [52].

322 Among the variants identified, those of several genes, such as *SMAD3*, *MEF2C*, *ADAMTS3*

323 and *PAM*, are interesting and may need further investigation. The first two genes are the hub

324 genes involved in half of the significant enriched gene sets. The third gene, *ADAMTS3,* contains

325 8 variants, including 3 rare or less common mis-sense mutations, which in total explains ~0.2%

326 of ALM variation. The fourth gene, *PAM*, has multiple lines of supporting evidence for its

327 regulatory roles, e.g., containing a mis-sense rare SNP rs78408340 (MAF=0.01%, $p=6.10\times10^{-10}$).

328 An intergenic variant, rs400596, is associated with the PAM expression level in whole blood

329    ($p=2.51\times10^{-21}$) [23] and associated with its protein level in peripheral blood tissue ($p=$

330    $p=2.51\times10^{-30}$) [25]. These genes may represent the next candidates for functional and

331    mechanistic analysis of lean mass regulation.

332    In summary, we performed a GWAS using ~half-million subjects for lean mass. Owing to its

333    high statistical power, our study identified a large number of variants mapped to GO terms with

334    functional relevance to musculoskeletal development. The explained variation of ~11% of lean

335    mass by the identified variants represents a significant leap in revealing the "hidden" heritability

336    of this complex trait using GWAS. Our findings' translational value is marked by lean mass'

337    importance to other complex diseases, such as type 2 diabetes and fracture, as our Mendelian

338    randomization analysis showed that ALM is a protective factor for these two diseases. Overall,

339    our study provides another example, where GWAS of substantially increased sample size may

340    lead a way to ultimately and thoroughly delineate genetic architecture of human complex traits.

341    This epitomizes the value of big data in genetic research of humans.

342

343 **Materials and Methods**

344 *Study participants*

345     Study sample came from the UK Biobank (UKB) cohort, which is a large prospective cohort

346 study of ~500,000 participants from across the United Kingdom, aged between 40-69 at

347 recruitment. Ethics approval for the UKB study was obtained from the North West Centre for

348 Research Ethics Committee (11/NW/0382), and informed consent was obtained from all

349 participants. This study (UKB project #41542) was covered by the general ethical approval for

350 the UKB study.

351     All the included subjects are those who self-reported as white (data field 21000). Subjects

352 who had a self-reported gender inconsistent with the genetic gender, who were genotyped but not

353 imputed or who withdraw their consents were removed. The final sample consisted of 450,580

354 subjects, including 244,945 females and 205,635 males.

355 *Phenotype and modeling*

356     Body composition was measured by bioelectrical impedance approach. Appendicular lean

357 mass (ALM) was quantified by the sum of fat-free mass at arms (data fields 23121 and 23125)

358 and at legs (data fields 23113 and 23117). Appendicular fat mass (AFM) was quantified by the

359 sum of fat mass at arms (data fields 23120 and 23124) and at legs (data fields 23112 and 23116).

360 In each gender, covariates including AFM, age, $age^2$, height and $height^2$ were tested for

361 significance in association with ALM using a step-wise linear regression model implemented in

362 the R function stepAIC. Raw ALM values were adjusted by the significant covariates, and the

363 residuals were normalized into inverse quantiles of standard normal distribution, which were

364 used for subsequent association analysis.

365      A small subset of 4,294 subjects also received a dual-energy X-ray absorptiometry (DXA)

366      body composition scan, and hence their DXA-derived ALM is also available. Therefore, raw

367      ALM derived from DXA and from electric impedance was compared in these subjects by

368      Pearson's correlation coefficient.

369      *Genotype quality control*

370      Genome-wide genotypes for all subjects were available at 784,256 genotyped autosome

371      markers, and were imputed into UK10K haplotype, 1000 Genomes project phase 3 and

372      Haplotype Reference Consortium (HRC) reference panels. A total of ~92 million variants were

373      generated by imputation. We excluded variants with MAF<0.1% and with imputation $r^2$<0.3. As

374      a result, ~19.4 million well imputed variants were retained for subsequent genetic association

375      analysis.

376      *Genetic association analysis*

377      In each gender group, we used BOLT-LMM to perform linear mixed model (LMM) analysis

378      [16]. As the LMM analysis can adjust for population structure and relatedness, we included all

379      eligible subjects into analysis, as recommended by BOLT [53]. We did not include principal

380      components (PCs) of ancestry as covariates in the LMM analysis.

381      After sex-specific associations were analyzed, we meta-analyzed the summary statistics of the

382      two genders by inverse-variance weighted fixed-effects model with METAL [54]. The

383      genome-wide significance (GWS) level was set at $\alpha=5\times10^{-9}$, to account for both common and

384      rare variants. The variants that passed this threshold in across-gender meta-analysis were then

385      checked for replicability across genders based on a suggestive significance level $5\times10^{-5}$ in each

386      gender. The suggestive level was set so as to account for multiple testing of presumed maximal

387      number of 1000 independent loci (0.05/1000). An association was defined as "replicated" if the

388  signal was significant at the GWS level ($p<5\times10^{-9}$) in the meta-analysis and was significant at

389  the suggestive level ($p<5\times10^{-5}$) in both genders.

390    This declaration of a replicated association was approximately same as a two-stage design,

391  where the first stage involves selecting variants at the suggestive level ($p<5\times10^{-5}$) in one gender

392  and the second stage involves replicating the selected variants at the same significance level

393  ($p<5\times10^{-5}$) in another gender. An association locus was defined as a genomic region of 500 kb to

394  both sides of a significant lead signal.

395    Difference in effect size between female and male was examined by a two-tailed p-value from

396  the z-score in the following equation

$$z = \frac{\beta_{female} - \beta_{male}}{\sqrt{var(\beta_{female}) + var(\beta_{male})}}$$

397  , where $\beta_{female}$ and $\beta_{male}$ are regression coefficients for females and males, and var($\cdot$) are their

398  variances, respectively.

399  ***Conditional association analysis***

400    To identify additional signals in regions of association, approximate joint and conditional

401  association analysis was performed in each region using the GCTA tool [55].

402

403  From the UKB sample, a reference sample of 100,000 unrelated subjects was generated for

404  estimating LD pattern for subsequent analyses. The unrelated subjects were inferred with KING

405  software [56], from whom the 100,000 subjects of the reference sample were randomly drawn.

406  Quality control (QC) procedures applied to the reference sample included Hardy-Weinberg

407  equilibrium ($p>1\times10^{-6}$) and MAF>0.1%.

408    A recursive conditional association analysis was performed. In each iteration, an approximate

409    conditional analysis conditioning on the current list of lead variants was performed in each

410    gender, followed by an across-gender meta-analysis to combine the gender-specific results.

411    Again, a significant replicated association was defined as achieving both a conditional

412    meta-analysis GWS signal ($p<5\times10^{-9}$) and a conditional suggestive signal ($p<5\times10^{-5}$) in both

413    genders. In addition, each such identified variant is required to be independent of all variants in

414    the lead SNP list (LD $r^2<0.1$). The variant with the lowest p-value among such identified ones

415    was added into the list of lead variants. Iterations of the conditional analysis were run until no

416    significant signal can be identified.

417    *Overlap with loci in previous GWAS of obesity traits*

418    GWAS summary statistics for 5 obesity traits, including body mass index (BMI) [57], waist

419    circumference (WC), WC adjusted for BMI ($WC_{adj}BMI$), waist-hip ratio (WHR) and WHR

420    adjusted for BMI ($WHR_{adj}BMI$) [58], were downloaded from the GIANT consortium website.

421    For each trait, SNPs located within all the 561 identified loci (lead SNP +500 kb flanking region

422    at each side) were extracted from the GWAS summary statistics. Significance level for the

423    obesity traits were set at the conventional level of $5.0\times10^{-8}$.

424    *Exome sequencing association analysis*

425    During the preparation of this manuscript, the UKB released exome-sequencing data on a

426    selected subset of ~50,000 participants. We compared the exome-sequencing based association

427    results with that based on genotype imputation. To accomplish this, we generated an unrelated

428    sample consisting of subjects who were both exome-sequenced and genotype-imputed.

429    As the QC procedure, we removed subjects who were not self-reported as white, whose

430    self-reported genders were inconsistent with their genetic genders, and who withdrew their

431    consents. The KING software was used to select unrelated subjects based on pairwise kinship

432    matrix for up to 2[nd] degree relatedness [56]. The final sample consisted of 45,554 participants,

433    including 24,740 females and 20,814 males.

434    Sequence variant coordinates, which were annotated to the GRCH38 assembly, were

435    converted          back          to          the          GRCH37          assembly          with          Liftover

436    (http://genome.ucsc.edu/cgi-bin/hgLiftOver). For each subject, variants that were missing in the

437    sequenced data were set to missing in the imputed data as well. In both datasets, genetic

438    association with normalized phenotype residuals was analyzed with PLINK2 [59]. The top 10

439    PCs were included as covariates to account for potential population stratification.

*Genetic architecture*

441    BOLT-REML was used to estimate heritability tagged by all the analyzed variants [18]. LD

442    score regression (LDSC) method was used to estimate the amount of genomic inflation due to

443    confounding factors such as population stratification and cryptic relatedness [17]. Pre-computed

444    LD scores from the 1000 Genomes project European subjects were used for estimation. The

445    relative contribution of confounding factors was measured by attenuation ratio (AR), which is

446    defined as $(\text{intercept-1})/(\text{mean chi}^2-1)$, where intercept and mean $\text{chi}^2$ are estimates of

447    confounding and the overall association inflation, respectively [17].

448    To compare AR with that estimated on unrelated subjects, a maximal subset of unrelated

449    subjects from the total sample being analyzed were generated. Specifically, KING was used to

450    extract a subset of unrelated subjects [56]. The resulting unrelated sample included 369,968

451    participants (198,989 females and 170,979 males). In each gender, PLINK2 was used to perform

452    genetic association analysis [59]. To account for genetic confounding, the top 10 PCs inferred

453    from UKB were used as the additional covariates.

454    To calculate the variance explained by all independent lead variants, individual variant effect

455    size was estimated with the formula $2f(1-f)\beta^2$, where $f$ is allele frequency and $\beta$ is regression

456    coefficient associated with the variant.

*Enrichment analysis*

458    Stratified LDSC was used to partition heritability from GWAS summary statistics into

459    different functional categories [21]. The analysis was based on the 'full baseline model' created

460    by Finucane et al. [21] from 24 publicly available main annotations that are not specific to any

461    cell type. Significance level of enrichment was set at $p < 2.08 \times 10^{-3}$ (0.05/24).

462    The stratified LDSC was used to also assess the enrichment of heritability into specific tissues

463    and cell types [22]. This method analyzes gene expression data together with GWAS summary

464    statistics, for which, the two pre-compiled gene expression datasets in LDSC were used. The first

465    one is the GTEx project v6p [23] and the second one is the Franke lab dataset [24]. The GTEx

466    dataset contains 53 tissues with an average of 161 samples per tissue. The Franke lab dataset is

467    an aggregation of publicly available microarray gene expression datasets comprising 37,427

468    human samples from 152 tissues. The total 205 (=53+152) tissues are classified into nine

469    categories for visualization. Significance was declared at a false discovery rate (FDR)<5%.

*Candidate gene prioritization*

471    In each associated locus, a set of credible risk variants (CRVs) were defined as those variants

472    in strong LD with the lead variant ($r^2>0.8$, including lead variant). LD $r^2$ measure was estimated

473    based on the 100,000 unrelated reference sample with LDstore [60]. Six sources of information

474    was used to evaluate a gene's causality: 1) being nearest to the lead CRV; 2) containing a

475    mis-sense coding CRV; 3) being a target gene for a cis-eQTL CRV; 4) being a target gene for a

476    cis- protein QTL (cis-pQTL) CRV; 5) being prioritized by DEPICT analysis [27] and 6) being

477    prioritized by SMR analysis [61].

478        Cis-eQTLs revealed by the GTEx (v7) project were accessed from the GTEx web portal

479    (www.gtexportal.org/) [23]. Cis-eQTL information is available for over 50 tissues. We selected

480    skeletal muscle and whole blood for our analysis. Cis-eQTL was searched within 500 kb distance

481    from a target gene. Significant cis-eQTL was declared at $p < 5 \times 10^{-5}$.

482        Cis-pQTL information was accessed from Sun et al. [25]. GWAS summary statistics for 3,284

483    proteins were downloaded from the study's website. Cis-pQTL was searched within 500 kb

484    distance from a target gene. Significant cis-eQTL was declared at $p < 5 \times 10^{-5}$.

485        DEPICT is an integrative tool that takes advantage of predicted gene functions to

486    systematically prioritize the most likely causal genes at loci of interest [27]. The input of

487    DEPICT includes a list of variant identifiers, and the output contains all genes located in the loci

488    and their p-values of being a candidate gene. All lead variants were submitted to DEPICT for

489    analysis. Significant genes were declared at a false discovery rate (FDR)<5%.

490        SMR (Summary data–based Mendelian Randomization) method [61] is another SNP

491    prioritization program that integrates summary-level data from GWAS with data from eQTL

492    studies to identify genes whose expression levels are associated with trait due to causal or

493    pleiotropy effects. Here, the pleiotropy effect means that a SNP is causally associated with both

494    gene expression and phenotypic variation. SMR uses SNPs as an instrumental variable and tests

495    the causal relation of gene expression to phenotype variation. The results are interpreted as the

496    effect of gene expression on phenotype free of confounding from non-genetic factors. We used a

497    pre-compiled eQTL dataset in whole blood tissue [62] for estimation.

498    ***Gene-based and gene set enrichment analyses***

499    Gene-based association analysis was performed with MAGMA v1.6 [63], as implemented on

500    the FUMA website (http://fuma.ctglab.nl/). GWAS meta-analysis summary statistics were

501    mapped to 19,427 protein-coding genes, resulting in 19,098 genes that were covered by at least

502    one SNP. Gene-based association test was performed taking into account the LD between

503    variants. Gene-based significance level was set at stringent Bonferroni corrected threshold

504    $2.62 \times 10^{-6}$, i.e., 0.05/19,098.

505    The generated gene-based summary statistics were further used to test for enrichment of

506    association to specific biological pathways or gene sets. A gene set's association signal was

507    evaluated by integrating all signals from the genes in the set with MAGMA [63]. A competitive

508    gene set analysis model was used to test whether the genes in a gene set are more strongly

509    associated with the phenotype than other genes.

510    Gene        sets        were        obtained        through        the        MSigDB        website

511    (http://software.broadinstitute.org/gsea/msigdb/index.jsp) [64]. Each gene was assigned to a gene

512    set as annotated by gene ontology (GO) , Kyoto encyclopedia of genes and genomes (KEGG),

513    Reactome and BioCarta gene set databases and other gene sets curated by domain experts or

514    biomedical literature [64]. A total of 10,651 (4,734 curated and 5,917 GO terms) gene sets were

515    used in this analysis. The significance level was set at a Bonferroni-corrected level of

516    $0.05/10,651 = 4.69 \times 10^{-6}$.

517    Protein-protein interaction network was constructed with STRING [65]. STRING uses

518    information based on gene co-expression, text-mining, and others, to construct protein interactive

519    network.

520    ***Polygenic risk score profiling***

521      To assess the capability of the GWAS finding to predict ALM, a polygenic scoring analysis

522      was conducted in the 369,968 unrelated subjects extracted from the main UKB sample. Three

523      quarters of the individuals (277,762 subjects, including 149,329 females) were randomly selected

524      as the training sample, and the remaining one quarter individuals (92,206 participants, including

525      49,660 females) as the validation sample. Female and male subjects were pooled together for

526      analysis.

527      Raw phenotype was adjusted by age, $age^2$, gender, height, $height^2$ and the top 10 PCs, and the

528      residuals were converted to the standard normal distribution quantiles for downstream analysis.

529      Genetic association analysis was performed with PLINK2 [59].

530      The same QC procedures as in the main analysis were used to process the variants. The

531      variants achieving a p-value of $<1\times10^{-5}$ in the training sample were selected and used for

532      prediction in the validation sample via LDpred approach [66]. LDpred infers the posterior mean

533      effect size of each marker by using a prior on effect sizes and LD information from an external

534      reference panel. Specifically, the validation sample with original genotypes was used as

535      reference panel for LD estimation. The number of SNPs used to adjust LD from each side of the

536      target SNP was set to 1000. Other software parameters were set to the default.

537      ***Genetic correlations with other traits***

538      To test whether lean mass has a shared genetic etiology with other diseases and relevant traits,

539      a genetic correlation analysis was performed with LDSC method [17]. An online web tool,

540      LDHub (http://ldsc.broadinstitute.org/ldhub/), was used to estimate the genetic correlation

541      between $ALM_{adj}$ and 49 complex traits and diseases. The standalone version of the software was

542      used to estimate between $ALM_{adj}$ and two additional traits, ALM and total body lean mass,

543      measured by the DXA scan, which are not available in the LDHub GWAS summary statistics

544    collections, and which were downloaded from the GEFOS consortium website

545    (http://www.gefos.org).

546    Both the LDHub and standalone analyses adopted same QC criteria. Specifically, only

547    HapMap3 autosomal SNPs were included to minimize poor imputation quality [17]. SNPs were

548    further removed given the following conditions: MAF<0.01, ambiguous strand (A/T or C/G),

549    duplicated identifier, or reported sample size less than 60% of total sample size. LD scores

550    pre-computed on the 1000 genomes project European individuals were used for calculation.

551    *Mendelian randomization analysis*

552    To investigate whether ALM (as exposure) is causally associated with complex diseases (as

553    outcome), a Mendelian randomization analysis with GSMR was performed [26] on selected 10

554    complex diseases, including fracture [36], type 2 diabetes (T2D) [37], asthma [38], insomnia

555    [39], inflammatory bowel disease (IBD) [40], smoking addiction [41], coronary artery disease

556    (CAD) [42], amyotrophic lateral sclerosis (ALS) [43], bipolar disorder [44] and autistic spectrum

557    disorder (ASD) [45].

558    GWAS summary statistics for these diseases were downloaded from the respective websites.

559    From the list of SNPs whose association signals with $ALM_{adj}$ were below $5 \times 10^{-8}$, qualified SNPs

560    were included based on the following criteria: concordant alleles between exposure and outcome

561    GWAS summary statistics, non-palindromic SNPs with certain strand, MAF>1%, and allele

562    frequency difference between exposure and outcome GWAS summary statistics <0.2.

563    Independent SNPs were further clumped with PLINK2 [59] with independence LD threshold

564    $r^2$<0.05 and 1 MB window size. The clumped independent SNPs were examined for their

565    pleiotropic effects to both exposure and outcome by the HEIDI test [26]. Significance level for

566    the HEIDI test was set to $\alpha = 1 \times 10^{-5}$. After removing pleiotropic SNPs, the remaining independent

567     SNPs were taken as instrumental variables to test for the causal effect of exposure to outcome.

568     The estimated causal effect coefficients are approximately equal to the natural log odds ratio (OR)

569     for a case–control trait. The MR analysis significance level was set to 0.005 (0.05/10).

570

577 **References**
578 1. Giles, J.T., et al., *Association of body composition with disability in rheumatoid arthritis:*
579 *impact of appendicular fat and lean tissue mass.* Arthritis Rheum, 2008. **59**(10): p.
580 1407-15.
581 2. Janssen, I., S.B. Heymsfield, and R. Ross, *Low relative skeletal muscle mass (sarcopenia)*
582 *in older persons is associated with functional impairment and physical disability.* J Am
583 Geriatr Soc, 2002. **50**(5): p. 889-96.
584 3. Miyakoshi, N., et al., *Prevalence of sarcopenia in Japanese women with osteopenia and*
585 *osteoporosis.* J Bone Miner Metab, 2013. **31**(5): p. 556-61.
586 4. Binkley, N., D. Krueger, and B. Buehring, *What's in a name revisited: should osteoporosis*
587 *and sarcopenia be considered components of "dysmobility syndrome?".* Osteoporos Int,
588 2013. **24**(12): p. 2955-9.
589 5. Wannamethee, S.G. and J.L. Atkins, *Muscle loss and obesity: the health implications of*
590 *sarcopenia and sarcopenic obesity.* Proc Nutr Soc, 2015. **74**(4): p. 405-12.
591 6. Evans, W.J., et al., *Cachexia: a new definition.* Clin Nutr, 2008. **27**(6): p. 793-9.
592 7. Brown, J.C., M.O. Harhay, and M.N. Harhay, *Sarcopenia and mortality among a*
593 *population-based sample of community-dwelling older adults.* J Cachexia Sarcopenia
594 Muscle, 2016. **7**(3): p. 290-8.
595 8. Arden, N.K. and T.D. Spector, *Genetic influences on muscle strength, lean body mass,*
596 *and bone mineral density: a twin study.* J Bone Miner Res, 1997. **12**(12): p. 2076-81.
597 9. Livshits, G., et al., *Contribution of Heritability and Epigenetic Factors to Skeletal Muscle*
598 *Mass Variation in United Kingdom Twins.* J Clin Endocrinol Metab, 2016. **101**(6): p.
599 2450-9.
600 10. Liu, X.G., et al., *Genome-wide association and replication studies identified TRHR as an*
601 *important gene for lean body mass.* Am J Hum Genet, 2009. **84**(3): p. 418-23.
602 11. Ran, S., et al., *Gene-based genome-wide association study identified 19p13.3 for lean*
603 *body mass.* Sci Rep, 2017. **7**: p. 45025.
604 12. Hai, R., et al., *Bivariate genome-wide association study suggests that the DARC gene*
605 *influences lean body mass and age at menarche.* Sci China Life Sci, 2012. **55**(6): p.
606 516-20.
607 13. Urano, T., et al., *Large-scale analysis reveals a functional single-nucleotide*
608 *polymorphism in the 5'-flanking region of PRDM16 gene associated with lean body mass.*
609 Aging Cell, 2014. **13**(4): p. 739-43.
610 14. Klimentidis, Y.C., et al., *Genetic Variant in ACVR2B Is Associated with Lean Mass.* Med Sci
611 Sports Exerc, 2016. **48**(7): p. 1270-5.
612 15. Zillikens, M.C., et al., *Large meta-analysis of genome-wide association studies identifies*
613 *five loci for lean body mass.* Nat Commun, 2017. **8**(1): p. 80.
614 16. Loh, P.R., et al., *Efficient Bayesian mixed-model analysis increases association power in*
615 *large cohorts.* Nat Genet, 2015. **47**(3): p. 284-90.
616 17. Bulik-Sullivan, B.K., et al., *LD Score regression distinguishes confounding from*
617 *polygenicity in genome-wide association studies.* Nature Genetics, 2015. **47**(3): p.
618 291-295.
619 18. Loh, P.R., et al., *Contrasting genetic architectures of schizophrenia and other complex*
620 *diseases using fast variance-components analysis.* Nat Genet, 2015. **47**(12): p. 1385-92.

621  19.  Medina-Gomez, C., et al., *Bivariate genome-wide association meta-analysis of pediatric*
622       *musculoskeletal traits reveals pleiotropic effects at the SREBF1/TOM1L2 locus.* Nat
623       Commun, 2017. **8**(1): p. 121.
624  20.  Hubel, C., et al., *Genomics of body fat percentage may contribute to sex bias in anorexia*
625       *nervosa.* Am J Med Genet B Neuropsychiatr Genet, 2018.
626  21.  Finucane, H.K., et al., *Partitioning heritability by functional annotation using*
627       *genome-wide association summary statistics.* Nat Genet, 2015. **47**(11): p. 1228-35.
628  22.  Finucane, H.K., et al., *Heritability enrichment of specifically expressed genes identifies*
629       *disease-relevant tissues and cell types.* Nat Genet, 2018. **50**(4): p. 621-629.
630  23.  Consortium, G.T., *Human genomics. The Genotype-Tissue Expression (GTEx) pilot*
631       *analysis: multitissue gene regulation in humans.* Science, 2015. **348**(6235): p. 648-60.
632  24.  Fehrmann, R.S., et al., *Gene expression analysis identifies global gene dosage sensitivity*
633       *in cancer.* Nat Genet, 2015. **47**(2): p. 115-25.
634  25.  Sun, B.B., et al., *Genomic atlas of the human plasma proteome.* Nature, 2018. **558**(7708):
635       p. 73-79.
636  26.  Zhu, Z., et al., *Causal associations between risk factors and common diseases inferred*
637       *from GWAS summary data.* Nat Commun, 2018. **9**(1): p. 224.
638  27.  Pers, T.H., et al., *Biological interpretation of genome-wide association studies using*
639       *predicted gene functions.* Nat Commun, 2015. **6**: p. 5890.
640  28.  Adzhubei, I.A., et al., *A method and server for predicting damaging missense mutations.*
641       Nat Methods, 2010. **7**(4): p. 248-9.
642  29.  Kumar, P., S. Henikoff, and P.C. Ng, *Predicting the effects of coding non-synonymous*
643       *variants on protein function using the SIFT algorithm.* Nat Protoc, 2009. **4**(7): p. 1073-81.
644  30.  Choi, Y., et al., *Predicting the functional effect of amino acid substitutions and indels.*
645       PLoS One, 2012. **7**(10): p. e46688.
646  31.  Shihab, H.A., et al., *Predicting the functional, molecular, and phenotypic consequences of*
647       *amino acid substitutions using hidden Markov models.* Hum Mutat, 2013. **34**(1): p.
648       57-65.
649  32.  Schmidt, K., et al., *Sox8 is a specific marker for muscle satellite cells and inhibits*
650       *myogenesis.* J Biol Chem, 2003. **278**(32): p. 29769-75.
651  33.  Ge, X., et al., *Lack of Smad3 signaling leads to impaired skeletal muscle regeneration.*
652       Am J Physiol Endocrinol Metab, 2012. **303**(1): p. E90-102.
653  34.  Estrella, N.L., et al., *MEF2 transcription factors regulate distinct gene programs in*
654       *mammalian skeletal muscle differentiation.* J Biol Chem, 2015. **290**(2): p. 1256-68.
655  35.  Frank, N.Y., et al., *Regulation of myogenic progenitor proliferation in human fetal*
656       *skeletal muscle by BMP4 and its antagonist Gremlin.* J Cell Biol, 2006. **175**(1): p. 99-110.
657  36.  Trajanoska, K., et al., *Assessment of the genetic and clinical determinants of fracture risk:*
658       *genome wide association and mendelian randomisation study.* BMJ, 2018. **362**: p. k3225.
659  37.  Xue, A., et al., *Genome-wide association analyses identify 143 risk variants and putative*
660       *regulatory mechanisms for type 2 diabetes.* Nat Commun, 2018. **9**(1): p. 2941.
661  38.  Moffatt, M.F., et al., *A large-scale, consortium-based genomewide association study of*
662       *asthma.* N Engl J Med, 2010. **363**(13): p. 1211-1221.

663  39.  Hammerschlag, A.R., et al., *Genome-wide association analysis of insomnia complaints*
664       *identifies risk genes and genetic overlap with psychiatric and metabolic traits.* Nat Genet,
665       2017. **49**(11): p. 1584-1592.
666  40.  Liu, J.Z., et al., *Association analyses identify 38 susceptibility loci for inflammatory bowel*
667       *disease and highlight shared genetic risk across populations.* Nat Genet, 2015. **47**(9): p.
668       979-986.
669  41.  *Genome-wide meta-analyses identify multiple loci associated with smoking behavior.*
670       Nat Genet, 2010. **42**(5): p. 441-7.
671  42.  van der Harst, P. and N. Verweij, *Identification of 64 Novel Genetic Loci Provides an*
672       *Expanded View on the Genetic Architecture of Coronary Artery Disease.* Circ Res, 2018.
673       **122**(3): p. 433-443.
674  43.  van Rheenen, W., et al., *Genome-wide association analyses identify new risk variants*
675       *and the genetic architecture of amyotrophic lateral sclerosis.* Nat Genet, 2016. **48**(9): p.
676       1043-8.
677  44.  Stahl, E.A., et al., *Genome-wide association study identifies 30 Loci Associated with*
678       *Bipolar Disorder.* bioRxiv, 2018: p. 173062.
679  45.  Grove, J., et al., *Identification of common genetic risk variants for autism spectrum*
680       *disorder.* Nat Genet, 2019. **51**(3): p. 431-444.
681  46.  Manolio, T.A., et al., *Finding the missing heritability of complex diseases.* Nature, 2009.
682       **461**(7265): p. 747-53.
683  47.  Zuk, O., et al., *The mystery of missing heritability: Genetic interactions create phantom*
684       *heritability.* Proc Natl Acad Sci U S A, 2012. **109**(4): p. 1193-8.
685  48.  Yang, J., et al., *Common SNPs explain a large proportion of the heritability for human*
686       *height.* Nat Genet, 2010. **42**(7): p. 565-9.
687  49.  Wood, A.R., et al., *Defining the role of common variation in the genomic and biological*
688       *architecture of adult human height.* Nat Genet, 2014. **46**(11): p. 1173-86.
689  50.  Consortium, S.W.G.o.t.P.G., *Biological insights from 108 schizophrenia-associated*
690       *genetic loci.* Nature, 2014. **511**(7510): p. 421-7.
691  51.  Karasik, D. and D.P. Kiel, *Genetics of the musculoskeletal system: a pleiotropic approach.*
692       J Bone Miner Res, 2008. **23**(6): p. 788-802.
693  52.  Blank, R.D., *Bone and Muscle Pleiotropy: The Genetics of Associated Traits.* Clin Rev Bone
694       Miner Metab, 2014. **12**(2): p. 61-65.
695  53.  Loh, P.R., et al., *Mixed-model association for biobank-scale datasets.* Nat Genet, 2018.
696       **50**(7): p. 906-908.
697  54.  Sanna, S., et al., *Common variants in the GDF5-UQCC region are associated with*
698       *variation in human height.* Nat Genet, 2008. **40**(2): p. 198-203.
699  55.  Yang, J., et al., *GCTA: a tool for genome-wide complex trait analysis.* Am J Hum Genet,
700       2011. **88**(1): p. 76-82.
701  56.  Manichaikul, A., et al., *Robust relationship inference in genome-wide association studies.*
702       Bioinformatics, 2010. **26**(22): p. 2867-73.
703  57.  Yengo, L., et al., *Meta-analysis of genome-wide association studies for height and body*
704       *mass index in approximately 700000 individuals of European ancestry.* Hum Mol Genet,
705       2018. **27**(20): p. 3641-3649.

706  58.  Shungin, D., et al., *New genetic loci link adipose and insulin biology to body fat distribution.* Nature, 2015. **518**(7538): p. 187-196.
708  59.  Chang, C.C., et al., *Second-generation PLINK: rising to the challenge of larger and richer datasets.* Gigascience, 2015. **4**: p. 7.
710  60.  Benner, C., et al., *Prospects of Fine-Mapping Trait-Associated Genomic Regions by Using Summary Statistics from Genome-wide Association Studies.* Am J Hum Genet, 2017. **101**(4): p. 539-551.
713  61.  Zhu, Z., et al., *Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets.* Nat Genet, 2016. **48**(5): p. 481-7.
715  62.  Westra, H.J., et al., *Systematic identification of trans eQTLs as putative drivers of known disease associations.* Nat Genet, 2013. **45**(10): p. 1238-1243.
717  63.  de Leeuw, C.A., et al., *MAGMA: generalized gene-set analysis of GWAS data.* PLoS Comput Biol, 2015. **11**(4): p. e1004219.
719  64.  Subramanian, A., et al., *Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles.* Proc Natl Acad Sci U S A, 2005. **102**(43): p. 15545-50.
722  65.  Szklarczyk, D., et al., *STRING v11: protein-protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets.* Nucleic Acids Res, 2019. **47**(D1): p. D607-D613.
725  66.  Vilhjalmsson, B.J., et al., *Modeling Linkage Disequilibrium Increases Accuracy of Polygenic Risk Scores.* Am J Hum Genet, 2015. **97**(4): p. 576-92.

727
728
729

730 **Figure Legends**

731 **Figure 1. Main association results.**

732 **Figure 1A**, Per allele effect size versus minor allele frequency (MAF). X-axis is MAF at the 717

733 identified variants and y-axis is per allele effect size (regression coefficient). **Figure 1B**, the Manhattan

734 plot of the meta-analysis combining both genders. The horizontal red line indicates the genome-wide

735 significance level (alpha= $5\times10^{-9}$) in $-\log_{10}$ scale. All novel loci were marked in green.

736

737 **Figure 2. Heritability enrichment in different functional annotations and tissues.**

738 **Figure 2A** is enrichment of genome-wide association signals in 24 main annotations using LDSC

739 regression. Y-axis represents the ratio of phenotypic variance explained by variants in a particular

740 annotation category against that explained in the remaining regions. Error bars represent jackknife

741 standard errors around the estimates of enrichment. A single asterisk indicates significance at $p<0.05$ after

742 Bonferroni correction for the 24 hypotheses tested, and two asterisks indicates significance at $p<0.01$.

743 **Figure 2B** is enrichment of genome-wide association signals in 206 cells/tissues from two different

744 databases (Franke lab dataset and GTEx consortium dataset). The total cells/tissues were divided into

745 9 categories. Each dot represents a specific cell/tissue and the tissues passing the cutoff of FDR < 5% at

746 $-\log10 (p) = 2.75$ were marked in large.

747

748 **Figure 3. Protein-protein interactional network.**

749 Thirty-four genes over-represented in 85 significant pathways were selected to construct a protein-protein

750 interaction network with STRING, which bases the construction on knowledge of gene co-expression,

751 text-mining, and others.

752

753 **Figure 4. Polygenic score prediction.**

754    A total of 277,762 subjects were randomly selected as the training sample, and another

755    independent 92,206 subjects were selected as the validation sample. The variants achieving a

756    p-value of $<1\times10^{-5}$ in the training sample were selected and used for prediction in the validation

757    sample via LDpred approach. Subjects in the two extreme tails of the predicted genome-wide

758    polygenic score (GPS) distribution were compared in terms of raw phenotype mean (after

759    correction). X-axis represents the fraction of subjects drawn from both extreme tails of the predicted

760    GPS distribution. Y-axis represents mean $ALM_{adj}$ ($\pm95\%$ confidence interval).

761

762    **Figure 5. Genetic overlap with other traits.**

763    Genetic correlations ($r_g$) between $ALM_{adj}$ and 51 traits and diseases were estimated. LD Score regression

764    tested genome-wide SNP associations for these participants against similar data for various other traits

765    and diseases containing Musculoskeletal system, anthropometrics, obesity, cognition, metabolism,

766    psychiatry, reproduction and neuropsychiatric outcomes. Error bars represent standard errors on these

767    estimates. Blue bars represent significantly positive correlation at the nominal level p<0.05; pink bars

768    represent significantly negative correlation (p<0.05); grey bars represent non-significant correlation.

769
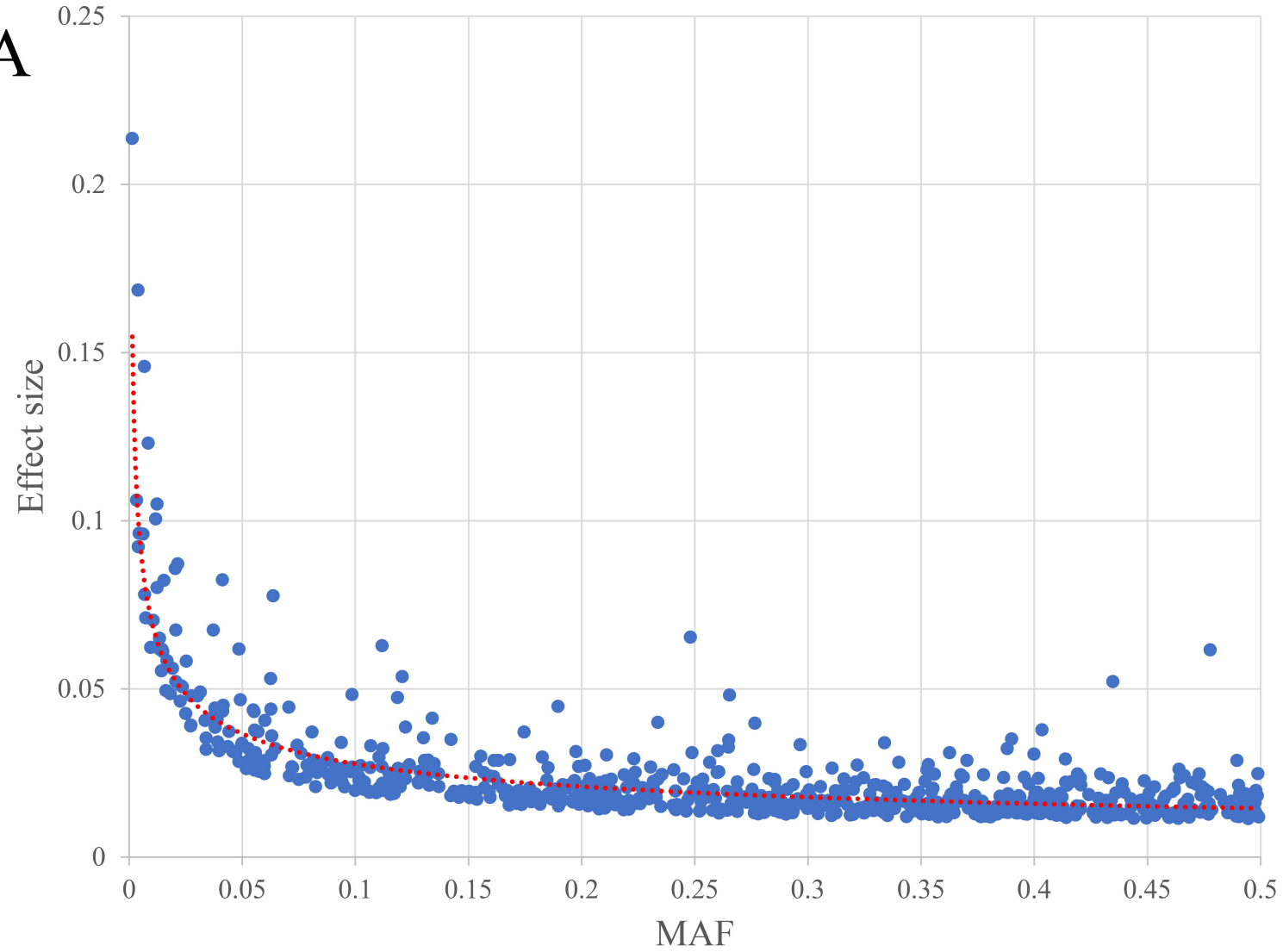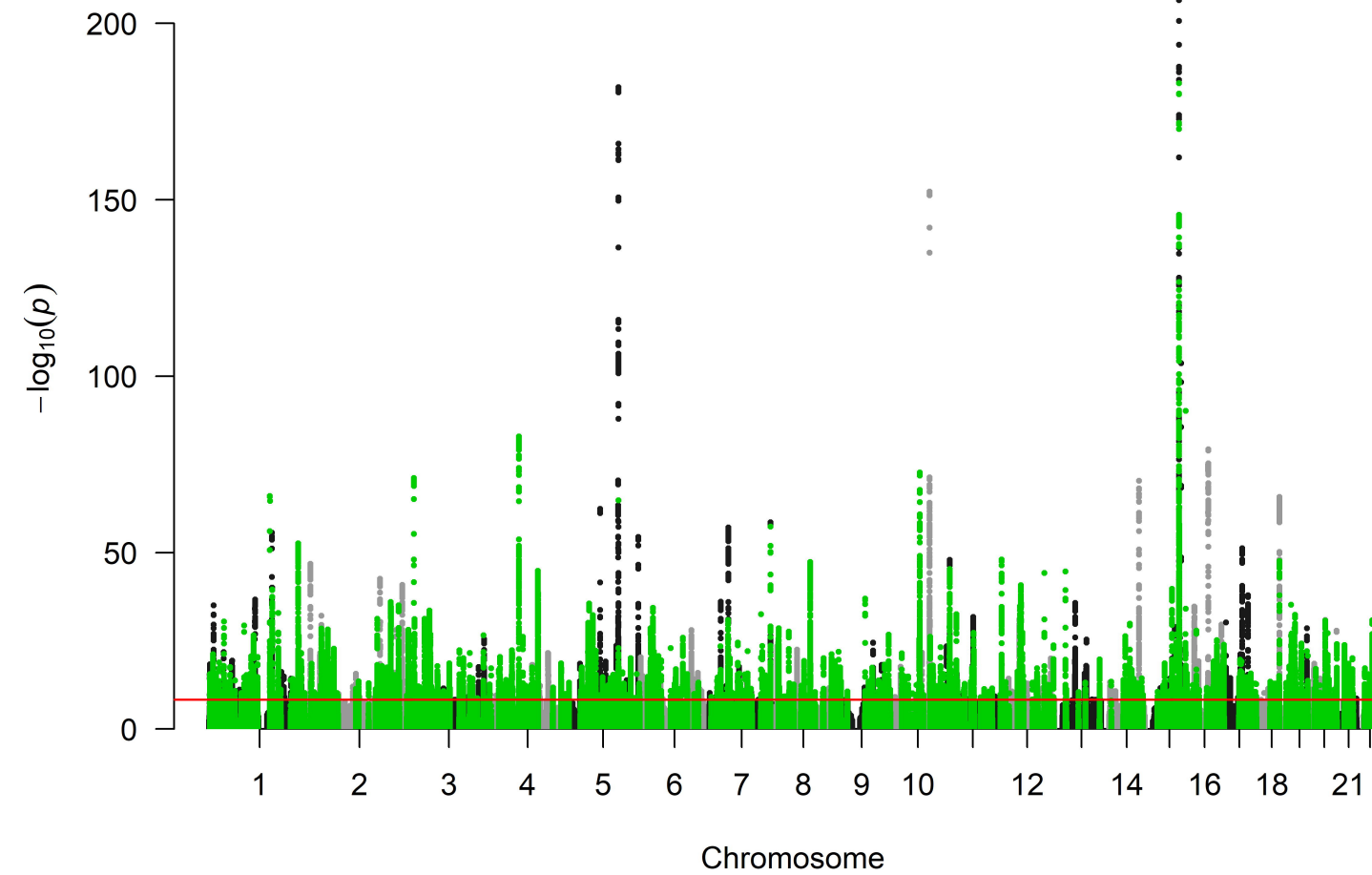
## Table 1. Association results of 42 mis-sense variants.

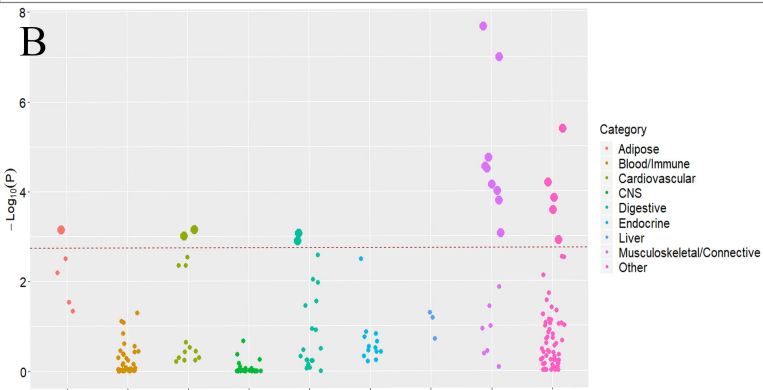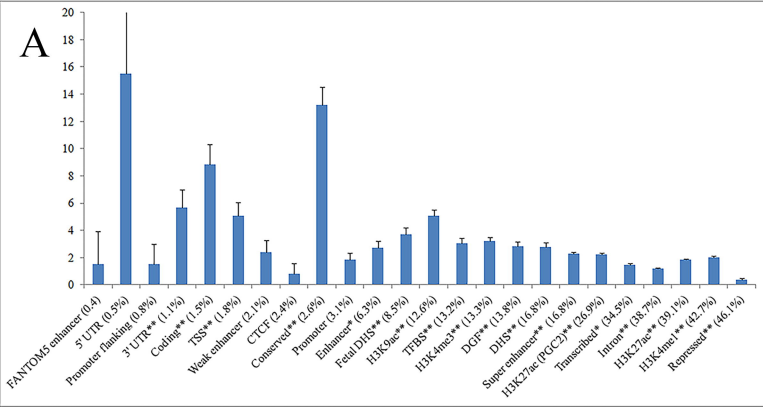| RSID | CHR | POS | BAND | Alleles (REF/ALT) | FRQ | Gene | Protein change | Condition | Female (N=244,945) | | | Male (N=205,635) | | | Meta-analysis (N=450,580) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | B | SE | P | B | SE | P | B | SE | P | H2 | I2 |
| Rare | | | | | | | | | | | | | | | | | | | |
| rs148330006 | 1 | 86048526 | 1p22.3 | C/G | 0.008 | CYR61 | p.Ser316Cys | Primary | 0.14 | 0.015 | $4.70\times10^{-22}$ | 0.1 | 0.016 | $1.80\times10^{-11}$ | 0.12 | 0.011 | $4.31\times10^{-30}$ | $2.52\times10^{-4}$ | 60.3 |
| rs200219556 | 2 | 241974013 | 2q37.3 | G/A | 0.007 | SNED1 | p.Arg224His | Secondary | 0.07 | 0.016 | $8.85\times10^{-6}$ | 0.07 | 0.017 | $1.13\times10^{-5}$ | -0.07 | 0.011 | $4.28\times10^{-10}$ | $7.45\times10^{-5}$ | 0 |
| rs141374503 | 4 | 73179445 | 4q13.3 | C/T | 0.004 | ADAMTS3 | p.Arg565Gln | Secondary | -0.17 | 0.021 | $4.50\times10^{-16}$ | -0.16 | 0.023 | $6.17\times10^{-13}$ | 0.17 | 0.016 | $2.02\times10^{-27}$ | $2.26\times10^{-4}$ | 0 |
| rs148833559 | 5 | 172755066 | 5q35.2 | C/A | 0.001 | STC2 | p.Arg44Leu | Primary | -0.17 | 0.035 | $8.90\times10^{-7}$ | -0.26 | 0.039 | $4.50\times10^{-12}$ | 0.21 | 0.026 | $2.34\times10^{-16}$ | $1.28\times10^{-4}$ | 63.1 |
| rs138940563 | 5 | 79375038 | 5q14.1 | C/T | 0.003 | THBS4 | p.Ala823Val | Primary | -0.09 | 0.023 | $3.70\times10^{-5}$ | -0.13 | 0.025 | $8.80\times10^{-8}$ | 0.11 | 0.017 | $2.10\times10^{-10}$ | $7.64\times10^{-5}$ | 36.9 |
| rs78727187 | 5 | 127668685 | 5q23.3 | G/T | 0.006 | FBN2 | p.His1381Asn | Secondary | 0.07 | 0.017 | $2.68\times10^{-5}$ | 0.12 | 0.019 | $6.81\times10^{-11}$ | -0.1 | 0.013 | $6.67\times10^{-14}$ | $1.14\times10^{-4}$ | 75.2 |
| rs78408340 | 5 | 102338739 | 5q21.1 | C/G | 0.009 | PAM | p.Ser539Trp | Secondary | 0.06 | 0.014 | $5.45\times10^{-6}$ | 0.06 | 0.015 | $2.70\times10^{-5}$ | 0.06 | 0.01 | $6.10\times10^{-10}$ | $7.33\times10^{-5}$ | 0 |
| rs139237114 | 9 | 110249816 | 9q31.2 | G/A | 0.004 | KLF4 | p.His287Tyr | Primary | -0.08 | 0.02 | $2.62\times10^{-5}$ | -0.11 | 0.022 | $5.60\times10^{-7}$ | 0.1 | 0.015 | $9.79\times10^{-11}$ | $8.12\times10^{-5}$ | 0 |
| Less common | | | | | | | | | | | | | | | | | | | |
| rs150270324 | 4 | 73178175 | 4q13.3 | T/C | 0.013 | ADAMTS3 | p.Asn585Ser | Secondary | -0.07 | 0.012 | $1.31\times10^{-9}$ | -0.06 | 0.013 | $2.97\times10^{-6}$ | -0.07 | 0.009 | $2.36\times10^{-14}$ | $1.12\times10^{-4}$ | 0 |
| rs139921635 | 4 | 73181637 | 4q13.3 | G/T | 0.024 | ADAMTS3 | p.Pro513Thr | Secondary | -0.06 | 0.009 | $5.58\times10^{-11}$ | -0.04 | 0.01 | $7.86\times10^{-6}$ | 0.05 | 0.007 | $4.06\times10^{-15}$ | $1.18\times10^{-4}$ | 20.9 |
| rs11722554 | 4 | 5016883 | 4p16.2 | G/A | 0.038 | CYTL1 | p.Arg136Cys | Primary | -0.05 | 0.007 | $1.60\times10^{-10}$ | -0.04 | 0.008 | $1.20\times10^{-8}$ | 0.04 | 0.005 | $5.02\times10^{-18}$ | $1.44\times10^{-4}$ | 0 |
| rs62621812 | 7 | 127015083 | 7q31.33 | G/A | 0.02 | ZNF800 | p.Pro103Ser | Primary | -0.1 | 0.01 | $6.90\times10^{-26}$ | -0.07 | 0.011 | $4.20\times10^{-10}$ | 0.09 | 0.007 | $2.74\times10^{-33}$ | $2.94\times10^{-4}$ | 80.3 |
| rs117874826 | 11 | 64027666 | 11q13.1 | A/C | 0.015 | PLCB3 | p.Glu564Ala | Primary | 0.07 | 0.013 | $1.40\times10^{-7}$ | 0.06 | 0.014 | $1.60\times10^{-5}$ | 0.06 | 0.009 | $2.23\times10^{-11}$ | $1.10\times10^{-4}$ | 0 |
| rs61688134 | 12 | 22017410 | 12p12.1 | C/T | 0.014 | ABCC9 | p.Val734Ile | Primary | 0.06 | 0.011 | $3.90\times10^{-7}$ | 0.05 | 0.013 | $1.50\times10^{-5}$ | -0.06 | 0.009 | $6.91\times10^{-11}$ | $8.65\times10^{-5}$ | 0 |
| rs78457529 | 16 | 24950880 | 16p12.1 | C/T | 0.012 | ARHGAP17 | p.Arg510Gln | Secondary | 0.12 | 0.012 | $1.08\times10^{-20}$ | 0.08 | 0.014 | $1.16\times10^{-9}$ | -0.1 | 0.009 | $4.37\times10^{-28}$ | $2.34\times10^{-4}$ | 69.8 |
| rs60782127 | 16 | 16142079 | 16p13.11 | G/T | 0.014 | ABCC1 | p.Arg433Ser | Primary | -0.07 | 0.011 | $1.60\times10^{-9}$ | -0.06 | 0.013 | $2.60\times10^{-5}$ | 0.06 | 0.009 | $2.91\times10^{-13}$ | $1.02\times10^{-4}$ | 0 |
| rs34934920 | 19 | 38976655 | 19q13.2 | C/T | 0.025 | RYR1 | p.Pro1787Leu | Primary | 0.05 | 0.009 | $1.10\times10^{-7}$ | 0.04 | 0.009 | $1.90\times10^{-5}$ | -0.04 | 0.006 | $1.72\times10^{-11}$ | $8.89\times10^{-5}$ | 0 |
| rs62621197 | 19 | 8670147 | 19p13.2 | C/T | 0.037 | ADAMTS10 | p.Arg62Gln | Primary | -0.07 | 0.007 | $1.20\times10^{-21}$ | -0.06 | 0.008 | $2.10\times10^{-15}$ | 0.07 | 0.005 | $5.60\times10^{-36}$ | $3.27\times10^{-4}$ | 0 |
| rs78648341 | 20 | 19915770 | 20p11.23 | G/A | 0.015 | RIN2 | p.Gly29Arg | Primary | -0.07 | 0.011 | $3.10\times10^{-9}$ | -0.05 | 0.013 | $1.20\times10^{-5}$ | 0.06 | 0.008 | $4.05\times10^{-13}$ | $1.10\times10^{-4}$ | 0 |
| Common | | | | | | | | | | | | | | | | | | | |
| rs3850625 | 1 | 201016296 | 1q32.1 | G/A | 0.118 | CACNA1S | p.Arg1539Cys | Primary | 0.03 | 0.004 | $1.20\times10^{-9}$ | 0.02 | 0.004 | $6.40\times10^{-6}$ | -0.02 | 0.003 | $1.26\times10^{-14}$ | $1.13\times10^{-4}$ | 7.3 |
| rs1047891 | 2 | 211540507 | 2q34 | C/A | 0.316 | CPS1 | p.Thr1412Asn | Primary | -0.03 | 0.003 | $6.70\times10^{-26}$ | -0.02 | 0.003 | $1.40\times10^{-6}$ | 0.02 | 0.002 | $8.41\times10^{-29}$ | $2.35\times10^{-4}$ | 92 |
| rs1260326 | 2 | 27730940 | 2p23.3 | C/T | 0.396 | GCKR | p.Leu446Pro | Primary | -0.02 | 0.003 | $1.60\times10^{-19}$ | -0.02 | 0.003 | $2.20\times10^{-14}$ | -0.02 | 0.002 | $7.83\times10^{-33}$ | $2.71\times10^{-4}$ | 0 |
| rs11545169 | 3 | 184020542 | 3q27.1 | G/T | 0.161 | PSMD2 | p.Glu313Asp | Primary | 0.03 | 0.004 | $7.00\times10^{-9}$ | 0.03 | 0.004 | $1.60\times10^{-15}$ | -0.03 | 0.003 | $2.56\times10^{-27}$ | $2.24\times10^{-4}$ | 0 |
| rs123509 | 3 | 42733468 | 3p22.1 | C/T | 0.248 | KLHL40 | p.Cys617Arg | Primary | 0.02 | 0.003 | $1.20\times10^{-9}$ | 0.02 | 0.003 | $9.90\times10^{-9}$ | -0.02 | 0.003 | $1.07\times10^{-15}$ | $1.23\times10^{-4}$ | 0 |
| rs34811474 | 4 | 25408838 | 4p15.2 | G/A | 0.231 | ANAPC4 | p.Arg465Gln | Primary | 0.02 | 0.003 | $1.10\times10^{-7}$ | 0.02 | 0.003 | $1.20\times10^{-9}$ | -0.02 | 0.002 | $1.13\times10^{-15}$ | $1.23\times10^{-4}$ | 0 |
| rs1291602 | 5 | 130766662 | 5q31.1 | C/T | 0.159 | CTC-432M15.3 | p.Gln1452Arg | Primary | -0.02 | 0.004 | $3.80\times10^{-7}$ | -0.02 | 0.004 | $2.90\times10^{-8}$ | -0.02 | 0.003 | $2.21\times10^{-13}$ | $1.03\times10^{-4}$ | 0 |
| rs351855 | 5 | 176520243 | 5q35.2 | G/A | 0.297 | FGFR4 | p.Gly388Arg | Primary | -0.04 | 0.003 | $4.90\times10^{-35}$ | -0.03 | 0.003 | $2.80\times10^{-24}$ | 0.03 | 0.003 | $3.16\times10^{-55}$ | $4.68\times10^{-4}$ | 0 |
| rs35523808 | 6 | 75834971 | 6q13 | A/T | 0.951 | COL12A1 | p.Glu2160Val | Primary | 0.05 | 0.006 | $1.60\times10^{-13}$ | 0.05 | 0.007 | $7.40\times10^{-12}$ | -0.05 | 0.005 | $9.93\times10^{-24}$ | $2.05\times10^{-4}$ | 0 |
| rs10283100 | 8 | 120596023 | 8q24.12 | G/A | 0.056 | ENPP2 | p.Ser493Pro | Primary | -0.04 | 0.006 | $2.70\times10^{-11}$ | -0.04 | 0.006 | $9.50\times10^{-10}$ | -0.04 | 0.004 | $1.20\times10^{-18}$ | $1.50\times10^{-4}$ | 0 |
| rs41307479 | 9 | 116082647 | 9q32 | C/G | 0.221 | WDR31 | p.Cys256Ser | Primary | 0.01 | 0.003 | $2.00\times10^{-5}$ | 0.02 | 0.004 | $3.70\times10^{-6}$ | 0.02 | 0.002 | $1.30\times10^{-9}$ | $7.04\times10^{-5}$ | 0 |
| rs10761129 | 9 | 94486321 | 9q22.31 | T/C | 0.331 | ROR2 | p.Val819Ile | Primary | -0.01 | 0.003 | $3.70\times10^{-7}$ | -0.02 | 0.003 | $9.60\times10^{-10}$ | 0.02 | 0.002 | $8.87\times10^{-15}$ | $1.16\times10^{-4}$ | 16.2 |
| rs2277339 | 12 | 57146069 | 12q13.3 | T/G | 0.104 | PRIM1 | p.Asp5Ala | Primary | 0.02 | 0.004 | $8.70\times10^{-7}$ | 0.02 | 0.005 | $1.90\times10^{-6}$ | 0.02 | 0.003 | $1.18\times10^{-11}$ | $8.85\times10^{-5}$ | 0 |
| rs12889267 | 14 | 21542766 | 14q11.2 | A/G | 0.167 | ARHGEF40 | p.Lys293Glu | Primary | 0.02 | 0.004 | $6.50\times10^{-6}$ | 0.03 | 0.004 | $4.50\times10^{-11}$ | 0.02 | 0.003 | $5.08\times10^{-14}$ | $1.09\times10^{-4}$ | 70.2 |
| rs117068593 | 14 | 93118229 | 14q32.13 | C/T | 0.19 | RIN3 | p.Arg204Cys | Primary | -0.04 | 0.003 | $8.30\times10^{-34}$ | -0.05 | 0.004 | $6.90\times10^{-39}$ | -0.05 | 0.003 | $4.40\times10^{-71}$ | $6.19\times10^{-4}$ | 50.9 |
| rs35874463 | 15 | 67457698 | 15q22.33 | A/G | 0.058 | SMAD3 | p.Ile170Val | Secondary | 0.02 | 0.006 | $2.64\times10^{-5}$ | 0.03 | 0.006 | $5.73\times10^{-6}$ | 0.03 | 0.004 | $7.11\times10^{-10}$ | $7.29\times10^{-5}$ | 0 |
| rs72755233 | 15 | 100692953 | 15q26.3 | G/A | 0.112 | ADAMTS17 | p.Thr446Ile | Primary | -0.07 | 0.004 | $4.70\times10^{-56}$ | -0.06 | 0.005 | $7.40\times10^{-38}$ | 0.06 | 0.003 | $6.16\times10^{-91}$ | $7.86\times10^{-4}$ | 0 |
| rs3817428 | 15 | 89415247 | 15q26.1 | G/C | 0.265 | ACAN | p.Asp2373Glu | Primary | -0.05 | 0.003 | $1.20\times10^{-64}$ | -0.05 | 0.003 | $9.70\times10^{-45}$ | -0.05 | 0.002 | $2.04\times10^{-104}$ | $9.06\times10^{-4}$ | 27 |
| rs36000545 | 17 | 79093822 | 17q25.3 | A/G | 0.396 | AATK | p.Phe1266Ser | Primary | 0.01 | 0.003 | $5.10\times10^{-6}$ | 0.01 | 0.003 | $2.40\times10^{-6}$ | 0.01 | 0.002 | $2.13\times10^{-10}$ | $7.96\times10^{-5}$ | 0 |
| rs61734651 | 20 | 61451332 | 20q13.33 | C/T | 0.071 | COL9A3 | p.Arg103Trp | Primary | 0.03 | 0.005 | $1.30\times10^{-10}$ | 0.06 | 0.006 | $1.10\times10^{-21}$ | -0.04 | 0.004 | $1.20\times10^{-28}$ | $2.61\times10^{-4}$ | 87.1 |
| rs1291212 | 20 | 62340115 | 20q13.33 | G/C | 0.081 | ZGPAT | p.Ser61Arg | Primary | 0.04 | 0.005 | $2.50\times10^{-18}$ | 0.03 | 0.005 | $6.30\times10^{-10}$ | 0.04 | 0.004 | $7.46\times10^{-25}$ | $2.06\times10^{-4}$ | 49.7 |
| rs1726513 | 20 | 39832628 | 20q12 | T/C | 0.199 | ZHX3 | p.Asn310Ser | Primary | 0.03 | 0.003 | $1.20\times10^{-17}$ | 0.02 | 0.004 | $4.90\times10^{-12}$ | 0.03 | 0.003 | $5.56\times10^{-28}$ | $2.32\times10^{-4}$ | 0 |
| rs2830585 | 21 | 28305212 | 21q21.3 | C/T | 0.16 | ADAMTS5 | p.Arg614His | Primary | -0.02 | 0.004 | $2.60\times10^{-8}$ | -0.03 | 0.004 | $3.40\times10^{-14}$ | 0.02 | 0.003 | $2.80\times10^{-19}$ | $1.56\times10^{-4}$ | 54.1 |

**Notes**: Abbreviations: rsID, based on dbSNP; CHR, chromosome; POS, base positions; BAND, chromosome band; REF, reference allele; ALT, alternative allele; FRQ, frequency of alternative allele; B,
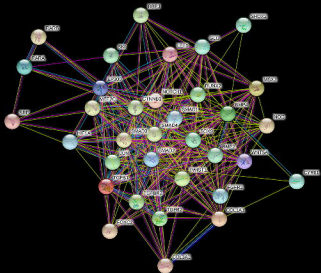
beta coefficient of linear mixed model; SE, standard error of beta; P, P-value; H2, proportion of variance explained by this SNP; I2, I2 statistics of SNP in meta-analysis. Under condition column,
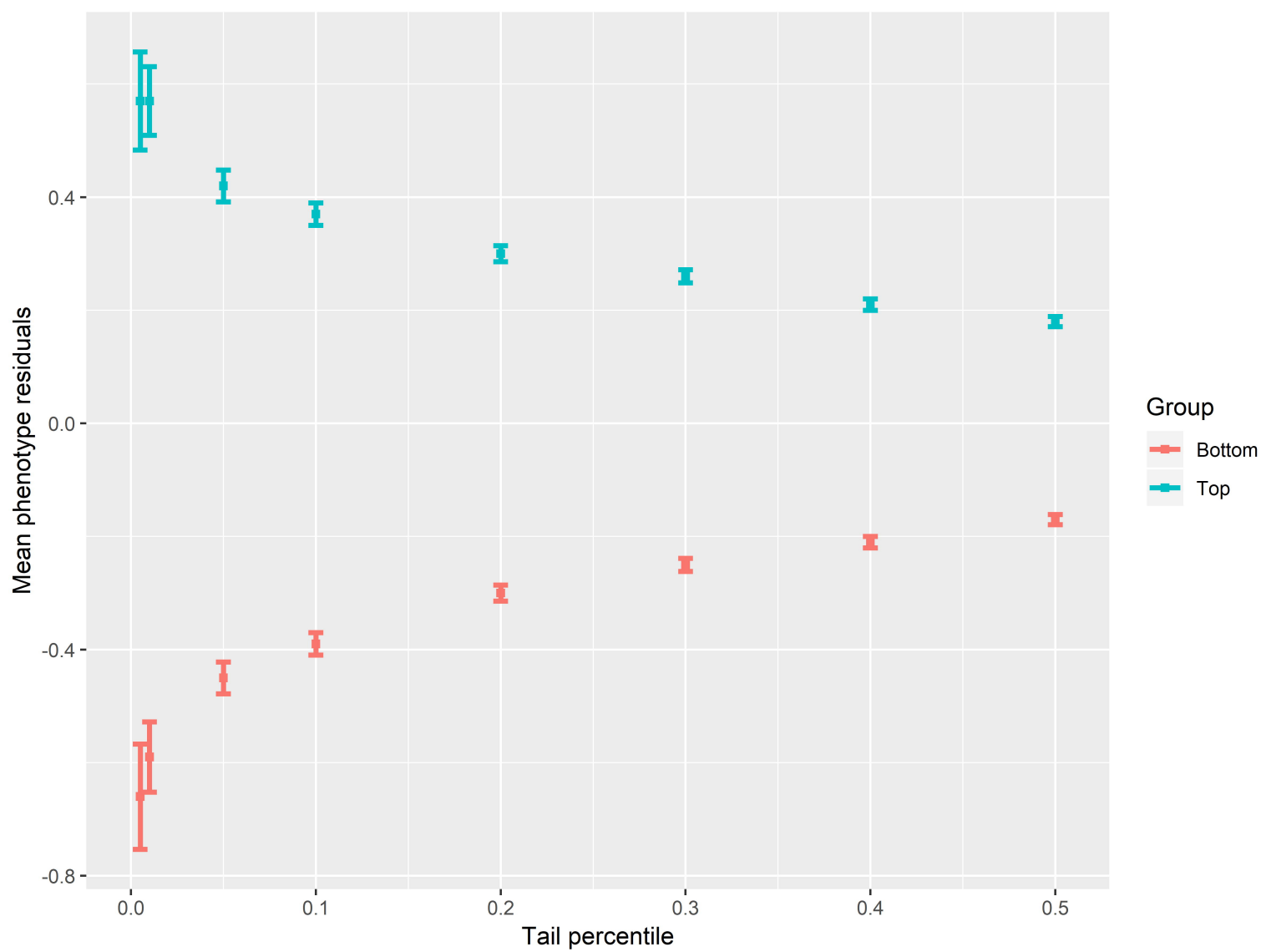
Primary means it is an independent lead SNP identified before conditioning analysis and Secondary means it is a SNP identified in conditional analysis. Genomic coordinate was based on human
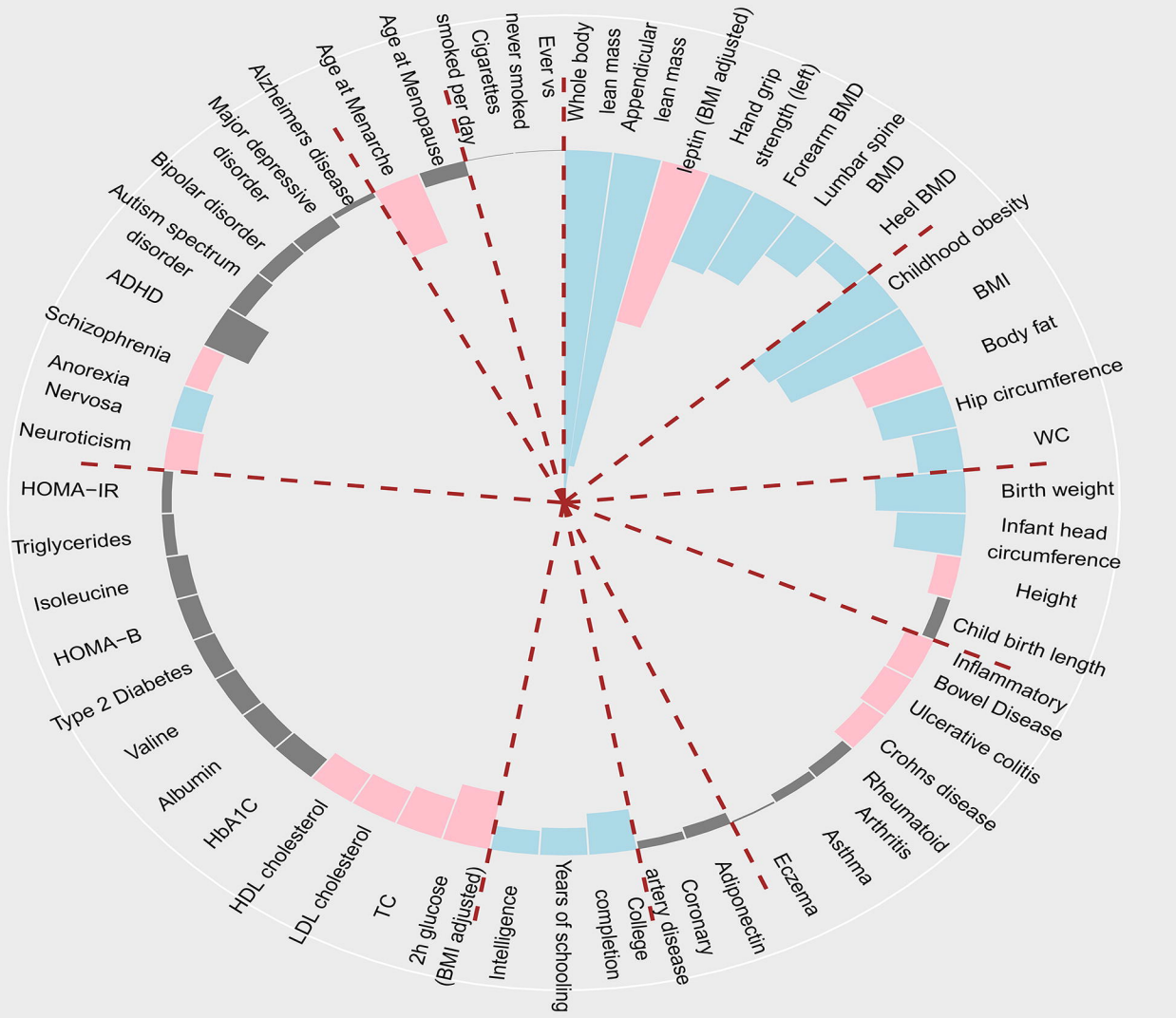
genome assembly build 37 (GRCh 37).

A



B