

1
2
3
4
5
6
7
8

Supplementary Information

**Genome-wide analyses of behavioural traits are subject to bias by
misreports and longitudinal changes**

Xue et al.

9 **Supplementary Notes**

10

11 **Supplementary Note 1: Identifying misreporting "never drinkers" in the UKB**

12 Following Klatksy et al.¹, we attempted to identify "unreliable" self-reported never drinkers using
13 follow-up questionnaires and medical records. The UKB had online follow-up questionnaires in 2017.
14 There were 11 questions related to "alcohol use" in the "mental health" category ($n = 157,365$). We
15 extracted the "frequency of drinking alcohol" (data-field ID: 20414) of 3,627 self-reported never
16 drinkers in the first assessment (2006-2010), but 335 of them (~9.2%) reported that they were not
17 never drinkers in this follow-up assessment (2017). Although these individuals could change drinking
18 status after a few years, it is reasonable to question the reliability of their reported drinking status in
19 the initial assessment. We also extracted the ICD 10 codes (data-field ID: 41202) of 14,488 self-
20 reported never drinkers. People with diagnosed alcohol-related diseases were very likely to have
21 misreported their drinking status. The diseases include E24.4: alcohol-induced pseudo-Cushing's
22 syndrome, F10: mental and behavioural disorders due to use of alcohol, G31.2: degeneration of
23 nervous system due to alcohol, G62.1: alcoholic polyneuropathy, G72.1: alcoholic myopathy, I42.6:
24 alcoholic cardiomyopathy, K29.2: alcoholic gastritis, K70: alcoholic liver disease, K85.2: alcohol-
25 induced acute pancreatitis, K86.0: alcohol-induced chronic pancreatitis, R78.0: finding of alcohol in
26 blood, T51: toxic effect of alcohol, Z50.2: alcohol rehabilitation, and Z72.1: alcohol use. There were
27 77 individuals diagnosed with these diseases; thus, their self-reported drinking status was also likely
28 to be unreliable.

29

30 **Supplementary Note 2: Simulation**

31 To validate our findings, we performed a series of simulations to mimic MLC due to disease
32 ascertainment. There were four simulation scenarios, as shown in **Supplementary Figure 4**. We
33 simulated 20,000 individuals and 100 causal variants affecting a behavioural phenotype (Y) and
34 another set of independent 100 causal variants affecting the liability of a disease (D). Both Y and D
35 were quantitative. The variance explained by the causal variants was 0.6 for both Y and D, *i.e.*, $h_Y^2 =$
36 $h_D^2 = 0.6$. The SNP effects were randomly drawn from $\mathcal{N}(0,1)$. The causal effect (b_{xy}) of Y on D was
37 set to 0.2.

38

39 We mimicked the disease ascertainment by reducing Y to a lower level if the corresponding D value
40 was high. In other words, those individuals with high D values (located in the 10, 20, 30 or 40% upper
41 tail of the distribution) were regarded as disease carriers, and their Y values were deducted by a
42 constant (1-5 standard deviations, *s.d.*). After the ascertainment, we rescaled Y and conducted GWAS
43 for Y and D, and then estimated the correlation of true SNP effects (r_b) between Y and D accounting

44 for errors in estimated SNP effects using a recently developed approach², and the causal effect (b_{xy}) of
45 Y on D using Mendelian Randomization (MR).

46

47 In model I, where Y and D were independent, and the SNPs were associated with Y only, the r_b and
48 b_{xy} estimates were expected to be 0 in the absence of ascertainment, consistent with our simulation
49 results (**Supplementary Figure 5a**). However, the ascertainment generated a negative correlation
50 between Y and D, leading to negative estimates of both r_b and b_{xy} (**Supplementary Figures 5-6**).

51

52 In model II, where Y had a causal effect on D, and the SNPs only had direct effects on Y, the \hat{r}_b only
53 slightly decreased with the increased strength of ascertainment, suggesting that the SNP effect
54 correlation estimate under a causal model was not heavily biased by the ascertainment
55 (**Supplementary Figure 5b**). Even when 10% of the individuals in the upper tail of the distribution of
56 D were reduced by 5 *s.d.* units in Y, the \hat{r}_b only decreased from 1.000 (*s.e.* = 0.003) to 0.929 (*s.e.* =
57 0.003). In the meanwhile, the causal effect estimated from MR analysis increased from 0.200 (*s.e.* =
58 0.002) to 0.390 (*s.e.* = 0.004). Notably, the number of index SNPs decreased as the ascertainment
59 became stronger (**Supplementary Figure 6b**), indicating that the ascertainment could reduce the
60 power to detect causal variants in GWAS.

61

62 In model III, where Y and D were independent, and the SNPs were associated with D only, the
63 ascertainment induced a negative correlation between Y and D (**Supplementary Figure 5c**), and
64 more genome-wide significant SNPs were detected to be associated with Y as the ascertainment
65 strength became larger (**Supplementary Figure 6c**).

66

67 In model IV, where Y had a causal effect on D with 100 SNPs affecting Y and another set of 100
68 SNPs affecting D, the \hat{r}_b gradually changed from positive to negative as the ascertainment became
69 stronger (**Supplementary Figure 5d**). In the MR analysis, when the ascertainment strength was
70 modest, the \hat{b}_{xy} was more robust than the \hat{r}_b (**Supplementary Figure 6d**).

71

72 The simulations above are all for longitudinal change; however, we can also simulate underreporting
73 using a similar procedure, *i.e.*, assigning a lower value to Y for individuals with large D. The only
74 difference between underreporting and longitudinal change in the simulation was the proportion of
75 individuals affected. We set the proportion of underreporting individuals from 2% to 8% of the upper
76 tail of the distribution of D based on that observed in the UKB. Our simulation results showed that the
77 effects of ascertainment bias from underreporting were smaller than those from longitudinal change
78 (**Supplementary Figures 7-8**).

79

80 **Supplementary Note 3: The relationship between AC and cardiovascular disease (CVD)**

81 To investigate the observed relationship between AC and CVD, we first performed logistic regression
82 analyses of cardiovascular disease on different AC intake levels as suggested in Wood et al.³. The
83 relationship was J-shaped where moderate drinking showed a lower disease risk and heavy drinking
84 showed a higher disease risk than that in the reference group ($0 < AC \leq 25$ grams/week)
85 (**Supplementary Figure 15a**). We performed the MLC corrections by excluding underreporting
86 individuals and individuals who reduced drinking because of illness or doctor's advice, and fitted
87 longitudinal change as the covariate in the logistic model. The J-shaped relationship remained but the
88 risk threshold (the point at which odds ratio, i.e. OR, of CVD becomes larger than 1 as AC increases)
89 shrank towards the left (**Supplementary Figure 15b**). However, when we removed only the
90 individuals who had reduced their drinking amount in the reference group, the relationship between
91 AC and CVD became monotonically increasing (**Supplementary Figure 15c**), suggesting an
92 enrichment of disease ascertained individuals in the reference group as demonstrated in
93 **Supplementary Figure 14**.

94
95 We performed a simulation to verify whether we would expect a J curve between the genetic
96 predictor of X and the raw phenotype of Y if the true relationship is a J curve. We first simulated X
97 and Y in 50,000 unrelated individuals. There were 160 causal variants for X (total $h^2 = 0.3$), 100
98 causal variants for Y (total $h^2 = 0.3$), and 40 pleiotropic variants for X and Y (total $h^2 = 0.1$ for both X
99 and Y). The SNP effects were randomly drawn from $\mathcal{N}(0,1)$. We simulated a J-shaped causal effect
100 of X on Y (formula: $Y \sim X^2 + X$). Individuals with the top 20% Y values were regarded as disease
101 carriers. We divided X into ten deciles and labelled the lowest 10% as the reference group and
102 estimated the effect of X and Y using the rest nine deciles against the reference group in a logistic
103 regression. We plotted the mean of X in each decile against the OR in each comparison as we did in
104 the real data and observed a J-shaped curve (**Supplementary Figure 16a**). Then we used genome-
105 wide significant variants ($P < 5 \times 10^{-8}$) for X to generate a genetic predictor of X, and then
106 estimated the effect of the genetic predictor of X on the Y in different quantiles of X. The simulation
107 was replicated 50 times and the relationship was still a J curve (**Supplementary Figure 16b**), which
108 suggests we would expect a J curve between a genetic predictor of X and Y if the true relationship is a
109 J curve.

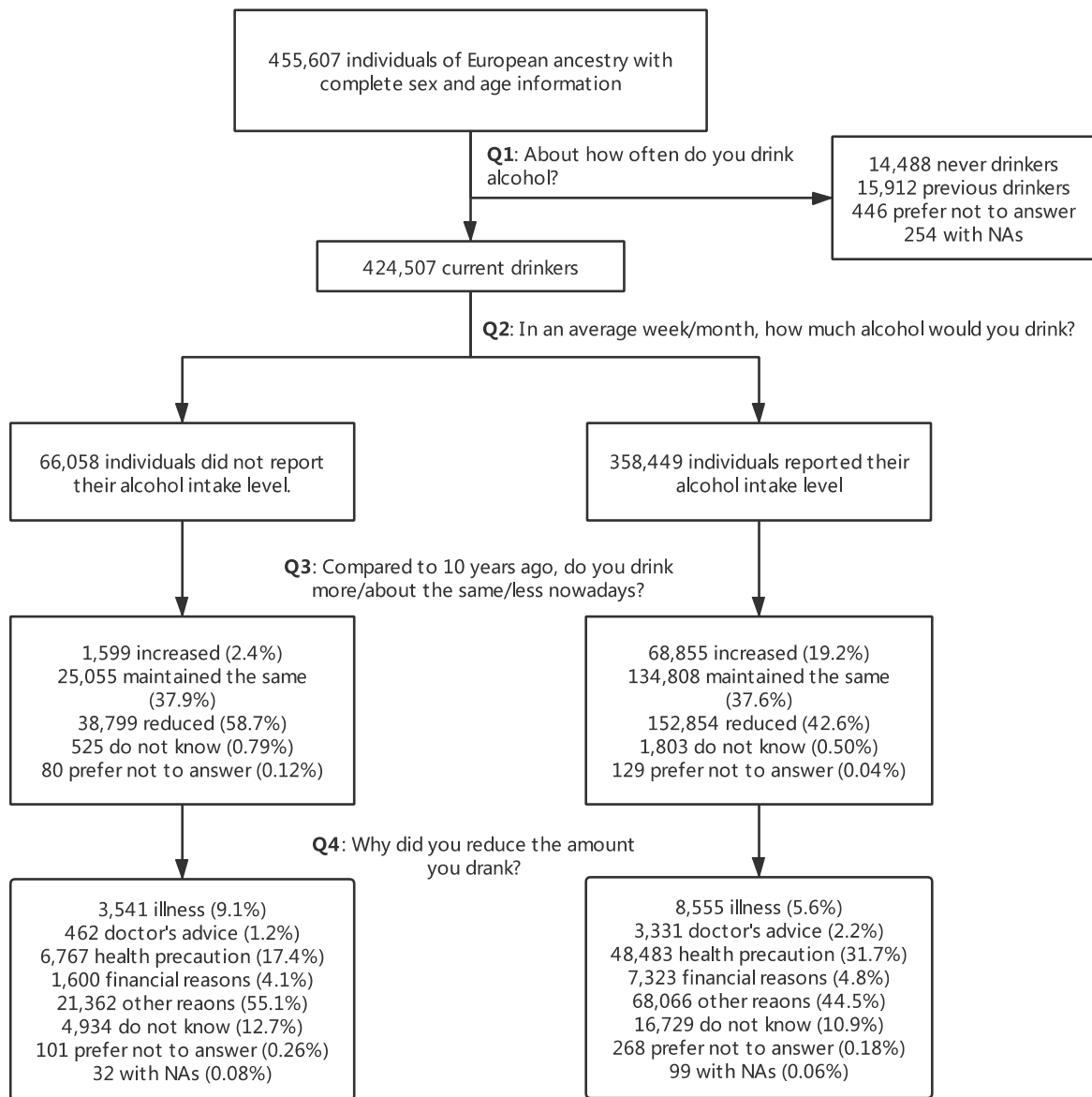
110
111 **Supplementary Note 4: MLC corrections for smoking intensity**

112 According to the self-reported records in the UKB (data-field ID: 20116), there were ~245,000 never
113 smokers, ~162,000 previous smokers and ~47,000 current smokers. The cigarettes per day (CPD) data
114 were collected among the current smokers who used manufactured cigarettes or hand-rolled cigarettes
115 (data-field ID: 3456). According to the self-reported longitudinal change information from 32,801
116 current cigarette smokers (data-field ID: 3506), 5,559 individuals increased their smoking intensity,

117 13,235 maintained the same intensity and 13,941 reduced their smoking intensity compared to 10
118 years ago (**Supplementary Data 12a**). We performed the MLC corrections for CPD by 1) partitioned
119 the current smokers into three longitudinal change groups, 2) excluded 3,061 individuals who chose
120 illness or doctor's advice as the reason for reducing smoking (data-field ID: 6158), 3) performed
121 GWAS in each group with standardised CPD and meta-analysed GWAS summary statistics from the
122 three groups. We compared the GWAS results for CPD with or without the MLC corrections
123 (**Methods**) and found that the estimate of genetic correlation between CPD before and after the MLC
124 corrections was not significantly different from 1 ($\hat{r}_g = 0.985$, *s. e.* = 0.015). Additionally, we did
125 not observe any large differences in the \hat{r}_g of CPD with diseases before and after the MLC corrections
126 (**Supplementary Data 13** and **Supplementary Figure 17**).

127
128
129

130 **Supplementary Figures**



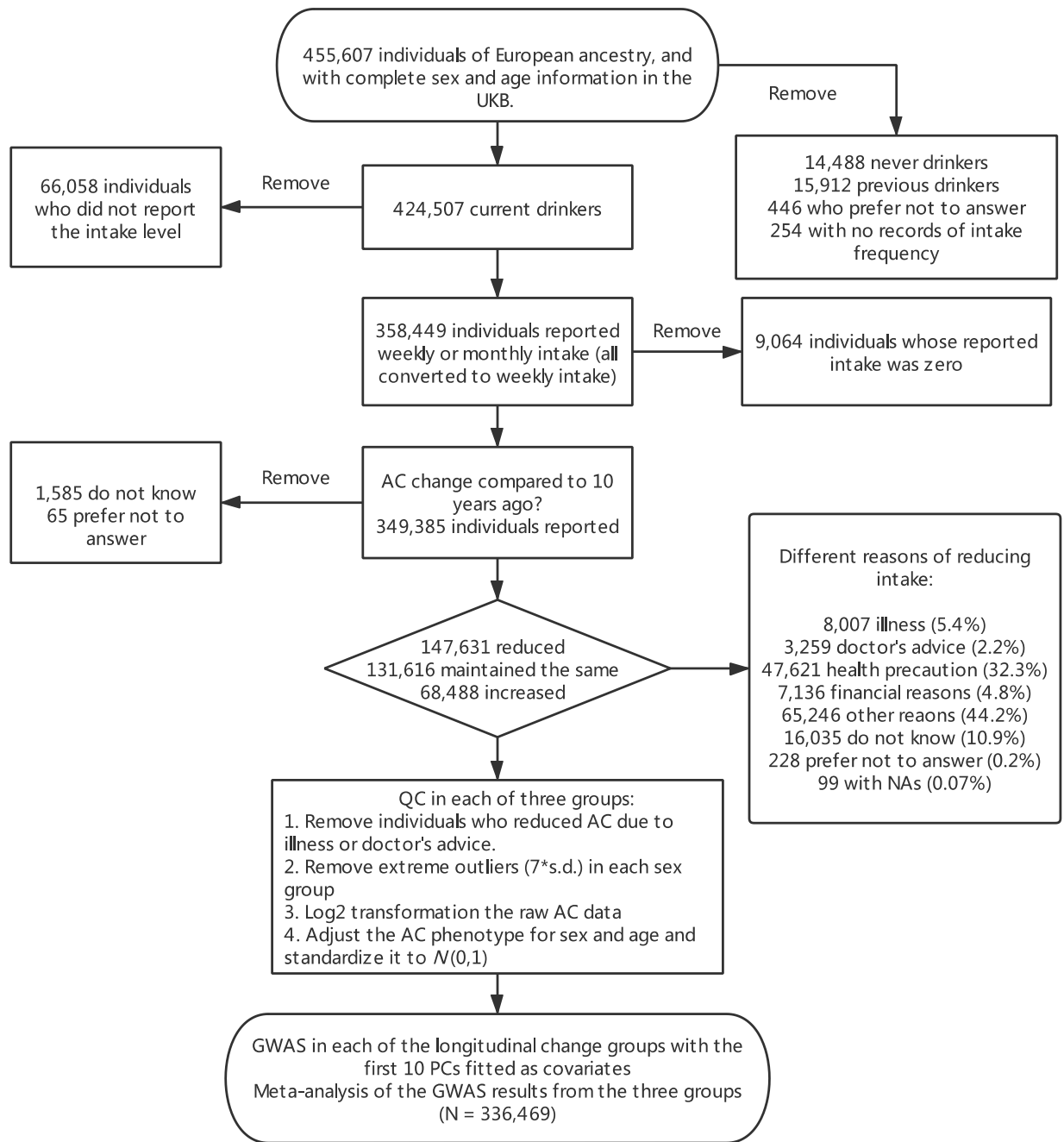
131

132 **Supplementary Figure 1. Flow chart of the alcohol-related questionnaire in the UK Biobank.**

133 The full questionnaire can be found in pages 35-38 at

134 <http://biobank.ndph.ox.ac.uk/showcase/showcase/docs/TouchscreenQuestionsMainFinal.pdf>.

135

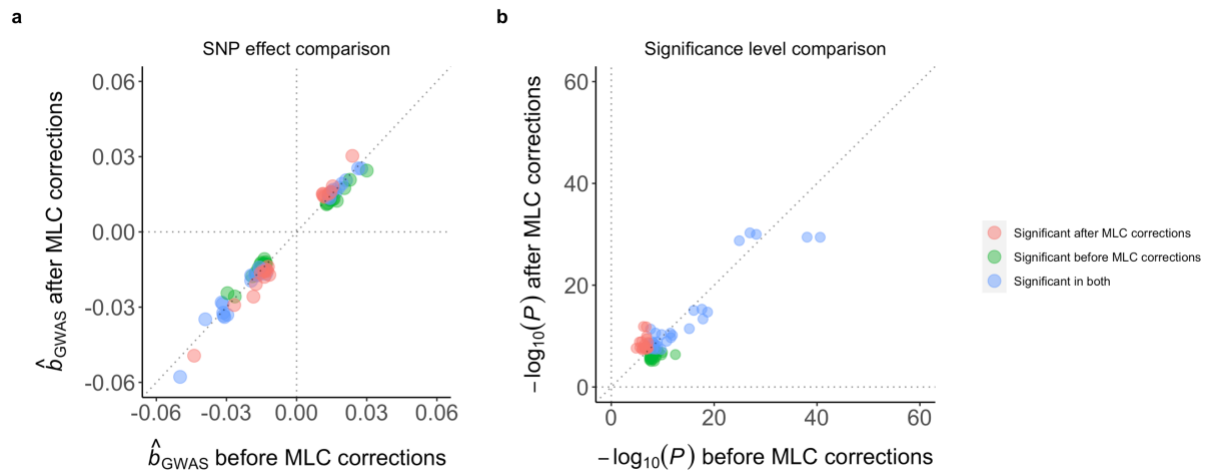


136

137 **Supplementary Figure 2. Flow chart of the MLC corrections for alcohol consumption.** UKB: UK

138 Biobank. AC: alcohol consumption. QC: quality control. PC: principal component.

139



140

141 **Supplementary Figure 3. Comparison between alcohol consumption GWAS results before and**

142 **after the MLC corrections. (a):** Effects of the AC-associated SNPs before and after the MLC

143 corrections. The red dots denote the SNPs that were not significantly associated with AC but became

144 significant ($P < 5 \times 10^{-8}$) after the MLC corrections. The green dots denote the SNPs that were

145 significant but became non-significant the MLC corrections. The blue dots indicate the SNPs that

146 were significant in both. (b): The $-\log_{10} P$ -values of the AC-associated SNPs before and after the

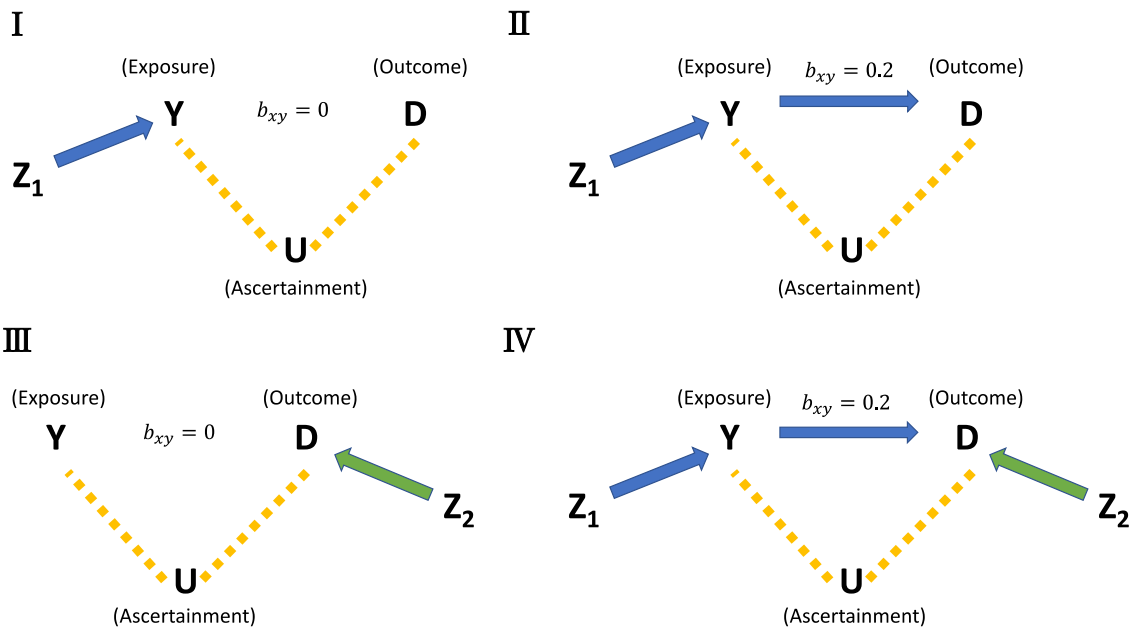
147 MLC corrections. The top SNP rs1229984 at the *ADH1B* locus is omitted due to its large effect size;

148 the effect of the T allele was -0.24 ($P = 4.10 \times 10^{-214}$) and -0.23 ($P = 1.04 \times 10^{-167}$),

149 respectively, before and after the MLC corrections. The P -value indicates the GWAS significant level

150 of each SNP with AC from BOLT-LMM analysis (two-sided χ^2 test).

151

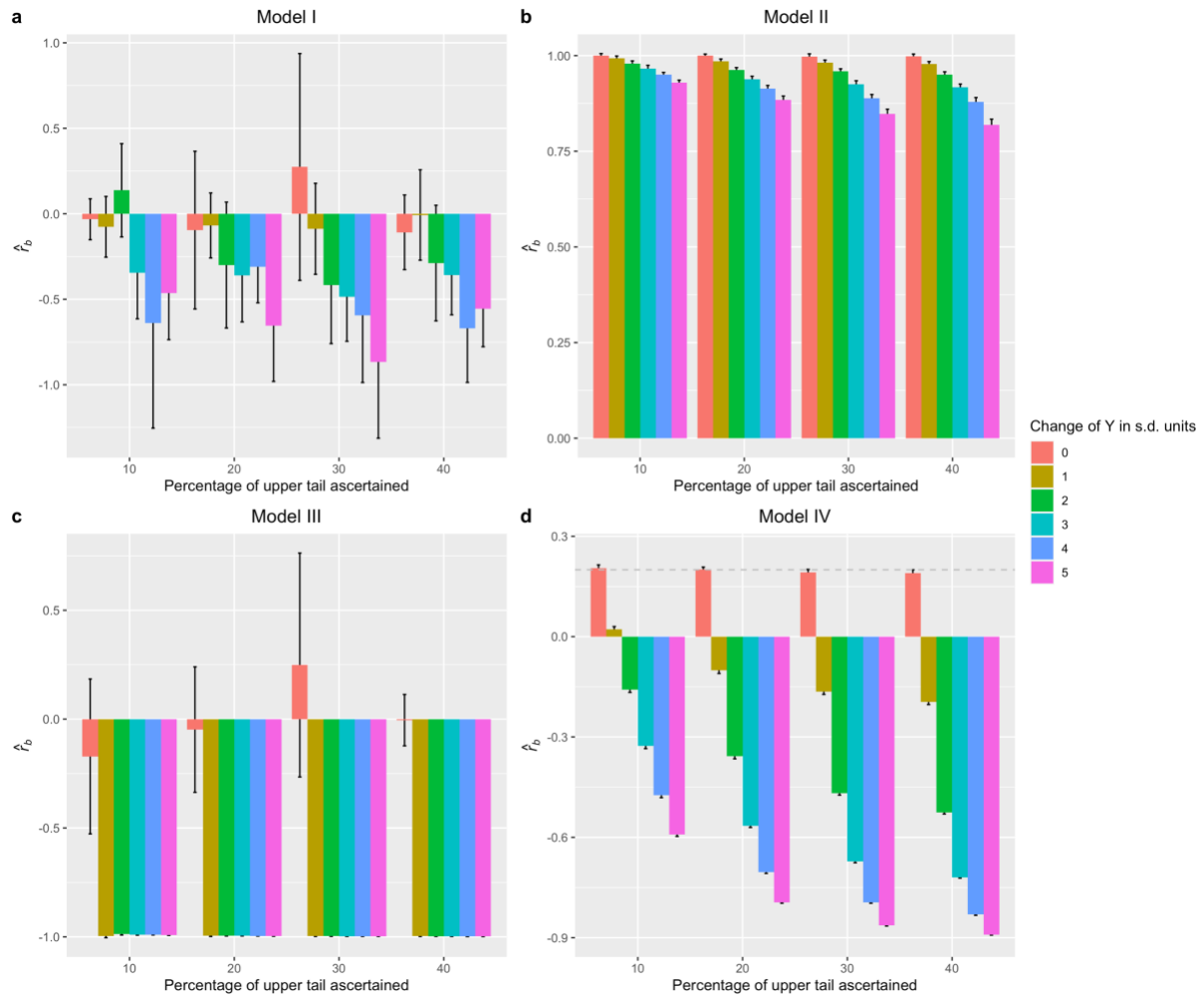


153

154 **Supplementary Figure 4. Four models used in the simulations to mimic disease ascertainment.**

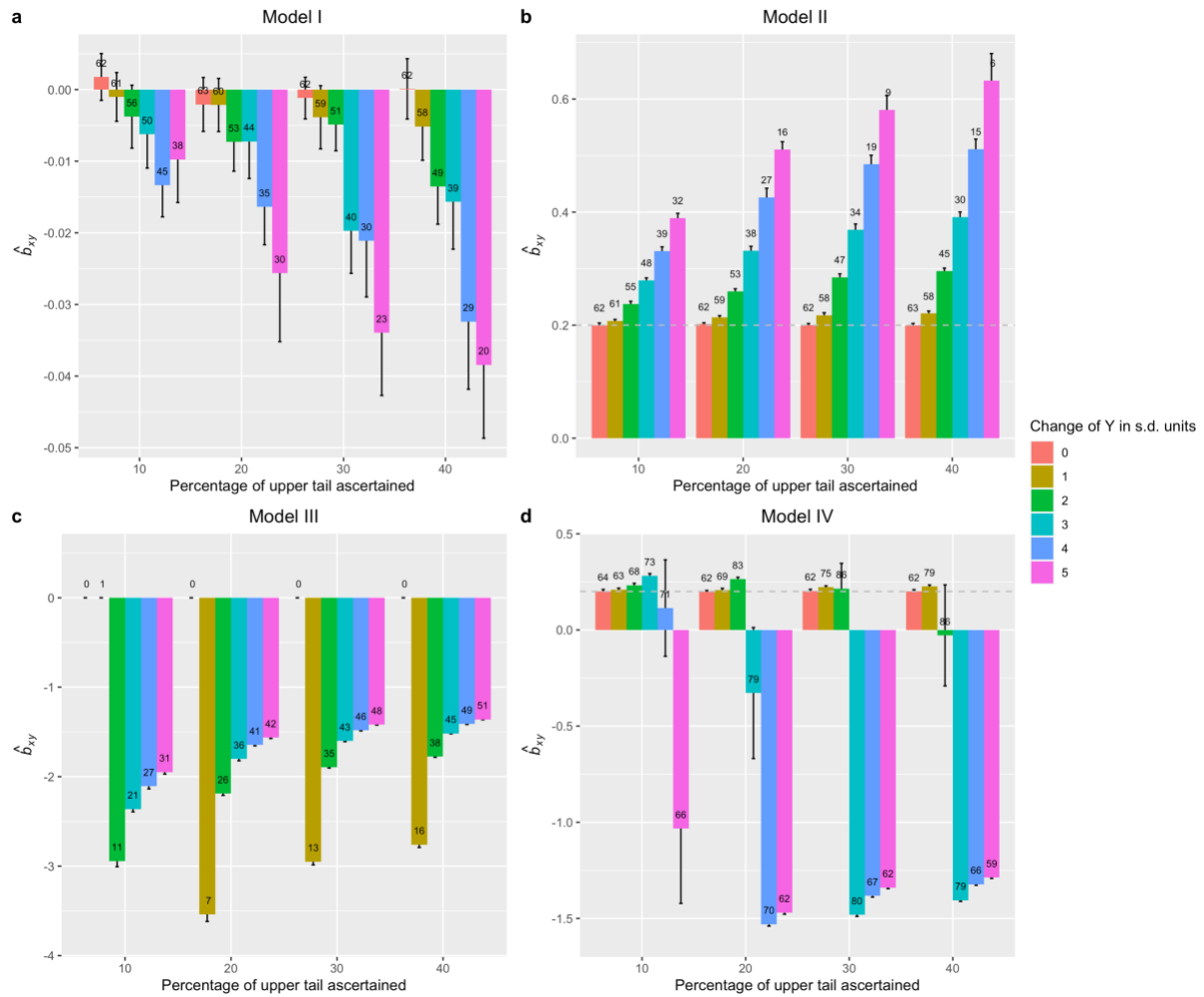
155 Y is a behavioural phenotype, D is the liability of a disease, Z_1 is a set of causal variants only for Y ,
 156 and Z_2 is a set of causal variants only for D . The yellow dashed line indicates the association between
 157 Y and D induced by the change of Y conditioning on D via ascertainment (U). Model I: Y and D are
 158 independent, and 100 SNPs are associated Y . Model II: Y had a causal effect on D , and 100 SNPs are
 159 associated with Y (and D mediated through Y). Model III: Y and D are independent, and 100 SNPs
 160 are associated with D . Model IV: Y had a causal effect on D , 100 SNPs are associated with Y (and D
 161 mediated through Y), and another set of 100 SNPs are associated with D directly.

162



163
 164
 165
 166
 167
 168
 169
 170
 171
 172

Supplementary Figure 5. Quantifying bias in the estimated SNP effect correlation due to longitudinal change by simulation. The four models are defined in **Supplementary Figure 4**. The x-axis indicates the percentage of ascertained individuals. The total sample size used in the simulation $n = 20,000$. The y-axis indicates the r_b estimates. The r_b is defined as Pearson's correlation between the effects of the genetic variants on Y and those on D accounting for errors in the estimated variant effects. The error bars indicate the 95% confidence interval of the r_b estimate. The colour of the bar indicates the strength of ascertainment (*i.e.*, the change of the phenotype Y in *s.d.* units). Change in *s.d.* = 0 means no ascertainment. The grey dashed line indicated $r_b = 0.2$.



173

174

Supplementary Figure 6. Quantifying bias in the estimated causal effect due to longitudinal

175

change by simulation. The four models are defined in **Supplementary Figure 4.** The x-axis

176

indicates the percentage of ascertained individuals. The total sample size used in the simulation $n =$

177

20,000. The y-axis indicates the causal effect estimates, \hat{b}_{xy} . The error bars indicate the 95%

178

confidence interval of the \hat{b}_{xy} . The colour of the bar indicates the strength of ascertainment (*i.e.*, the

179

change in *s.d.* = 0 means no ascertainment. The number

180

labelled on the bar indicates the number of genome-wide significant SNPs of Y. Some of the bars are

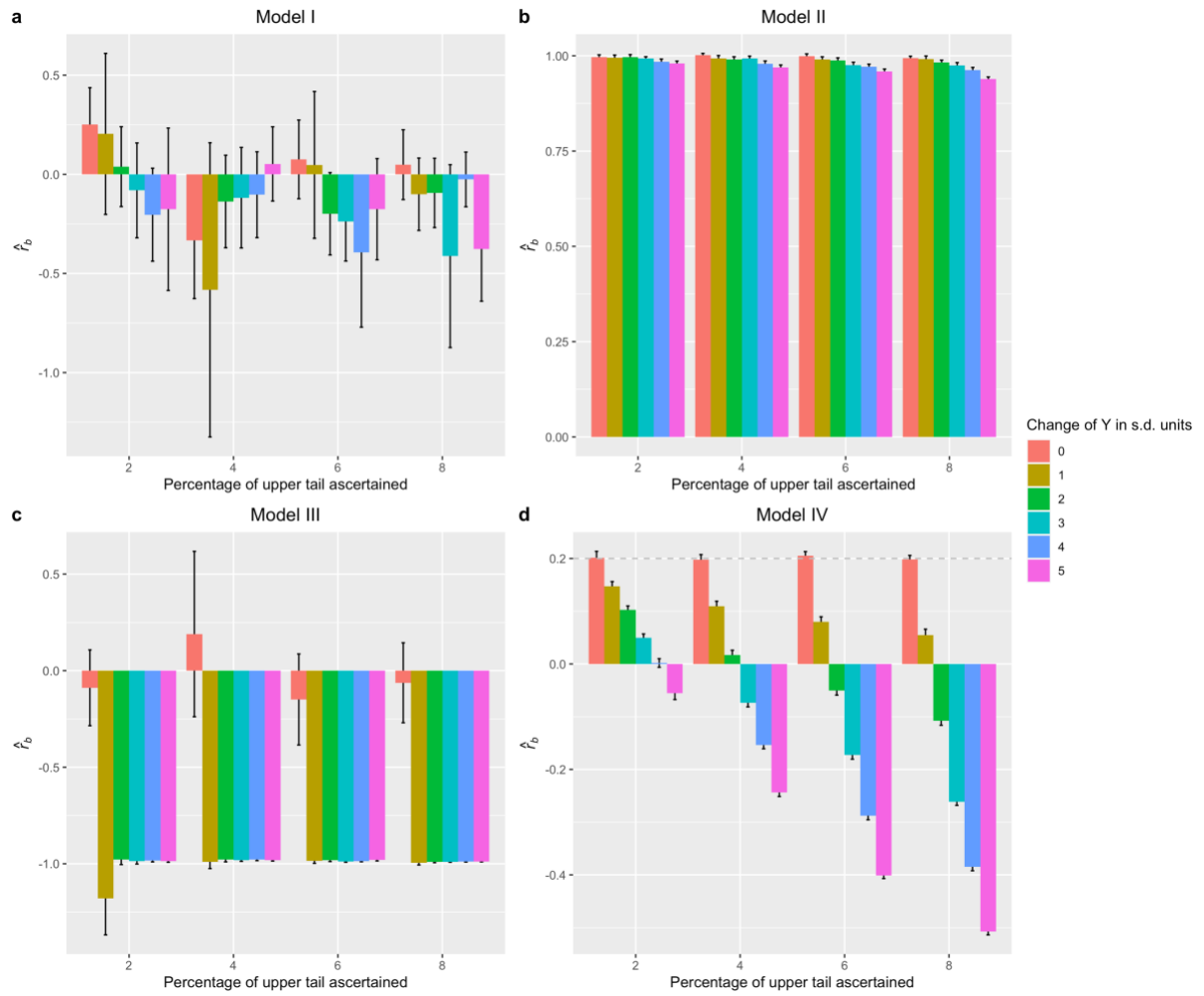
181

missing in panel C because there were not enough instrumental SNPs to perform the GSMR analysis.

182

The grey dashed line indicated $b_{xy} = 0.2$.

183



184

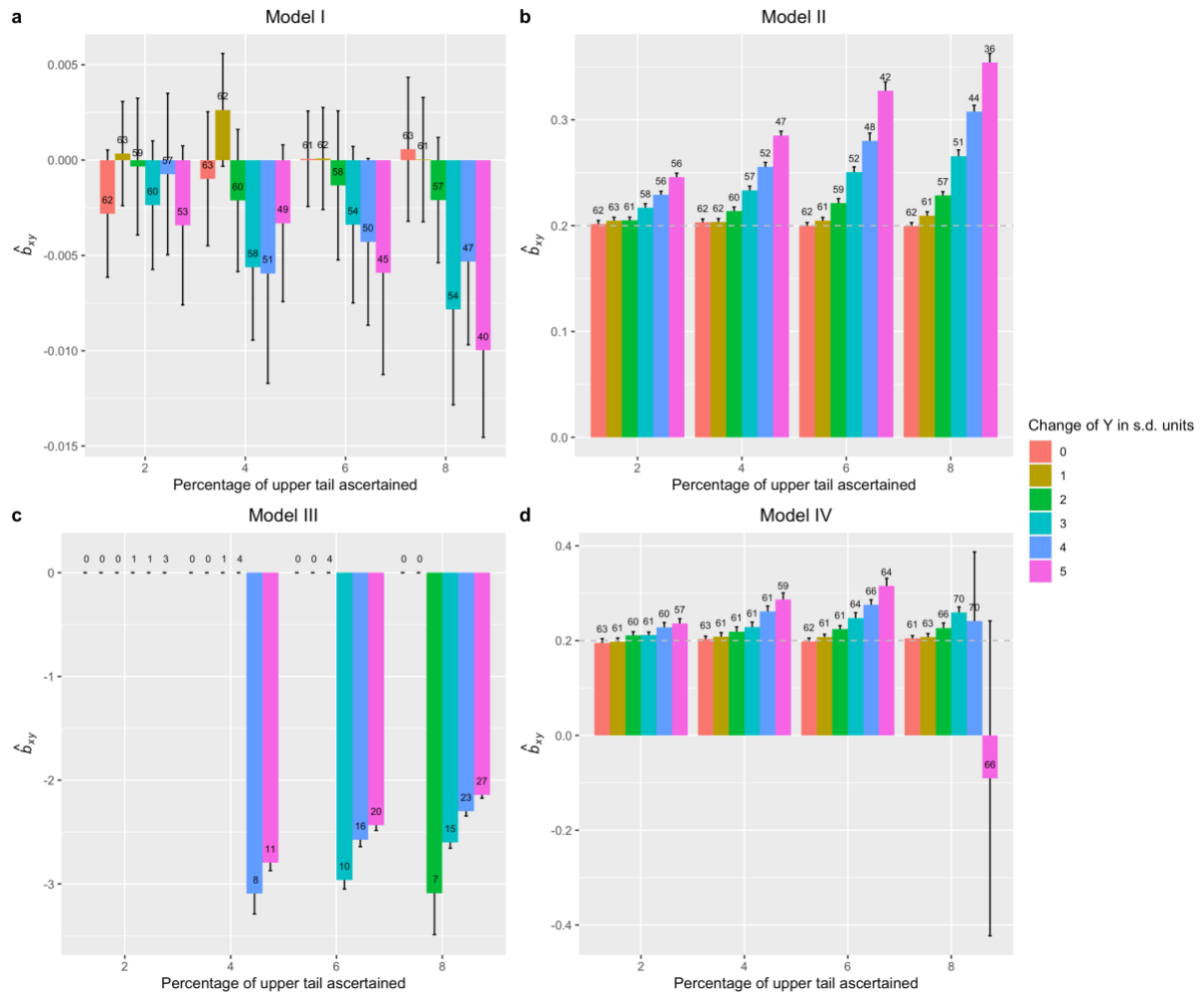
185 **Supplementary Figure 7. Quantifying bias in the estimated SNP effect correlation due to**

186 **misreporting by simulation.** The total sample size used in the simulation $n = 20,000$. The error bars

187 indicate the 95% confidence interval of the r_b estimates. All the labels and colour code are the same as

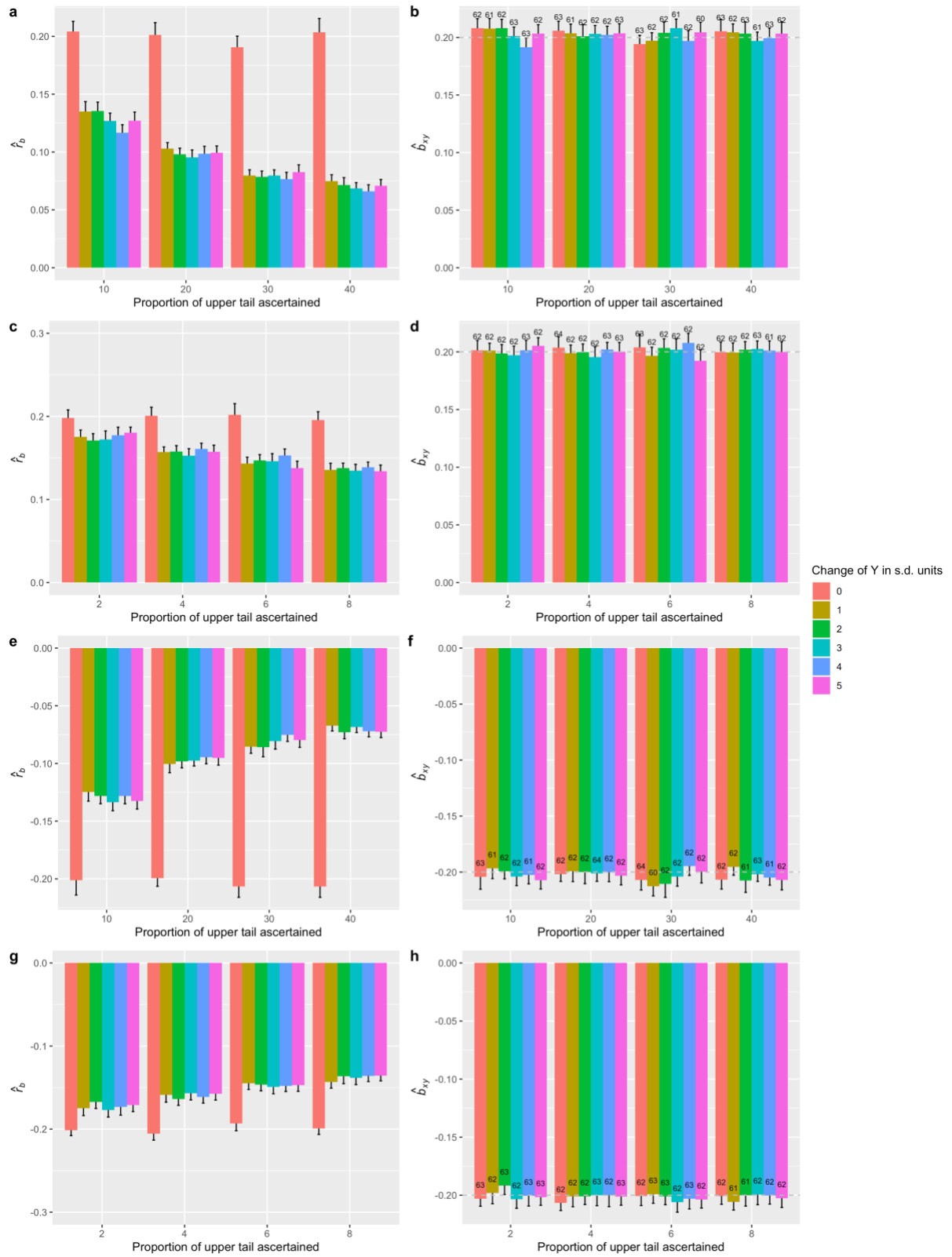
188 those in **Supplementary Figure 5.**

189

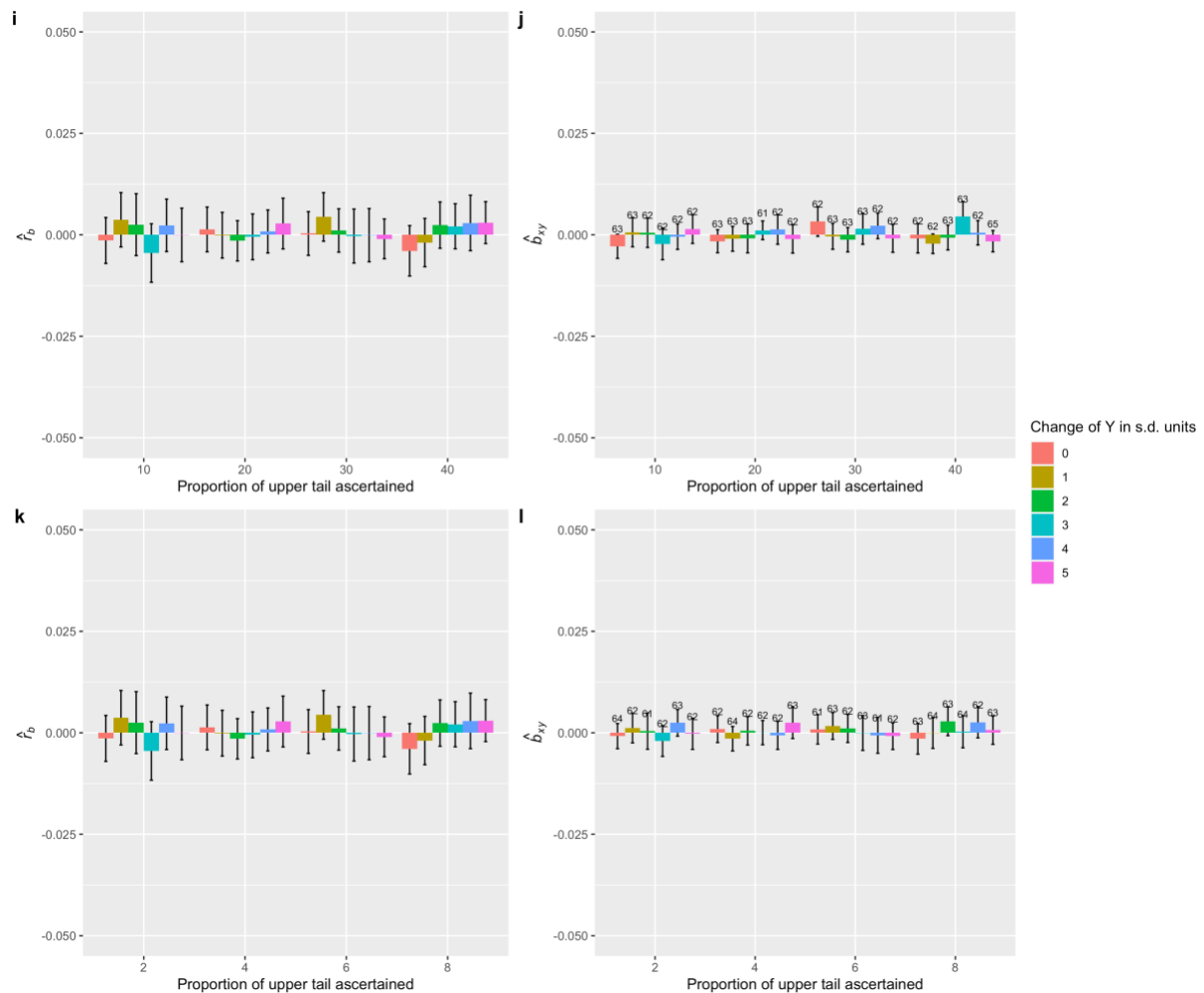


190

191 **Supplementary Figure 8. Quantifying bias in the estimated causal effect due to misreporting by**
 192 **simulation.** The total sample size used in the simulation $n = 20,000$. All the labels and colour code
 193 are the same as those in **Supplementary Figure 6**. The error bars indicate the 95% confidence
 194 interval of the $\hat{\delta}_{xy}$. Change in *s.d.* = 0 means no ascertainment. The number labelled on the bar
 195 indicates the number of genome-wide significant SNPs of Y. Some of the bars are missing in panel C
 196 because there were not enough instrumental SNPs to perform the GSMR analysis.



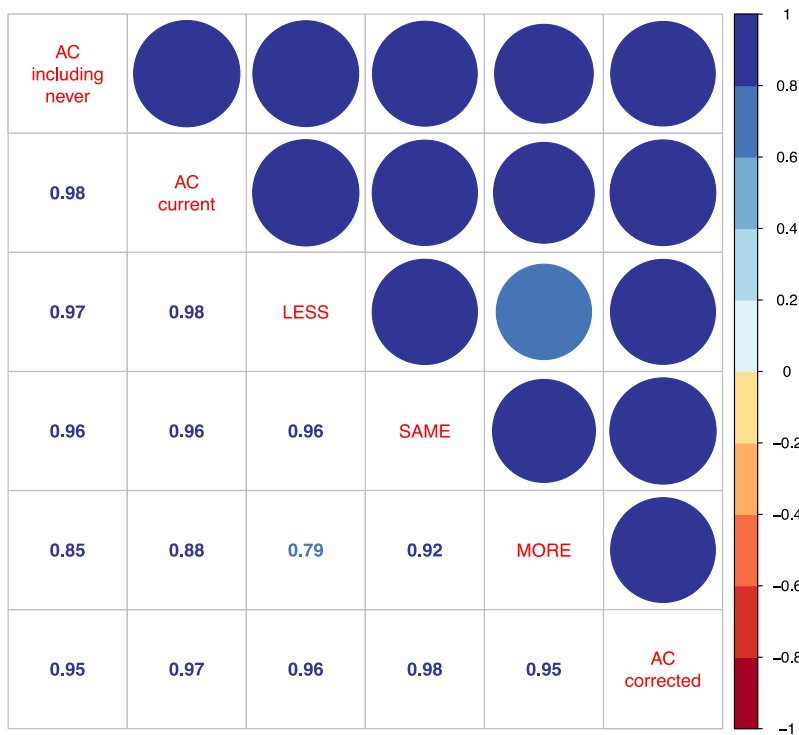
197
198



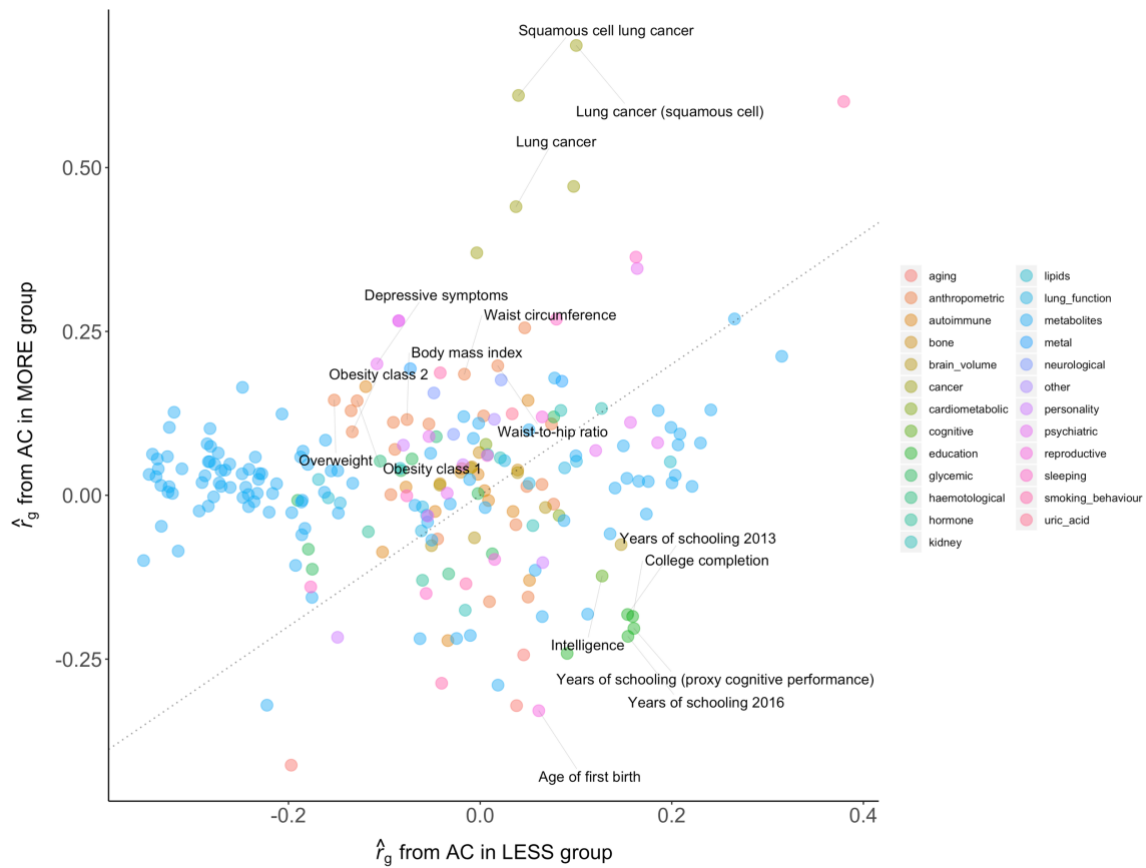
199

200 **Supplementary Figure 9. Estimates of SNP effect correlation and causal effects in simulations**

201 **after the MLC corrections.** The total sample size used in the simulation $n = 20,000$. All the labels
 202 and colour code are the same as those in **Supplementary Figures 5 and 6.** Only the data simulated
 203 based on Model IV were analysed here. Panels (a) and (b) show the r_b and b_{xy} estimates after the MLC
 204 corrections in the presence of longitudinal change, respectively. Panels (c) and (d) show the r_b and b_{xy}
 205 estimates after the MLC corrections in the presence of underreporting, respectively. The error bars
 206 indicate the 95% confidence interval of the r_b or b_{xy} estimates. Panels (e) to (h) are based on the same
 207 simulation setting as those for panels A to D except for that b_{xy} is set to -0.2 . Panels (i) to (l) are
 208 based on the same simulation settings as those for panels (a) to (d) except for that b_{xy} is set to 0 .

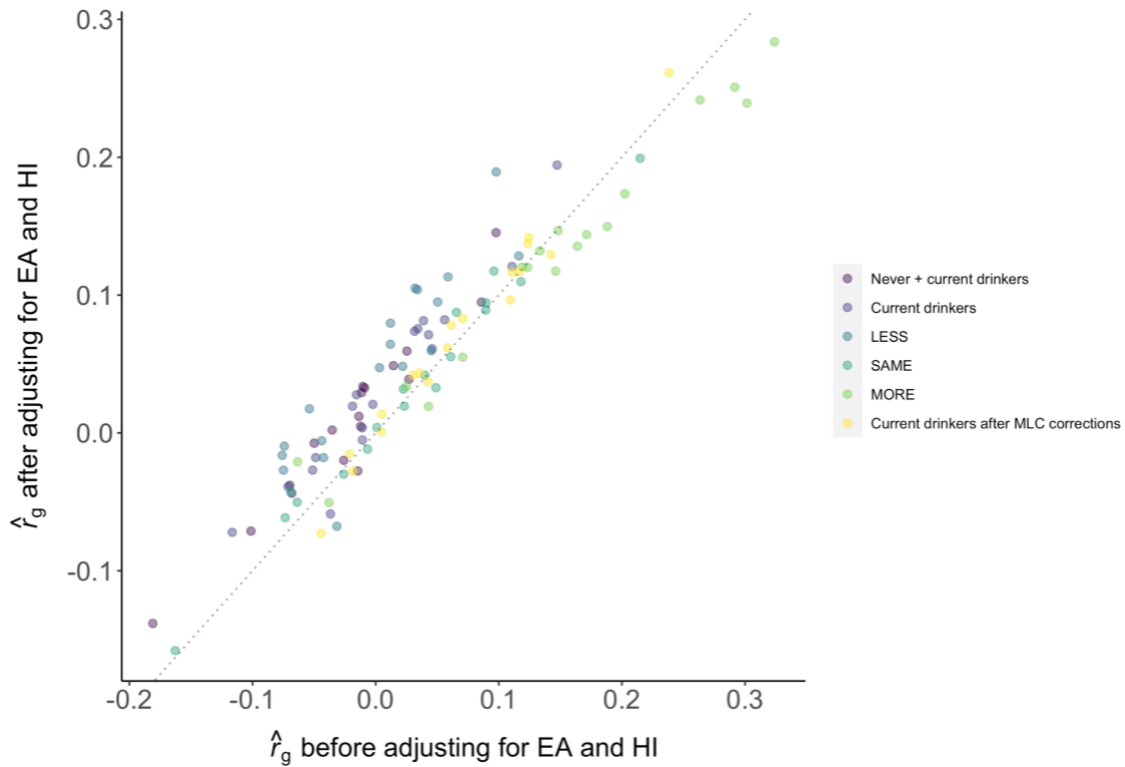


209 **Supplementary Figure 10. Estimates of genetic correlation between different AC groups.** The
 210 value in each cell below the diagonal denotes the r_g estimate from a bivariate LDSC analysis. The
 211 circle in each cell above the diagonal shows the r_g estimate visually: larger circle size and darker color
 212 indicate higher r_g estimate. "AC including never" represents alcohol consumption in current and never
 213 drinkers. "AC current" represents alcohol consumption in current drinkers. LESS, SAME, and MORE
 214 represent current drinkers whose AC levels were reduced, maintained the same, and increased,
 215 respectively, compared to 10 years ago. "AC corrected" represents alcohol consumption in current
 216 drinkers after the MLC corrections.



217
 218
 219
 220
 221
 222
 223

Supplementary Figure 11. Estimates of genetic correlation between AC and 234 traits in LD Hub. The x-axis indicates the r_g estimates using AC from the LESS group, and the y-axis indicates the r_g estimates using AC from the MORE group. The traits with large differences in r_g estimate between the LESS and MORE groups are annotated. The colours of the dots indicate the trait categories defined as defined in LD Hub.

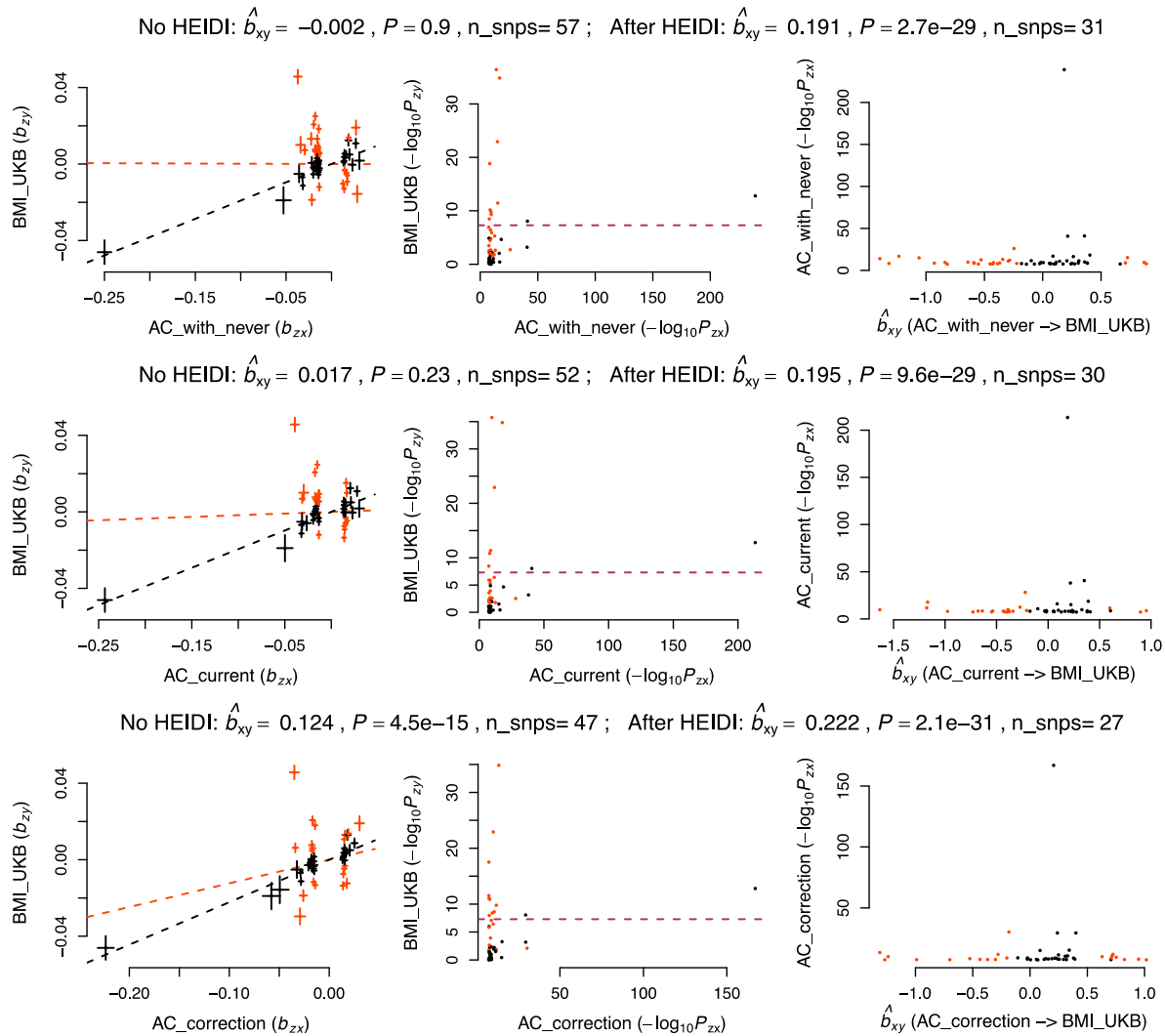


224

225 **Supplementary Figure 12. Comparison of the estimates of genetic correlation between AC and**
 226 **18 common diseases before and after adjusting for socio-economic status (SES).** The x-axis and
 227 y-axis indicate the r_g estimates between AC and common diseases before and after adjusting two SES
 228 traits, educational attainment (EA) and household income (HI). The color of the circle indicates six
 229 different scenarios of AC GWAS. LESS, SAME, and MORE represent current drinkers whose AC
 230 levels were reduced, maintained the same, or increased, respectively, compared to 10 years ago. The
 231 r_g estimates are largely consistent before and after adjusting for EA and HI (Pearson's correlation $r =$
 232 0.951). The Pearson's correlations in the subgroups are 0.967, 0.966, 0.895, 0.991, 0.987, 0.988,
 233 respectively, from the top to the bottom as shown in the legend.

234

235



236

237 **Supplementary Figure 13. GSMR diagnostic analysis of the causal association between AC and**

238 **BMI in the UKB.** The genetic instruments, which were detected by the HEIDI-outlier test as

239 pleiotropic outliers, are highlighted in red. The three panels on the left show the estimated effects of

240 the genetic instruments (index SNPs) of AC (x-axis) against those for BMI (y-axis). The error bars

241 indicate the standard errors of the SNP effect estimates. The slope of the red and black dashed line

242 indicates \hat{b}_{xy} (GSMR estimate of the causal effect of AC on BMI) before and after the HEIDI-outlier

243 filtering, respectively. The panels in the middle shows a plot of $-\log_{10}(P_{zx}$ or $P_{zy})$ for the effect of an

244 index SNPs on the exposure (x-axis) against that for the outcome (y-axis). The panels on the right

245 show the \hat{b}_{xy} estimated using each index SNP (x-axis) against $-\log_{10}(P_{zx})$ for the SNP effect on the

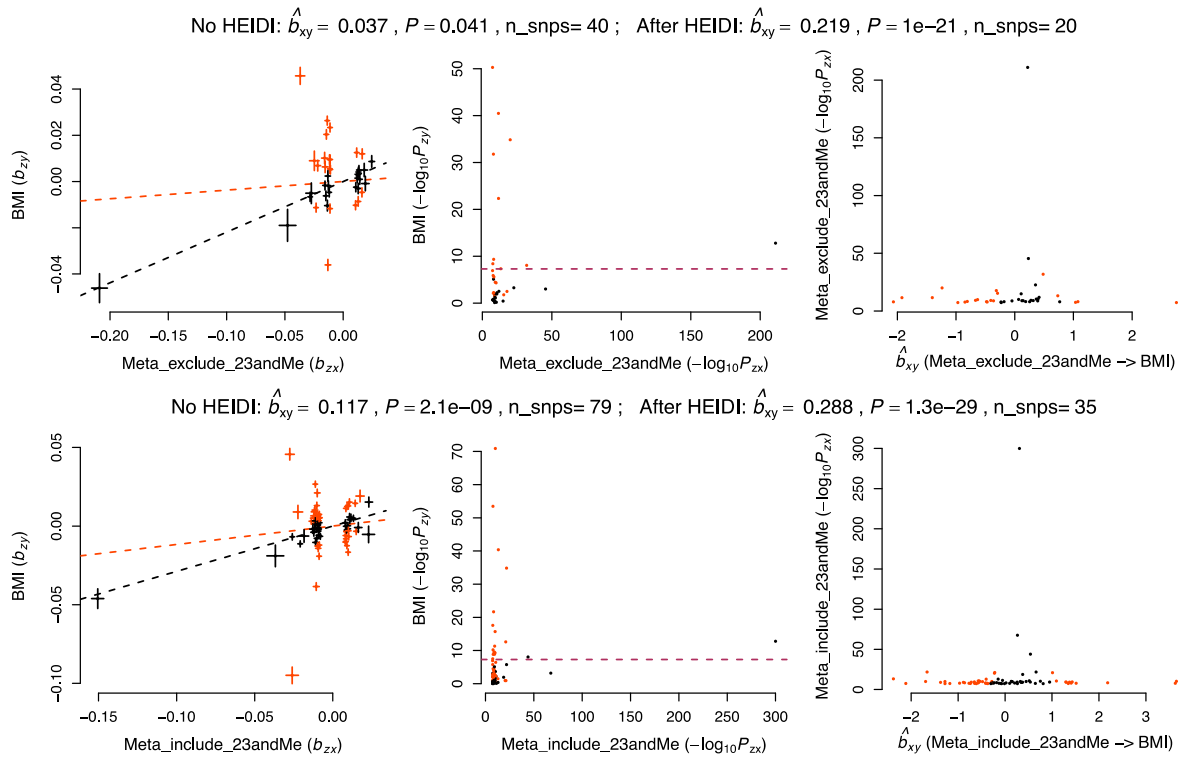
246 exposure (y-axis). "AC_with_never": AC of current and never drinkers; "AC_current": AC of current

247 drinkers; "AC_correction": AC after the MLC corrections. The b_{xy} is estimated based on a two-step

248 least squares (2SLS) approach, and P -value indicates the significant level of the \hat{b}_{xy} in GSMR

249 analysis (two-sided test).

250



251

252

Supplementary Figure 14. GSMR diagnostic analysis of the causal association between AC and

253

BMI using the UKB and GSCAN data. The genetic instruments, which were detected by the

254

HEIDI-outlier test as pleiotropic outliers, are highlighted in red. The two panels on the left show the

255

estimated effects of the genetic instruments (index SNPs) of AC (x-axis) against those for BMI (y-

256

axis). The error bars indicate the standard errors of the SNP effect estimates. The slope of the red and

257

black dashed line indicates \hat{b}_{xy} (GSMR estimate of the causal effect of AC on BMI) before and after

258

the HEIDI-outlier filtering, respectively. The panels in the middle shows a plot of $-\log_{10}(P_{zx}$ or $P_{zy})$ for

259

the effect of index SNPs on the exposure (x-axis) against that for the outcome (y-axis). The panels on

260

the right show the \hat{b}_{xy} estimated using each index SNP (x-axis) against $-\log_{10}(P_{zx})$ for the SNP effect

261

on the exposure (y-axis). "Meta_exclude_23andMe" and "Meta_include_23andMe" represent the

262

GSCAN data⁴ of AC excluding and including 23andMe cohort, respectively. The b_{xy} is estimated

263

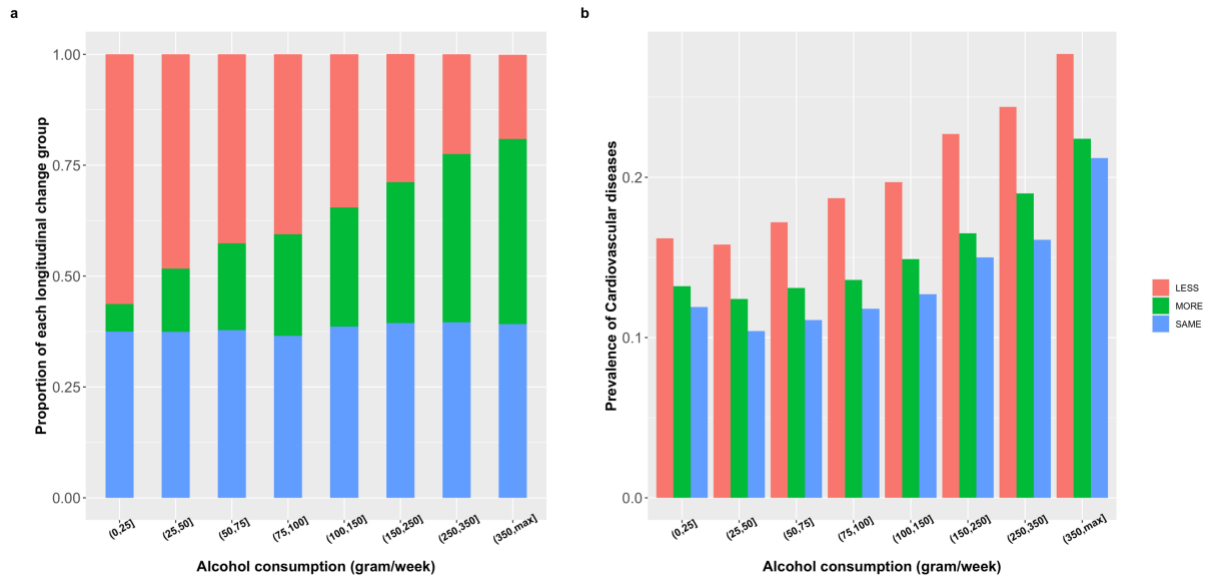
based on a two-step least squares (2SLS) approach, and P -value indicates the significant level of the

264

\hat{b}_{xy} in GSMR analysis (two-sided test).

265

266



267

268

269

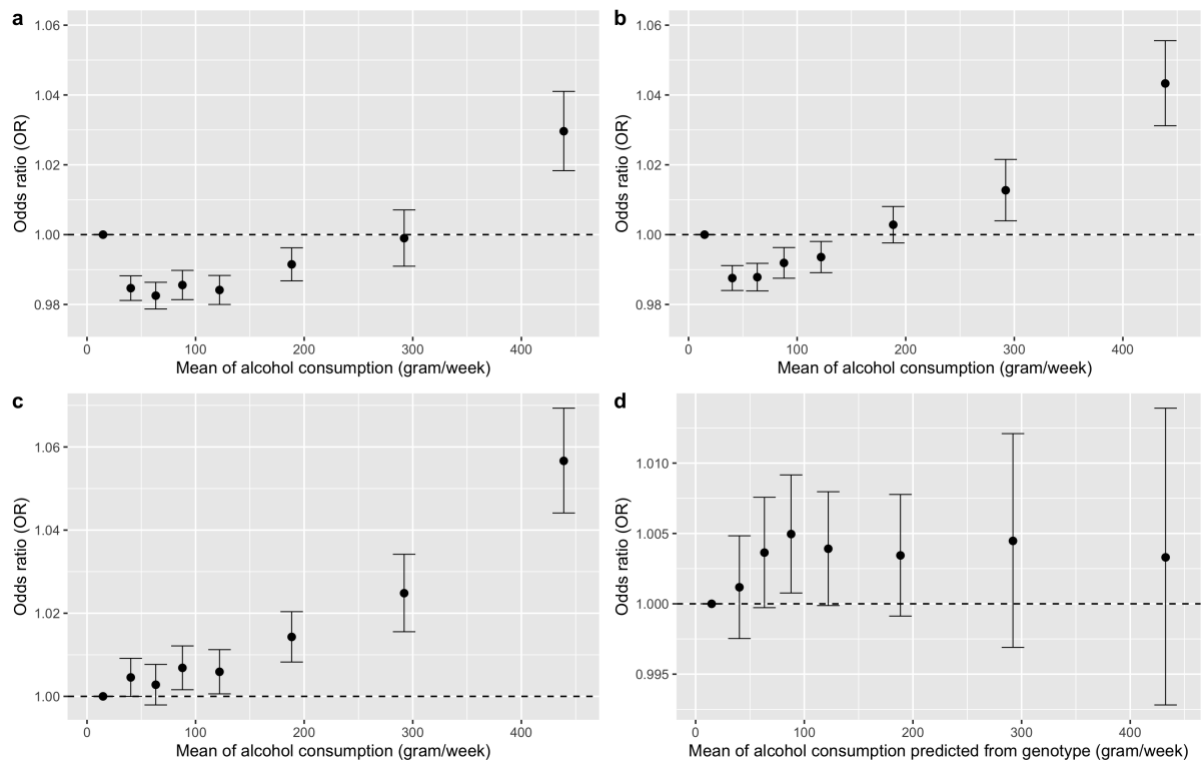
270

271

272

273

Supplementary Figure 15. Proportion of longitudinal change patterns and CVD prevalence in different AC level groups. (a) The x-axis shows eight AC level groups (measured by grams/week) as defined by the criteria in Wood et al.³. The y-axis shows the proportion of each longitudinal change group. (b) The y-axis denotes the prevalence of cardiovascular diseases. This x-axis is the same as in panel (a).



274

275

276

277

278

279

280

281

282

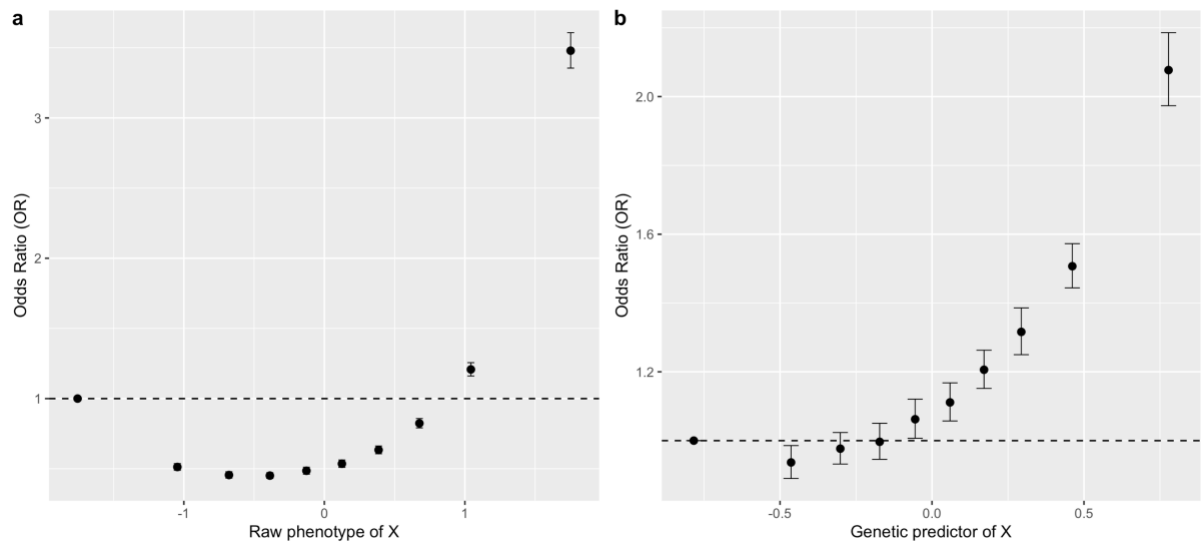
283

284

285

286

Supplementary Figure 16. The relationship between alcohol consumption and cardiovascular disease risk. The x-axes in panels a-c denote the mean alcohol consumption (gram/week) in each intake level group. The y-axes in all the panels denote the cardiovascular disease risk, measured by odds ratio (OR), against the reference group ($0 < \text{intake level} \leq 25$ grams/week). (a) The regression was performed in all current drinkers ($n = 356,138$). (b) The individuals likely to underreport AC or reduce intake due to illness or doctor's advice were removed, and the logistic regression was adjusted for the longitudinal changes ($n = 347,356$). (c) Individuals from the LESS group were removed from the reference group ($n = 319,320$). (d) The x-axis denotes the genetically predicted alcohol consumption ($n = 347,329$). Each dot indicates the OR estimated against the reference group, and the error bars in all the panels indicate the 95% confidence intervals of the estimates.



287

288

289

290

291

292

293

294

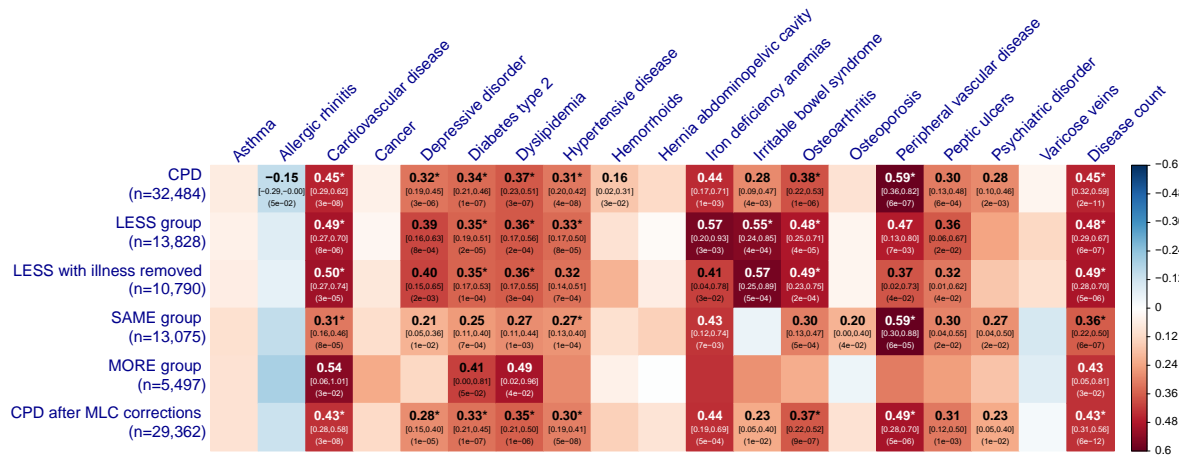
295

296

297

298

Supplementary Figure 17. J curve relationship between the genetic predictor of a trait X and the raw phenotype of a trait Y if the true relationship between X and Y is J-shaped. This figure demonstrates that if we predict X from genotypes, it is expected to have a J curve relationship with the raw phenotype of Y if the true relationship is a J curve. The total sample size used in the simulation $n = 50,000$. (a) The x-axis represents the simulated phenotypic value of X divided into 10 deciles. The y-axis represents the disease risk, as measured by odds ratio (OR) of each decile against the reference group (the first decile of X. (b) The x-axis represents the genetic predictor of X divided into 10 deciles. The y-axis represents the OR of each decile against the first decile. Each dot indicates the OR estimated against the reference group, and the error bars in all the panels indicate the 95% confidence intervals of the estimates.



300

301

302

303

304

305

306

307

308

309

310

311

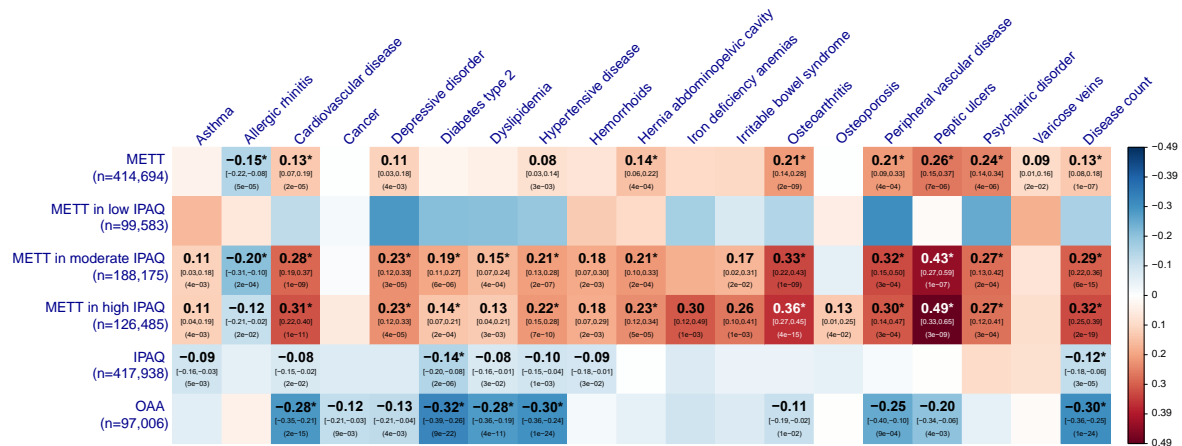
312

313

Supplementary Figure 18. Estimates of genetic correlation between cigarettes per day and common diseases in the UKB. The rows denote 6 GWAS summary data sets for cigarettes per day (CPD). The columns are 18 common diseases as well as disease count. The nominally significant estimates (P -value < 0.05) are labelled with the \hat{r}_g [95% confidence interval] (P -value), and the significant estimates after multiple corrections (P -value $< 0.05/114$) are labelled with an additional asterisk. The genetic correlation is estimated from cross-trait LD score regression. The P -value shown in the block is the original P -value for \hat{r}_g (two-sided χ^2 test). CPD represents the CPD in all current smokers; LESS, SAME, and MORE groups represent the CPD within the group who reduced, maintained the same, or increased the amount of smoking, respectively, compared to 10 years ago. "LESS with illness removed" represents the CPD in the LESS group excluding individuals who reduced smoking due to illness or doctor's advice. "CPD after MLC corrections" represents the CPD after the MLC corrections.



314 **Supplementary Figure 19. Estimates of genetic correlation between physical activity traits.** The
 315 value in each cell below the diagonal denotes the r_g estimate from a bivariate LDSC analysis. The
 316 circle in each cell above the diagonal shows the r_g estimate visually: larger circle size and darker color
 317 indicate higher r_g estimate. METT: Metabolic Equivalent Task in Total. IPAQ: International Physical
 318 Activity Questionnaire, short form. METT_low/moderate/high: METT in each of the three IPAQ
 319 categories. OAA: overall acceleration average measured by wrist-worn accelerometers. The estimates
 320 with P -value > 0.05 are annotated with a cross.



321

322 **Supplementary Figure 20. Estimates of genetic correlation between physical activity traits and**

323 **18 common diseases.** METT: Summed MET minutes per week for all activities.

324 METT_low/moderate/high IPAQ: METT in each of the three IPAQ categories. IPAQ: International

325 Physical Activity Questionnaire, short form. OAA: overall acceleration average. The columns are 18

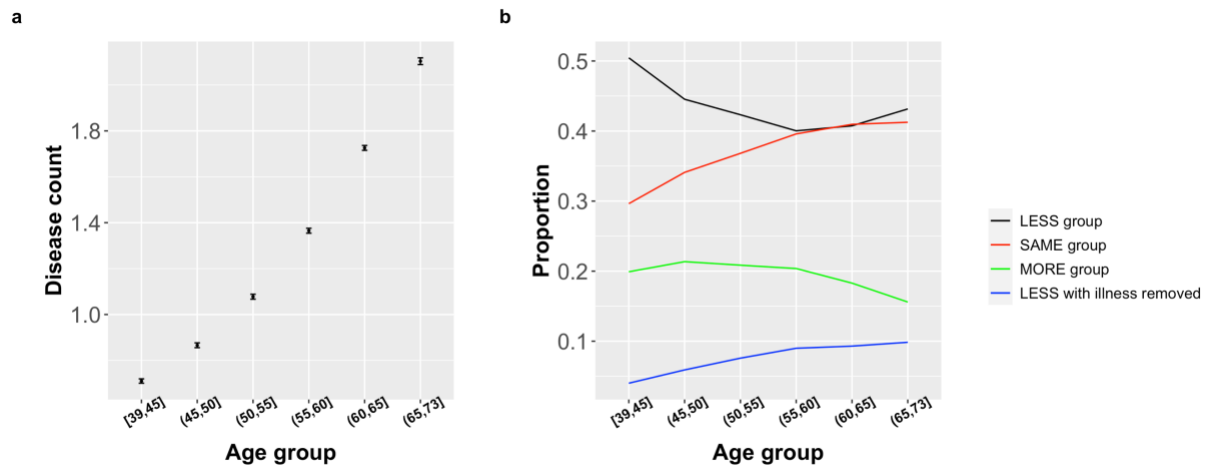
326 common diseases along with disease count. The nominally significant estimates (P -value < 0.05) are

327 labelled with the $\hat{\rho}_g$ [95% confidence interval] (P -value), and the significant estimates after multiple

328 corrections (P -value < 0.05/114) are labelled with an additional asterisk. The genetic correlation is

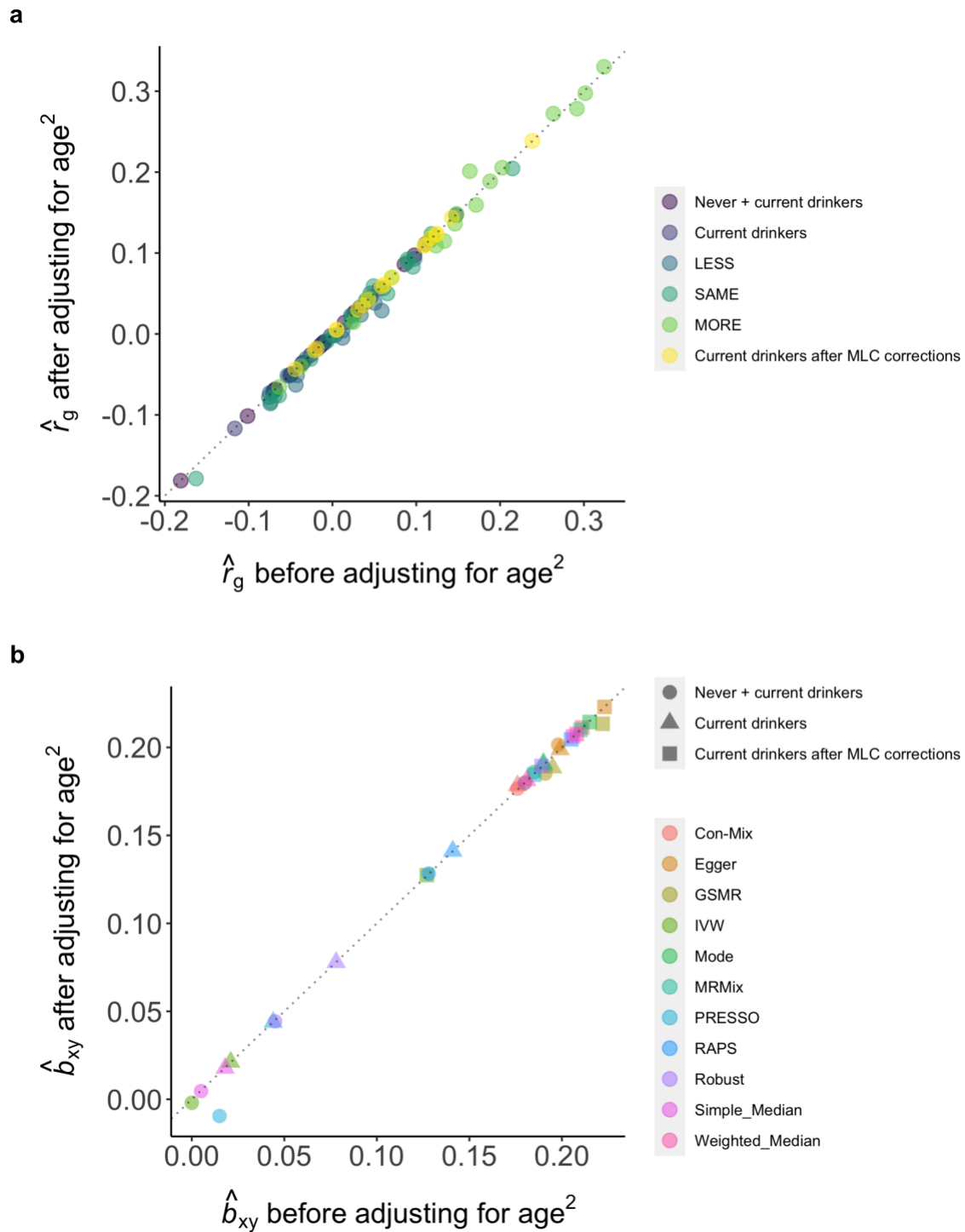
329 estimated from cross-trait LD score regression. The P -value shown in the block is the original P -value

330 for $\hat{\rho}_g$ (two-sided χ^2 test).

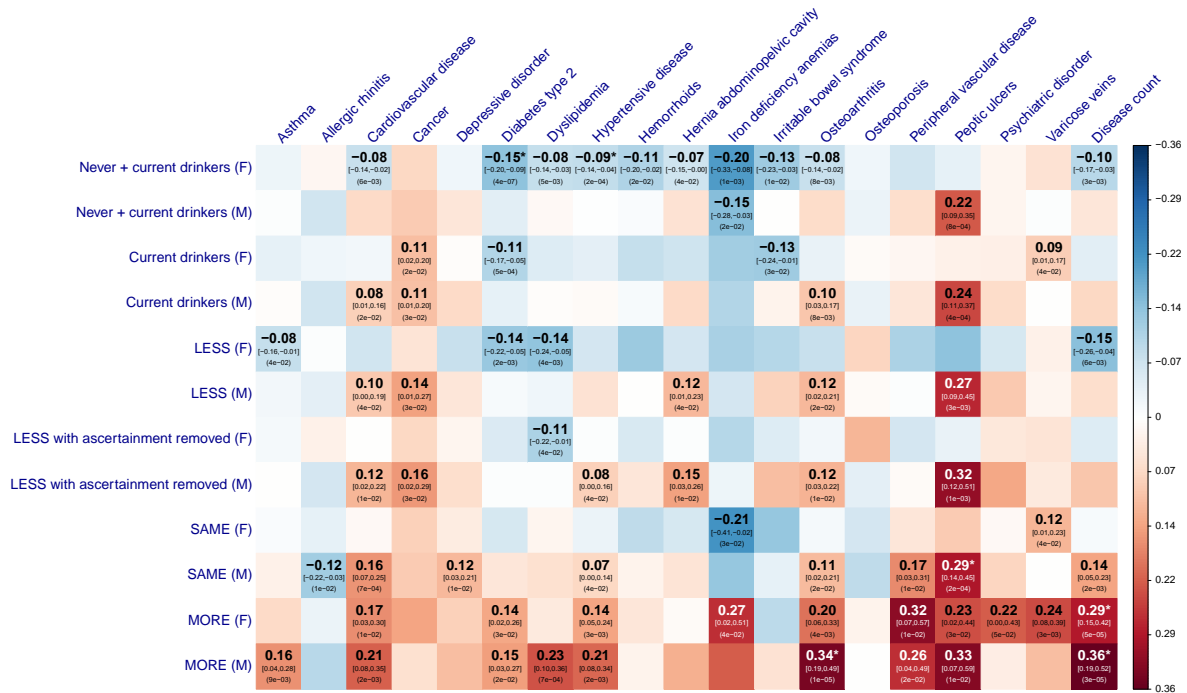


331

332 **Supplementary Figure 21. Disease count and ascertainment are age dependent.** The x-axis
 333 indicates 6 different age groups. (a) The y-axis indicates the average disease count in each age group.
 334 Each dot indicate the mean disease count in each age group, with the error bars indicating 95%
 335 confidence intervals. (b) The y-axis indicates the proportion of each longitudinal change group. The
 336 four groups are annotated by different colours. "LESS with illness removed" represents individuals
 337 who reduced drinking because of illness or doctor's advice, compared to 10 years ago.



338
 339 **Supplementary Figure 22. Comparison of the estimates of genetic correlation or causal**
 340 **association between AC and BMI before and after adjusting for age².** Panels (a) and (b) shows the
 341 results for genetic correlation and causal effect, respectively. The grey dashed line is the diagonal line.
 342 In panel (a), the color of the dot indicates different GWAS scenarios. In the panel (b), the shape of the
 343 point indicates different GWAS scenarios, and the color indicates different MR methods.
 344
 345



346
347
348
349
350
351
352
353
354
355
356
357

Supplementary Figure 23. Estimates of genetic correlation between AC and 18 common diseases in males and females separately under different scenarios. The rows denote 12 GWAS summary data sets for AC with the sex group labelled in the bracket ("F" for females and "M" for males). The columns are 18 common diseases along with the disease count. The nominally significant effects ($P < 0.05$) are labelled with r_g [95% confidence interval] (P -value), and the significant effects passing multiple testing correction ($P < 0.05/228$) are labelled with an additional asterisk. The genetic correlation is estimated from cross-trait LD score regression. The P -value shown in the block is the original P -value for \hat{r}_g (two-sided χ^2 test). LESS, SAME, and MORE represent current drinkers whose AC levels were reduced, maintained the same, and increased, respectively, compared to 10 years ago. "LESS with ascertainment removed" represents the LESS group excluding the participants who reduced their AC intake level due to illness or doctor's advice.

358 **Supplementary References**

- 359 1. Klatsky, A.L., Gunderson, E.P., Kipp, H., Udaltsova, N. & Friedman, G.D. Higher prevalence
360 of systemic hypertension among moderate alcohol drinkers: an exploration of the role of
361 underreporting. *Journal of Studies on Alcohol* **67**, 421-8 (2006).
- 362 2. Qi, T. *et al.* Identifying gene targets for brain-related traits using transcriptomic and
363 methylomic data from blood. *Nat Commun* **9**, 2282 (2018).
- 364 3. Wood, A.M. *et al.* Risk thresholds for alcohol consumption: combined analysis of individual-
365 participant data for 599 912 current drinkers in 83 prospective studies. *Lancet* **391**, 1513-1523
366 (2018).
- 367 4. Liu, M. *et al.* Association studies of up to 1.2 million individuals yield new insights into the
368 genetic etiology of tobacco and alcohol use. *Nat Genet* **51**, 237-244 (2019).
- 369