1  **Using genetic variation to disentangle the complex relationship between food intake and**

2  **health outcomes.**

3

4  **Authors**

5  Nicola Pirastu[1†], Ciara McDonnell[1,2*],Eryk J. Grzeszkowiak[1*], Ninon Mounier[3, 4], Fumiaki

6  Imamura[5], Jordi Merino[6,7,8], Felix R. Day[5], Jie Zheng[9], Nele Taba[10,11], Maria Pina Concas[11],

7  Linda Repetto[1], Katherine A. Kentistou[1,2], Antonietta Robino[12], Tõnu Esko[11,8],Peter K.

8  Joshi[1], Krista Fischer[11], Ken K. Ong[5],Tom R. Gaunt[9],  Zoltan Kutalik[3,4], John R. B. Perry[5],

9  James F. Wilson[1,13].

10

11  **Affiliations**

12  [1] Centre for Global Health Research, Usher Institute, University of Edinburgh, Teviot Place,

13  Edinburgh, EH8 9AG, Scotland.

14  [2] Centre for Cardiovascular Sciences, Queen's Medical Research Institute, University of

15  Edinburgh, Royal Infirmary of Edinburgh, Little France Crescent, Edinburgh EH16 4TJ,

16  Scotland

17  [3] Center for Primary Care and Public Health, University of Lausanne, Lausanne, Switzerland

18  [4] Swiss Institute of Bioinformatics, Lausanne, Switzerland

19  [5] MRC Epidemiology Unit, Institute of Metabolic Science, Cambridge Biomedical Campus,

20  University of Cambridge School of Clinical Medicine, Box 285, Cambridge, CB2 0QQ, UK

21  [6] Diabetes Unit and Center for Genomic Medicine, Massachusetts General Hospital, Boston,

22  MA, USA

23    [7] Program in Medical and Population Genetics, Broad Institute, Cambridge, MA, USA

24    [8] Department of Medicine, Harvard Medical School, Boston, MA, USA

25    [9] MRC Integrative Epidemiology Unit, Bristol Medical School, Bristol, UK

26    [10] Estonian Genome Center, Institute of Genomics, University of Tartu, Tartu, Riia 23b,

27    51010, Estonia

28    [11] Institute for Maternal and Child Health - IRCCS "Burlo Garofolo", Trieste, Italy

29    [12] Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Nobels väg

30    12A, SE-171 77 Stockholm, Sweden

31    [13] MRC Human Genetics Unit, Institute of Genetic and Molecular Medicine, University of

32    Edinburgh, Western General Hospital, Crewe Road, Edinburgh, EH4 2XU, Scotland

33    [14] Institute of Molecular and Cell Biology, University of Tartu, Tartu, Riia 23, 51010, Estonia

34

35    *Authors contributed equally to this work.

36    †Correspondence should be addressed to Nicola Pirastu, nicola.pirastu@ed.ac.uk

37

38

39

40

41

42

43

**Abstract:**

**Despite food choices being one of the most important factors influencing health, efforts to identify individual food groups and dietary patterns that cause disease have been challenging, with traditional nutritional epidemiological approaches plagued by biases and confounding. After identifying 302 individual genetic determinants of dietary intake in 445,779 individuals in the UK Biobank study, we develop a statistical genetics framework that enables us, to directly assess the impact of food choices on health outcomes. We show that the biases which affect observational studies extend also to GWAS, genetic correlations and causal inference through genetics, which can be corrected by applying our methods. Finally, by applying Mendelian Randomization approaches to the corrected results we identify some of the first robust causal associations between eating patterns and cancer, heart disease, obesity, and several other health related risk factors, distinguishing between the effects of specific foods or dietary patterns.**

## Introduction

Given their profound impact on human well-being, diet is one of the most studied human behaviours. Quality, quantity, and patterns of consumed foods are associated with a wide range of medical conditions such as metabolic, inflammatory, or mental health diseases[1]. However, despite the growing number of studies reporting associations between diet and health outcomes, it has been challenging to establish causal relationships due methodological limitations such as measurement error, confounding, and reverse causation. To date, several methods have been devised to try to account for intrinsic limitations in nutritional studies such as calibration of food records[2] or the implementation of domiciled feeding studies (ie. the PREDICT study[3]) in which participants are instructed to eat only the food provided by the study. Although these methods have helped in addressing some the limitations related to food consumption measurement, problems still remain especially when it comes to measure the effects of food on health over a long period of time.

In this context genetics may represent an alternative approach through the use of Mendellian Randomization. Mendelian Randomization (MR) is a methodological approach in which genetic variants associated with a phenotype of interest are used as instrumental variables to measure the "life-long effect of an exposure" to an outcome.[4] To date, several MR studies have been designed to investigate the associations between the consumption of single food groups, such as alcoholic beverages[5] , coffee[6],  milk[7–9] and specific health outcomes, but a systematic study investigating the overall role of diet is missing. In addition, previous MR studies have not accounted for the fact that genetic variants associated with reported dietary intake may be primarily associated with other risk factors or social determinants of health which may confound the causal estimates if used. In addition, previous studies on single food groups have not accounted for inter-relationships between different foods thus limiting the interpretability of the findings.

92    Given the complex number of factors that are driving the association between diet and health

93    outcomes, the present study was designed to initially identify the genetic variants associated

94    with reported food consumption, and then to leverage a causal inference statistical framework

95    to systematically investigate the causal effects of dietary factors on health outcomes, while

96    accounting for the effects that health determinants have on habitual dietary intake reporting.

97    **Methods**

98    **Study population and genome-wide association for dietary intake**

99    The UK Biobank[10] is a large population-based cohort including 500 000 adults aged between

100    40 and 69 years at baseline across 22 assessments centers in the United Kingdom. Data were

101    collected based on clinical examinations, assays of biological samples, detailed information

102    on self-reported health characteristics, and genome-wide genotyping. Dietary intake in UK

103    Biobank was assessed using a food frequency questionnaire which included questions about

104    the frequency of consumption specific foods and beverages over the past year. The number of

105    samples used for each trait can be found in table S1 while a detailed description of the

106    phenotypes, can be found in the in the supplementary methods 1.2 and table S2.

107    We used the BOLT-LMM software[11] to assess the association between the genetic variants

108    across the human genome and 29 food phenotypes. Analyses were conducted on genetic data

109    release version 3 imputed to the HRC panel[12], as provided by the UK Biobank

110    (http://www.ukbiobank.ac.uk/wp-

111    content/uploads/2014/04/UKBiobank_genotyping_QC_documentation-web.pdf). Population

112    stratification was assessed using LD-score regression as implemented in LD Hub[13,14] using

113    the LD scores provided with the software. Table S15 reports for each food trait the LD

114    regression intercept and heritability estimation using ldsc. Cluster analysis conducted on the

115    foods  identified 5 main groups of traits (see additional online methods paragraph 1.8 and 2.2

116    for details of group definition) and we thus set the genome-wide significance threshold at

117    $1 \times 10^{-8}$. Work within was conducted under UKB application 19655. Participants enrolled in

118    UK Biobank have signed consent forms. Replication analyses for identified signals

119    associated with food phenotypes were conducted independently by using genetic and dietary

120    data from the EPIC-Norfolk Study[15] and the Fenland Study[16]. Details additional online

121    methods 1.4.

122    **Investigating the effect of health outcomes on reported food intake using MR.**

123    Univariable MR analyses were initially conducted to measure the causal effect of health

124    outcomes on food consumption using the TwoSampleMR[17] R package. Exposures of interest

125    were selected amongst those for which nutritional advice is given and included body mass

126    index (BMI), low density lipoprotein cholesterol (LDLc), high density lipoprotein cholesterol

127    (HDLc), Total cholesterol, Triglycerides, Diastolic and Systolic blood pressure, Type 2

128    diabetes, and coronary artery disease. In addition, we included educational attainment as a

129    proxy of socio-economic status which is likely to affect food consumption. The full list of

130    studies from which the summary statistics were derived is detailed in Table S6. For each

131    exposure we selected all SNPs with $p < 5 \times 10^{-8}$ and $r^2 < 0.001$ to be used as instruments in the

132    MR analysis. After performing stepwise heterogeneity pruning we performed MR analysis

133    using the inverse variance method[18]. We then tested if the intercept from the MR-Egger[19]

134    regression was different from zero ($p < 0.05$). If this was the case, MR-Egger was used for the
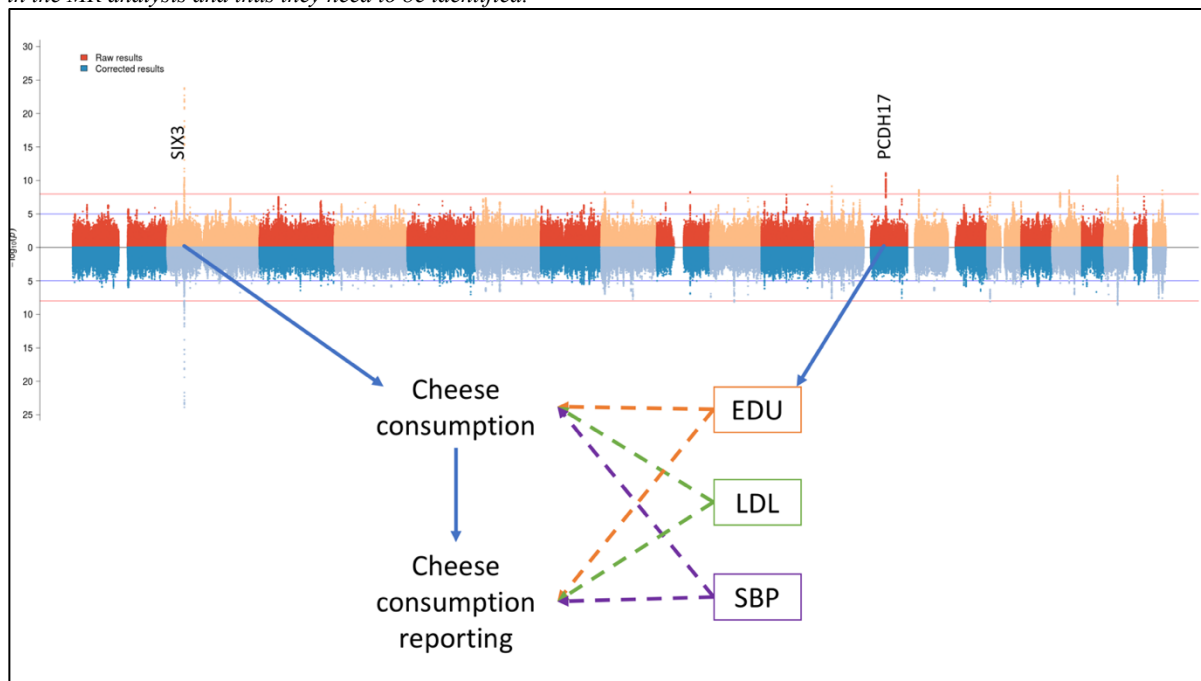
135    analysis instead.

136

137    **Measuring the direct effects of food types on health outcomes and identifying genetic**

138    **variants with predominantly direct-effects**

6

139     One of the most important assumptions in MR is that the effect of the instrument on the

140     outcome must be mediated only through the exposure of interest (sometimes referred as

141     exclusion restriction criteria)[20]. In this light the instruments whose effect on food is mediated

142     through the health outcomes or through educational attainment may violate this assumption

143     acting as confounders in the relationship between the exposure and the outcome. Moreover if

144     the mediating trait is acting on the reporting of food consumption and not food consumption

145     itself it would mean that the genetic variant  is not truly associated to food consumption and it

146     would thus not be a valid instrument. It is thus important to estimate the direct effect(i.e., the

147     effect that acts directly on food intake rather than is mediated through other factors see Figure

148     1)  the SNPs are exerting on actual food consumption in order to properly select the genetic

149     variants to be used as instrumental variables.

150     To this end we use a modified version of the method implemented in bGWAS[21]. This method

151     consists of a first step were the phenotype of interest (i.e., food consumption) is used as

152     outcome in multivariable MR. Next, exposures of interest are selected using a forward step

153     wise regression selection algorithm where each exposure is added until their p-value is less

154     than 0.05. The method provides a corrected estimate for each genetic variant of its effect on

155     the outcome trait once all mediated effects are removed. Further details can be found in

156     supplementary methods 1.6. In order to identify genetic variants with only a direct effect on

157     the phenotype of interest we defined the corrected to uncorrected ratio (CUR) as the ratio

158     between the corrected and the uncorrected effects (see additional methods 1.7 for a detailed

159     explanation).

160     ***Fig. 1 Direct and indirect SNP effects.*** *The plot shows the causal path of exemplar genes identified for cheese consumption.*
161     *In the multivariable MR model cheese consumption is causally influenced by educational attainment (EDU), low density*
162     *lipoprotein cholesterol levels (LDL) and systolic blood pressure (SBP). The effect of PDCH17 and is mediated through*
163     *educational attainment, while SIX3 has a direct effect on cheese consumption. The mediated effects cannot be used reliably*
164     *as MR instruments as they could be affecting either consumption or its reporting. Moreover, they could act as confounders*

165     *in the MR analysis and thus they need to be identified.*



166

167     The threshold to define genetic variants with non-mediated effects (CUR=1±0.05) is based on

168     simulations provided in the supplementary note 2.1 and on the genetic variants with known

169     biological function (ie. bitter receptors). We defined as "non-mediated" those SNPs whose

170     CUR fell within the defined ranges while "uncertain" the others.  We applied bGWAS to all

171     29 food phenotypes. As potential mediators, we used the same cardiometabolic phenotypes as

172     before except total cholesterol to avoid collinearity issues with LDL and HDL cholesterol,

173     and we added summary statistics from Crohn's disease and ulcerative colitis as they are

174     likely to affect dietary patterns. A Detailed discussion of this approach can be found in

175     supplementary methods 1.6.

176     **Genome-wide genetic correlations between corrected dietary intake and health**

177     **outcomes.**

178     We used LD-score regression implemented in LD Hub[13,14] to estimate genome-wide genetic

179     correlations between dietary intake phenotypes and 844 health outcomes and intermediary

180     phenotypes. Genetic correlations were estimated both with the corrected and uncorrected

181     GWAS summary statistics using the bivariate LD-score regression model. Stratified LD-
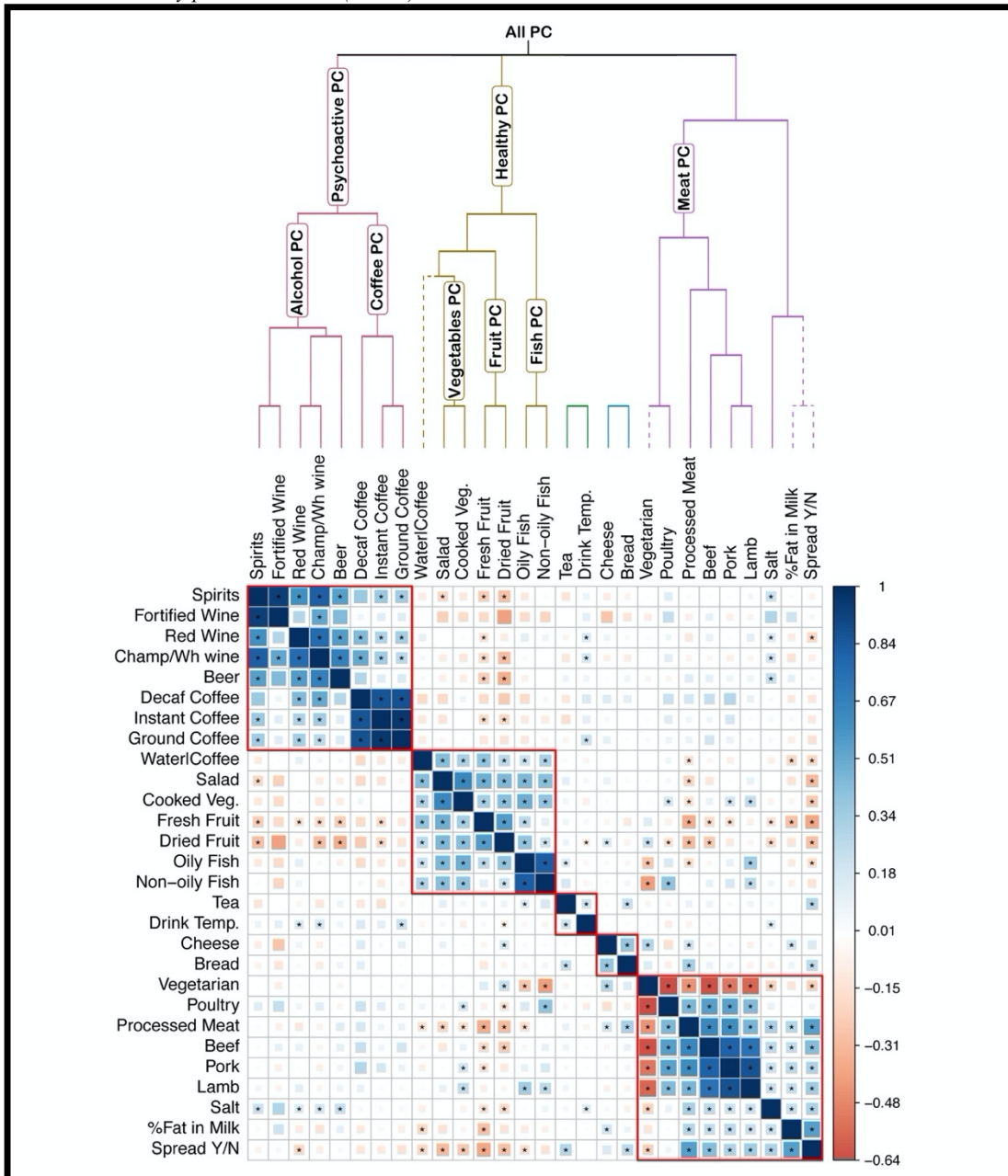
182    score regression[22] analyses were implemented using ldsc and the annotation files available on

183    the ldsc website.

**Definition of food group variables**

185    In order to define measures of dietary patterns we first performed cluster analysis of the 29

186    food items applying iCLUST[23] to the corrected genetic correlation matrix between the

187    different foods. iCLUST clusters items in different groups based on a hierarchical structure

188    (Details additional methods 1.8). Figure 2 shows the resulting dendrogram and its

189    comparison with the genetic correlation matrix.

190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232

233 ***Fig2 Clustering of the food traits and definition of measures of dietary patterns.*** *The plot reports the genetic correlation*
234 *plot amongst the food traits after applying the correction. The stars report the Bonferroni-corrected significant correlations.*
235 *The dendrogram and the boxes represent the clustering according to the ICLUST algorithm. The labels on the dendrogram*
236 *branches show the traits used to define each measure of dietary pattern. The dashed line represents the traits excluded from*
237 *the estimation of the dietary patterns traits. The "Vegetarian" trait was excluded from the "Meat PC" trait but was included*
238 *in the overall dietary pattern measure (All PC).*



239

240    We then defined based on the resulting structure several measures of dietary pattern at

241    different levels of the dendrogram as shown in Figure 2. For each measure we performed

242    principal component analysis of the items which participated to each group. The rotation

243    matrix was derived from the eigen decomposition of the correlation matrix of the foods in the

244    PC trait of interest. For example for the Coffee PC measure we performed principal

245     component analysis of "Ground Coffee", "Instant Coffee" and "Decaf Coffee". Once the

246     rotation matrix was estimated for each SNP its effect on the new measure was estimated as

247     the linear combination of the effect on each food trait using as weights the loadings on each

248     PC. A correlation plot of the loadings of each item onto the PC traits can be found in figure

249     S3.

250     **MR analyses to assess causal relationships between food intake and health outcomes**

251     MR analyses were conducted to estimate the effects of the food phenotypes on 79 health

252     related phenotypes (see table S17 for details) available in MR-base.[17] Genetic instruments for

253     each exposure of interest included independent genetic variants ($p < 5 \times 10^{-8}$ and pruning for

254     LD ($r^2 < 0.001$)). For dietary patterns exposures SNPs were selected as outlined in additional

255     methods 1.12. For the main analysis we restricted the genetic instruments to those that only

256     had evidence of a direct effect (i.e., not affecting the main exposure through a different

257     pathway; CUR $1 \pm 0.05$). Discussion of the relationship with other methods can be found in

258     supplementary note 2.7. Weights for the genetic instruments were based on the uncorrected

259     effects. To verify the effects of using only direct effect only SNPs on MR, all the analyses

260     were also conducted without applying the CUR filtering.

261     After selecting the genetic instruments, exposure and outcome data were harmonised. The

262     MR estimates were tested for heterogeneity and outliers were removed using the MR-Radial

263     method.[24] MR analyses were based on the inverse variance weighted method, which

264     estimates the causal effect of an exposure on an outcome by combining ratio estimates using

265     each variant. A random effect model was used if significant heterogeneity between the

266     different estimates was detected. We then tested for the presence of directional pleiotropy

267     using the intercept from the MR-Egger regression. MR median and MR-Raps were used as

268     sensitivity analyses. All results have been made available through an online app (

269     https://npirastu.shinyapps.io/Food_MR/) and can be found in additional table S18.

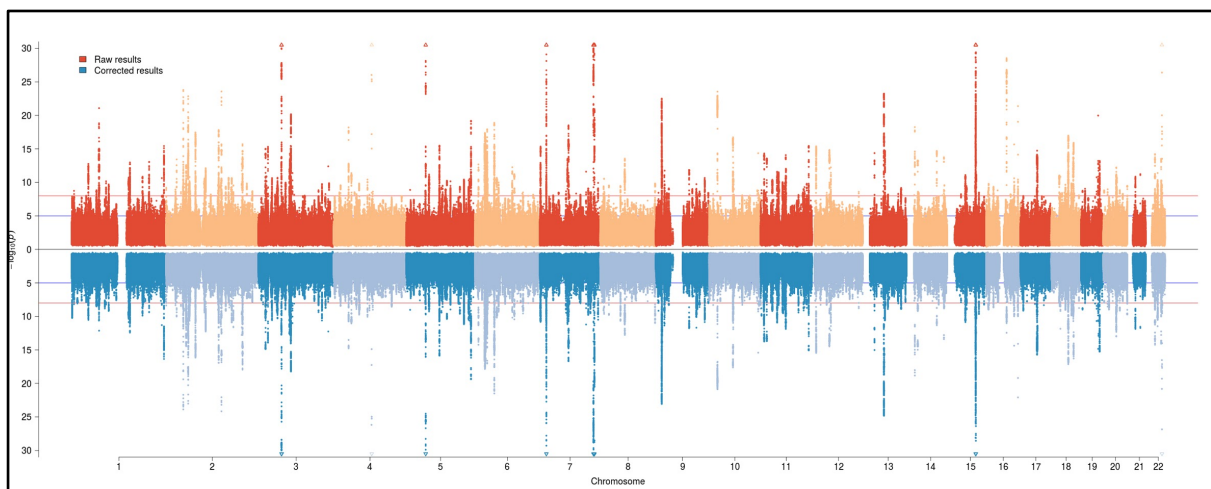270 **Patient and public involvement**

271 This research did not involve patients or the public as it uses data from the UK Biobank study

272 that were previously obtained from a cohort of people who had already been recruited. As

273 such, no patients or member of the public were involved in the design or implementation of

274 this study or the research questions addressed.

275 **Results**

276 **Genetic variants associated with food intake**

277 In a GWAS of 29 food phenotypes we identified 414 genetic associations in 260 independent

278 loci (Fig 3 and additional table S4) at Bonferroni corrected level of significance (P< $1\times10^{-8}$).

279 ***Fig. 3 302 independent genomic loci associate with food choices.*** *Results for both univariate (260 loci) and multivariate*
280 *(additional 42 loci see paragraph S2.3) analyses are included. For each SNP the lowest p-value for all traits was plotted.*
281 *The upper panel represents the unadjusted GWAS associations while the lower panel represents the association with food*
282 *choices, after adjustment for mediating traits, such as health status.*



283

284 Replication was sought in two additional UK-based cohorts including up to 32,779

285 participants. Despite relatively limited power in replication cohorts, concordant direction of

286 effect was observed for 82% of the signals (p=$7.82\times10^{-35}$, Binomial test; Table S5), and

287 nominal significance was achieved by 32% of the signals (p=$9.47\times10^{-54}$). Gene prioritization

288 is described in supplementary methods 1.10 while biological annotation, network analysis

289 and tissue enrichment analysis are discussed in additional paragraphs 1.11, 2.4 and 2.5.

290    Several of the identified loci have been previously associated with BMI. However, contrary

291    to our expectations, the BMI-raising allele was consistently associated with lower reported

292    consumption of energy-dense foods such as meat or fat, and higher reported intake of low-
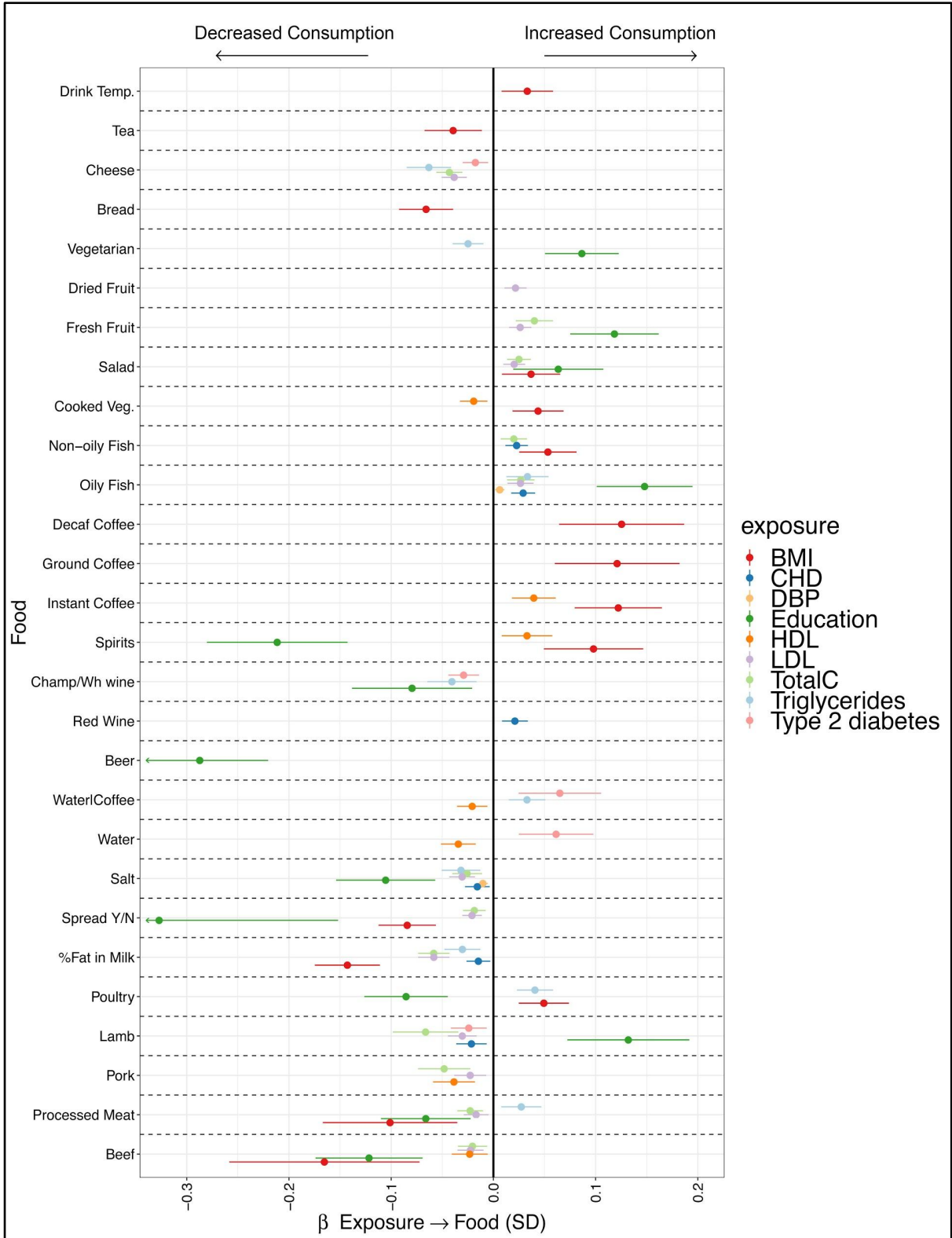
293    calorie foods.

294    **Genetic variants associated with food intake are strongly influenced by other**

295    **phenotypes**

296    In univariable MR we identified 81 instances in which health-related traits significantly

297    influencing food intake (Fig. 4 additional table S7). In particular BMI and Educational

298    attainment influenced more than 50% of the food traits. Similar effects extend to a broad

299    range of traits, for example LDL and triglycerides influenced 15 and 18 traits respectively.

300    Higher genetically-determined CAD associates with higher consumption of fish and red

301    wine, and lower consumption of whole milk, salt and lamb. These findings suggest that some

302    of the signals identified in GWAS for reported food phenotypes are not directly associated

303    with food intake but are mediated through a wide range of potential confounders.

304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321

322

323  **Fig 4. Health status influences reported food choices.** *The plot reports only the univariable MR results which were*
324  *significant at FDR<0.05. For each food outcome the effect estimate (β) is reported in standard deviations of the exposure*
325  *trait, together with 95% confidence intervals. Each colour represents a different exposure. BMI, body mass index; CHD,*
326  *coronary heart disease; DBP, diastolic blood pressure; HDL, high density lipoprotein cholesterol; LDL, low density*
327  *lipoprotein cholesterol; TotalC, total cholesterol. Champ/Wh wine, champagne, white wine. Temp, temperature.*
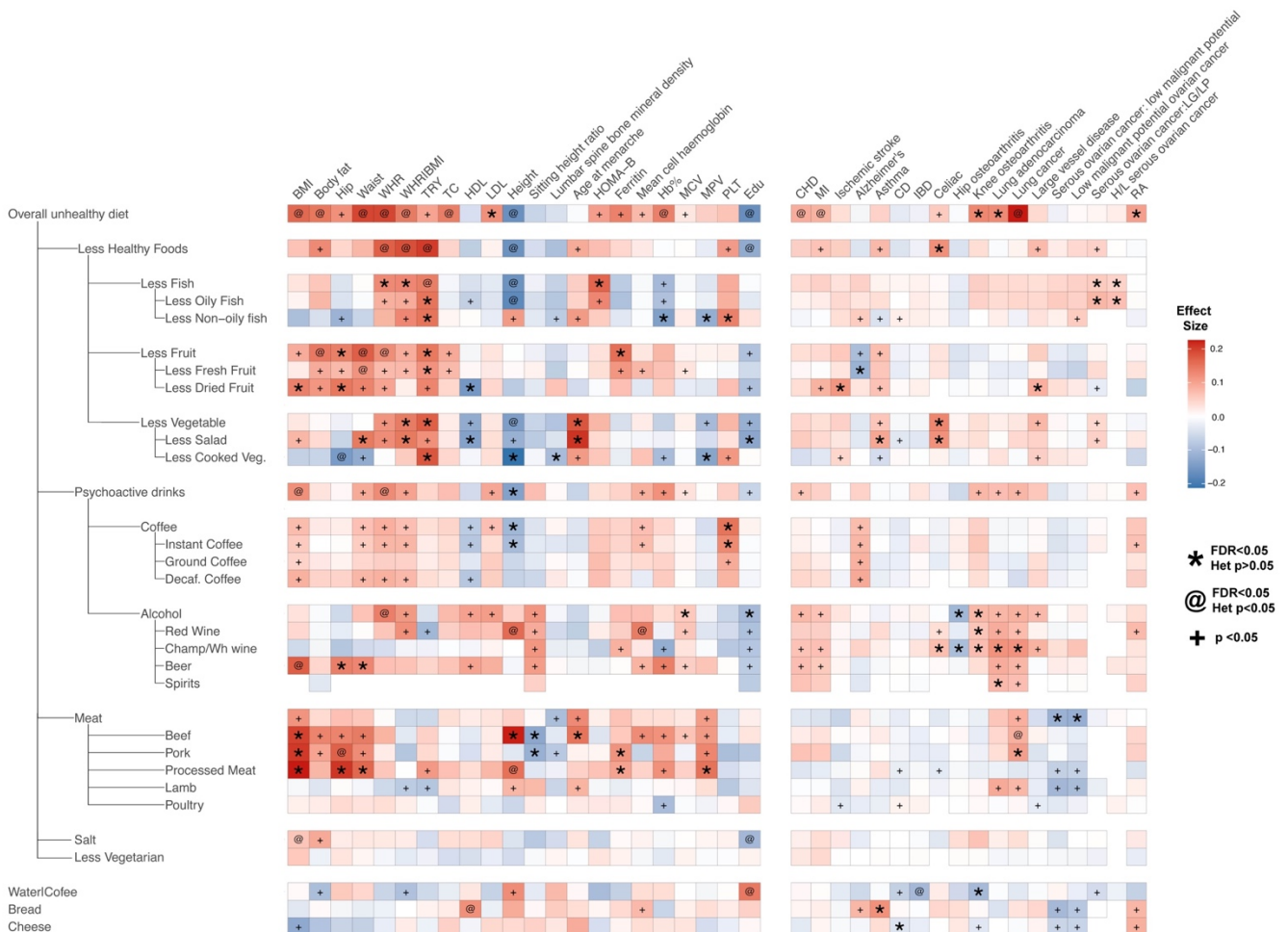


328

329 The Multivariable MR confirmed the univariable MR results (Supplementary Fig S4 panel A

330 and Supplementary Table S8). The percentage of genetic variance for the reported food

331 phenotypes explained by health determinants ranged from 42% for cheese to ~0% for

332 fortified wine and white wine/champagne (Supplementary Fig S3 panel B and Supplementary

333 Table S16). We systematically compared the estimated effect sizes of each genetic variants

334 influencing food consumption before and after correcting for the effect of health determinants

335 and showed that in many loci the variant initially identified for food phenotypes changed

336 dramatically after taking into account the effect of health factors (Fig. 3, see Supplementary

337 file 1 for trait-specific plots). For example, the effect size of the lead *FTO* variant

338 (rs55872725, $p$=2x10$^{-29}$) on milk fat percentage chosen decreased three-fold after accounting

339 for the mediated effects. To further explore the magnitude of this indirect effect on food

340 intake phenotypes, we compared the correlation patterns between the 29 food phenotypes and

341 832 phenotypes present in the LD hub[14] database identifying great differences. For example,

342 low fat milk intake was correlated with a beneficial effect on body fat percentage ($r_G$ = -0.43)

343 but this association diminished to near zero ($r_G$ = -0.04) after accounting for indirect effects

344 (Supplementary Data 2.2 and additional table S10). The effects of the correction procedure

345 on the genetic correlation amongst the traits and with the 844 health traits are discussed in

346 supplementary note 2.2 while full results can be found at in table S9 and browsed at

347 https://npirastu.shinyapps.io/rg_plotter_2/. These findings highlight the relevance of biases

348 and confounding in genetic correlation studies, and provide the framework to study complex

349 physiological relationships.

350 **Causal inference analyses for diet phenotypes and health outcomes**

351 A total of 230 out of 414 genetic variants initially associated with food phenotypes

352 (corresponding to 169/260 loci) were categorized as "non-mediated" associations (Table S3).

353 The balance of uncertain to non-mediated genetic associations varied by food group, ranging

15

354    from none uncertain for tea, spirits and processed meat, to all uncertain for percentage fat in

355    milk and adding spread to bread (Table S3).

356    In two-sample MR analyses we found 141 significant associations between food phenotypes

357    and health outcomes after multiple test correction (pFDR < 0.05, Table S18).

358    Of these 89 showed no sign of heterogeneity amongst the estimates (heterogeneity test p

359    >0.05). Figure 5 reports full results for all significant food exposure trait outcome pairs.

360    *Fig 5. Significant effects of food choice on disease related traits.* *The heatmap reports the results for all significant food*
361    *trait exposure trait outcome. Only dietary pattern exposures summarising the overall group consumption (PC1) have been*
362    *reported. All exposures have been aligned to have a positive loading onto the "overall unhealthy diet" measure. Significant*
363    *food/trait association are indicated with * if they show no sign of heterogeneity while @ if they show significant*
364    *heterogeneity. To facilitate meaningful visualisation and maximise the appearance of signal rather than noise, we applied a*
365    *shrinkage method - imposing a bayesian prior assumption on the distribution of beta (mean 0, SD 0.1), and conjugating that*
366    *with the likelihood of our results and then taking mean beta from the resulting distribution, thus shrinking estimates with*
367    *larger SEs more towards 0. Abbreviations: BMI Body Mass Index, WHR Waist to Hip Ratio, TRY tryglicerides, TC total*
368    *cholesterol, HDL HDL cholesterol, LDL LDL cholesterol, Hb% Haemoglobin percentage, MCV Mean Corpuscolar Volume,*
369    *MPV Mean Platelet Volume, PLT Platelet count, Edu Educational attainment, CHD Coronary Heart Disease, MI*
370    *Myocardial Infarction, CD Chron's Disease, IBD Inflammatory Bowel Disease, Serous ovarian cancer:LG/LP low grade*
371    *low potential. H/L serous ovarian cancer High and Low grade serous ovarian cancer, RA Rheumatoid Arthritis.*



372

373  Overall we found evidence supporting the beneficial effect of a healthy diet on health

374  outcomes. For example, for obesity/adiposity outcomes, genetically-determined unhealthy

375  diet leads to very similar effects across, increasing obesity measurements. For lipid-related

376  outcomes, the overall unhealthy diet is associated with higher levels of LDLc with no

377  significant heterogeneity, but no association with any of the other dietary traits. The overall

378  unhealthy diet was also strongly associated with Lung adenocarcinoma (OR 1.4xSD CI 1.2-

379  1.9) which seemed to be driven mostly by alcoholic beverages.

380  We identified 51 instances in which we would have not detected a significant result without

381  filtering out the non-direct effect instruments such as  the effect of increased fruit

382  consumption on triglycerides levels (estimated uncorrected effect= -0.03 (SE=0.05) vs.

383  estimated corrected effect = -0.17 (SE=0.05) ) or the effect of increased beef consumption on

384  height (uncorrected effect  = -0.02 (-0.17, 0.13) vs corrected effect = -0.52 (0.29, 0.74). In

385  addition, we found 124 food/trait relationships which were not significant after applying

386  CUR filtering, showing that either confounding effects or reduced power explain the lack of

387  association (see additional note 2.6). For example, red wine consumption was initially

388  associated with increased BMI (uncorrected effect =0.22 (SE 0.05)) and waist circumference

389  (uncorrected beta= 0.26 (SE 0.07), but after correcting for CAD liability, both effects

390  disappeared (corrected effect for BMI 0.05 (SE 0.06), corrected effect for WC 0.005 (0.08)).

391  On the flip side, we showed that the effect of red wine on mean corpuscular volume remains

392  substantially unchanged when applying the filtering approach (beta 0.07 (SE 0.02)

393  uncorrected and 0.065 (SE 0.02) corrected), suggesting that our approach could precisely

394  identify relevant biological relationships.

395  A full description of our findings are found in table S18 and have been made available

396  through an online app ( https://npirastu.shinyapps.io/Food_MR/).

397

**Discussion**

In this study we have provided quantitative data about the complex interplay between diet and health outcomes showing that the causal path from food intake to adverse health outcomes is not unidirectional and may be influenced by reverse causation and confounding even when MR is used. We showed that genetic correlations and causal inference can be improved by leveraging statistical approaches that take into account this mediated effects and identify genetic variants that have a only non-mediated effects on the exposure of interest. This information allowed us to perform causal inference analyses that helped identifying more reliable potential causal effects of food on health outcomes.

**Results in context**

Previous MR studies have mainly focused on specific food groups such as coffee, alcohol and milk consumption while none has comprehensively investigated the role of different food groups on health outcomes. Our results support previous observations such as the effect of alcohol consumption on coronary artery diseases reported in previous MR studies. In addition, we were able to confirm similar previous results detecting no evidence of an effect on IBD and CD[25] , ovarian cancer[26] or rheumatoid arthritis[27].

Findings from this study also suggest that the same biases that affect measures of food consumption such as reporting bias, confounding and reverse causation are reflected also in studies focusing on genetic associations. We have shown that these issues extend beyond obesity and socio-economic status including a broader range of intermediate factors. For example blood LDL and triglycerides concentration influence a wide variety of food traits thus being important factors to be considered as potential sources of bias, yet to our knowledge this is the first time this has been reported. For our analyses we have used UK biobank in which participants were aged between 40 and 60 at the time of the questionnaire,

18

422     it is likely that a younger cohort will suffer less from some of these (ie. LDL cholesterol or

423     blood pressure) as it is unlikely that they will display pathological level of these traits.

424     Our results are in contradiction to some previous studies in which no evidence of reverse

425     causation influencing genetic susceptibility for dietary patterns was reported.[28,29] We believe

426     that this difference is due to our novel approach, which is not based on using the potential

427     mediators as covariates, but rather exploits MR, which should be able to distinguish the

428     forward and reverse effects when the causal relationship is bidirectional. We have thus shown

429     that it is possible, through the use of available data and methods, to disentangle these

430     different colliding effects and to select the instrumental variables which show a non-mediated

431     effect, thus enabling the use of MR for the assessment of causal relationships between food

432     and health.

433      Many studies have looked at the relationship between nutritional composition and health

434     outcomes. One of the most salient examples is the relationship between saturated fat intake

435     and cardiovascular disease and all-cause mortality, in which recent studies suggest that food

436     sources of saturated fatty acids are more important than saturated fat content per se[Citation error].

437     Our study provide a new angle on the importance of food sources by providing evidence that

438     foods with similar nutrient profile, for example cheese and meat, which are both relatively

439     high in saturated fat and protein, have opposite effects on some metabolic risk factors such as

440     BMI (Figure S24 A) but there is no difference in other phenotypes such as blood lipids. A

441     similar conclusion can be drawn if we look at the foods which have the greatest effect on

442     triglycerides, fruit, vegetables and fish; all with very similar lowering effects (Figure S24 B),

443     which have relatively different macronutrient compositions. While the findings require

444     further investigations in mechanisms and related behaviours, our genetic evidence lends

445     support for the importance of studying foods in their complexity and not as a mere mixture of

446     nutrients. This approach, in fact, does not consider that the sources of the nutrients are not

447     equal due to the food matrix, the different preparations and that foods are seldom consumed

448     by themselves but in patterns which are likely to modify the effects on health.

449     Our findings illustrate that the effect of diet on health outcomes is complex, and components

450     of specific food groups have a differential association with health. In this case, although fish

451     and fruit and vegetables have a very different macronutrient composition it was impossible to

452     separate their effect on triglyceride concentrations. This suggests that at least in this case the

453     macronutrient composition is not as important as the an overall tendency to eat certain foods

454     and it highlights the importance of always including the assessment of dietary patterns before

455     claiming health effects of single foods or nutrients.

456     Some of the effects we have identified are more complex to explain and will need different

457     sources of evidence to be understood. For example we have found that the overall unhealthy

458     diet is associated to a higher risk of both lung andenocarcinoma and lung cancer. When

459     looking more closely to which of food explain this association the most we can see that

460     Alcohol seems to be driving the overall effect. One possibility is that this relationship is

461     confounded by smoking through a common tendency to addictive behaviours. However a

462     recent GWAS on cigarette smoking in Japan Biobank[30] reported a strong association between

463     the ALDH2 gene and number of cigarettes per day smoked which has also been associated to

464     differences in alcohol consumption[31], suggesting a causal effect of increased alcohol

465     consumption on increased smoking thus predisposing to lung cancer. Regardless of the

466     interpretation this example shows how complex the interpretation of MR results are when

467     behavioural traits are involved as they influence each other constantly creating a complex net

468     of interrelationships. This also points to the need of extreme care when claiming beneficial

469     health effects of food and multiple sources of evidence and approaches should always be

470     used before translating these findings into public policies.

471 Our study has several potential limitations. First, the number of items available in the dietary

472 questionnaire in the UK BioBank is limited, and therefore it limited our ability to capture

473 overall diet or specific food groups not detailed. The inclusion of white and relatively healthy

474 and educated participants from UK Biobank may have limited the generalisability of our

475 findings. Estimated effect sizes could be inflated because of the underestimation of the SNP

476 effects on the actual food trait consumption, rather than its self-report, if so, this will have

477 inflated our estimates of the effects of food on health, due to the noise in the questionnaire

478 responses, and warrants further statistical investigations. Even so, our method should not

479 have falsely identified a causal effect or reversed its direction, but further studies are needed

480 to assess the precise effect sizes.

481 In conclusion, our findings show that overall what is generally considered a healthy diet leads

482 to many favourable health outcomes and to reducing a wide range of risk factors broadly

483 agreeing with current guidelines aimed at reducing meat and alcohol consumption while

484 increasing fruit vegetables and fish. We also show that some of these effects are mostly

485 reconductable to specific food or group of foods which however are not characterized by

486 common nutrient composition thus adding granularity to our knowledge on the effect of diet

487 on health. This information can be useful to inform the design and implementation of future

488 studies to reduce the burden of diet-related diseases.

489 **Author Contributions**

490 NP,JFW,JRBP,ZK,EJG,FRD,KKO contributed to the study design.

491 JFW,TE,JRBP,AR,TG,FI,KKO,FRD contributed data.

492 NP,CMD,EJG,NM,FI,JZ,NT,KAK,MPC, performed the statistical analyses. NP, JFW,ZK,

493 JRBP, JM,TE, NT,KF,CMD,LR,EJG,FI,KKO,FRD contributed to the interpretation of the

494 results. All authors contributed to writing and editing of the text.

## Acknowledgements

**Data Availability**

All GWAS results will be made available through GWAS catalog at the time of publication.

All results from the MR analyses have been shared in the additional tables.

**References**

1. Misra, A. & Khurana, L. Obesity and the Metabolic Syndrome in Developing Countries. *The Journal of Clinical Endocrinology & Metabolism* vol. 93 s9–s30 (2008).

2. Naska, A., Lagiou, A. & Lagiou, P. Dietary assessment methods in epidemiological research: current state of the art and future prospects. *F1000Res.* **6**, 926 (2017).

3. Berry, S. E. *et al.* Human postprandial responses to food and potential for precision nutrition. *Nat. Med.* **26**, 964–973 (2020).

4. Zheng, J. *et al.* Recent Developments in Mendelian Randomization Studies. *Curr Epidemiol Rep* **4**, 330–345 (2017).

5. Millwood, I. Y. *et al.* Conventional and genetic evidence on alcohol and vascular disease aetiology: a prospective study of 500 000 men and women in China. *Lancet* **393**, 1831–1842 (2019).

6. Cornelis, M. C. & Munafo, M. R. Mendelian Randomization Studies of Coffee and Caffeine Consumption. *Nutrients* **10**, (2018).

7. Bergholdt, H. K. M., Nordestgaard, B. G., Varbo, A. & Ellervik, C. Milk intake is not associated with ischaemic heart disease in observational or Mendelian randomization analyses in 98,529 Danish adults. *Int. J. Epidemiol.* **44**, 587–603 (2015).

8. Vissers, L. E. T. *et al.* Dairy Product Intake and Risk of Type 2 Diabetes in EPIC-InterAct: A Mendelian Randomization Study. *Diabetes Care* **42**, 568–575 (2019).

23

542   9.   Hartwig, F. P., Horta, B. L., Smith, G. D., de Mola, C. L. & Victora, C. G. Association

543       of lactase persistence genotype with milk consumption, obesity and blood pressure: a

544       Mendelian randomization study in the 1982 Pelotas (Brazil) Birth Cohort, with a

545       systematic review and meta-analysis. *Int. J. Epidemiol.* **45**, 1573–1587 (2016).

546   10.  Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes of a

547       wide range of complex diseases of middle and old age. *PLoS Med.* **12**, e1001779 (2015).

548   11.  Loh, P.-R. *et al.* Efficient Bayesian mixed-model analysis increases association power in

549       large cohorts. *Nat. Genet.* **47**, 284–290 (2015).

550   12.  McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat.*

551       *Genet.* **48**, 1279–1283 (2016).

552   13.  Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from

553       polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).

554   14.  Zheng, J. *et al.* LD Hub: a centralized database and web interface to perform LD score

555       regression that maximizes the potential of summary level GWAS data for SNP

556       heritability and genetic correlation analysis. *Bioinformatics* vol. 33 272–279 (2017).

557   15.  Day, N. *et al.* EPIC-Norfolk: study design and characteristics of the cohort. European

558       Prospective Investigation of Cancer. *Br. J. Cancer* **80 Suppl 1**, 95–103 (1999).

559   16.  Lotta, L. A. *et al.* Integrative genomic analysis implicates limited peripheral adipose

560       storage capacity in the pathogenesis of human insulin resistance. *Nat. Genet.* **49**, 17–26

561       (2017).

562   17.  Hemani, G. *et al.* The MR-Base platform supports systematic causal inference across the

563       human phenome. *Elife* **7**, (2018).

564   18.  Burgess, S., Butterworth, A. & Thompson, S. G. Mendelian randomization analysis with

565    multiple genetic variants using summarized data. *Genet. Epidemiol.* **37**, 658–665 (2013).

566    19.  Bowden, J., Davey Smith, G. & Burgess, S. Mendelian randomization with invalid

567    instruments: effect estimation and bias detection through Egger regression. *Int. J.*

568    *Epidemiol.* **44**, 512–525 (2015).

569    20.  Davies, N. M., Holmes, M. V. & Davey Smith, G. Reading Mendelian randomisation

570    studies: a guide, glossary, and checklist for clinicians. *BMJ* **362**, k601 (2018).

571    21.  Mounier, N. & Kutalik, Z. bGWAS: an R package to perform Bayesian Genome Wide

572    Association Studies. *Bioinformatics* (2020) doi:10.1093/bioinformatics/btaa549.

573    22.  Finucane, H. K. *et al.* Partitioning heritability by functional category using GWAS

574    summary statistics. doi:10.1101/014241.

575    23.  Revelle, W. Hierarchical cluster analysis and the internal structure of tests. *Multivariate*

576    *Behav. Res.* **14**, 57–74 (1979).

577    24.  Bowden, J. *et al.* Improving the visualization, interpretation and analysis of two-sample

578    summary data Mendelian randomization via the Radial plot and Radial regression. *Int. J.*

579    *Epidemiol.* **47**, 2100 (2018).

580    25.  Georgiou, A. N., Ntritsos, G., Papadimitriou, N., Dimou, N. & Evangelou, E. Cigarette

581    smoking, coffee consumption, alcohol intake, and risk of crohn's disease and ulcerative

582    colitis: A Mendelian randomization study. *Inflamm. Bowel Dis.* (2020)

583    doi:10.1093/ibd/izaa152.

584    26.  Zhu, J., Jiang, X. & Niu, Z. Alcohol consumption and risk of breast and ovarian cancer:

585    A Mendelian randomization study. *Cancer Genet.* **245**, 35–41 (2020).

586    27.  Jiang, X. & Alfredsson, L. Modifiable environmental exposure and risk of rheumatoid

587    arthritis-current evidence from genetic studies. *Arthritis Res. Ther.* **22**, 154 (2020).

588    28.  Meddens, S. F. W., de Vlaming, R., Bowers, P. & Burik, C. A. P. Genomic analysis of

589        diet composition finds novel loci and associations with health and lifestyle. *bioRxiv*

590        (2018).

591    29.  Cole, J. B., Florez, J. C. & Hirschhorn, J. N. Comprehensive genomic analysis of dietary

592        habits in UK Biobank identifies hundreds of genetic loci and establishes causal

593        relationships between educational attainment and healthy eating. doi:10.1101/662239.

594    30.  Matoba, N. *et al.* GWAS of smoking behaviour in 165,436 Japanese people reveals

595        seven new loci and shared genetic architecture. *Nat. Hum. Behav.* **3**, 471–477 (2019).

596    31.  Matoba, N. *et al.* GWAS of 165,084 Japanese individuals identified nine loci associated

597        with dietary habits. *Nat. Hum. Behav.* **4**, 308–316 (2020).