

1 **Integrating untargeted metabolomics, genetically informed causal inference, and pathway**  
2 **enrichment to define the obesity metabolome**

3 Running title: Combine metabolomics and genetics to study obesity

4  
5 Yu-Han H. Hsu\*<sup>1-3</sup>, Christina M. Astley\*<sup>2,3</sup>, Joanne B. Cole<sup>2-4</sup>, Sailaja Vedantam<sup>2,3</sup>, Josep M.  
6 Mercader<sup>3,4</sup>, Andres Metspalu<sup>5</sup>, Krista Fischer<sup>5,6</sup>, Kristen Fortney<sup>7</sup>, Eric K. Morgen<sup>7</sup>, Clicerio  
7 Gonzalez<sup>8,9</sup>, Maria E. Gonzalez<sup>8,9</sup>, Tonu Esko<sup>5,3</sup>, Joel N. Hirschhorn<sup>#1-3</sup>

8  
9 <sup>1</sup> Department of Genetics, Harvard Medical School, Boston, Massachusetts, United States of  
10 America

11 <sup>2</sup> Division of Endocrinology and Center for Basic and Translational Obesity Research, Boston  
12 Children's Hospital, Boston, Massachusetts, United States of America

13 <sup>3</sup> Programs in Metabolism and Medical & Population Genetics, Broad Institute of Harvard and  
14 MIT, Cambridge, Massachusetts, United States of America

15 <sup>4</sup> Diabetes Unit and Center for Genomic Medicine, Massachusetts General Hospital, Boston,  
16 Massachusetts, United States of America

17 <sup>5</sup> Estonian Genome Center, Institute of Genomics, University of Tartu, Tartu, Estonia

18 <sup>6</sup> Institute of Mathematics and Statistics, University of Tartu, Tartu, Estonia

19 <sup>7</sup> BioAge Labs, Richmond, CA, United States of America

20 <sup>8</sup> Instituto Nacional de Salud Publica, Cuernavaca, Morelos, Mexico

21 <sup>9</sup> Centro de Estudios en Diabetes, Mexico City, Mexico

22

23 \* These authors contributed equally to this work.

24

25 # Corresponding author:

26 Joel N. Hirschhorn

27 Address: Boston Children's Hospital, 3 Blackfan Circle, CLS 16028, Boston, MA 02115, USA

28 Phone: (617) 919-2129

29 Email: Joel.Hirschhorn@childrens.harvard.edu

30

31 Competing Interests:

32 K.Fortney and E.K.M. are affiliated with BioAge Labs, Inc.; J.N.H. serves on the Scientific

33 Advisory Board of Camp4 Therapeutics.

34 **ABSTRACT**

35 **Background:** Obesity and its associated diseases are major health problems characterized  
36 by extensive metabolic disturbances. Understanding the causal connections between these  
37 phenotypes and variation in metabolite levels can uncover relevant biology and inform novel  
38 intervention strategies. Recent studies have combined metabolite profiling with genetic  
39 instrumental variable (IV) analyses to infer the direction of causality between metabolites and  
40 obesity, but often omitted a large portion of untargeted profiling data consisting of unknown,  
41 unidentified metabolite signals.

42 **Methods:** We expanded upon previous research by identifying body mass index (BMI)-  
43 associated metabolites in multiple untargeted metabolomics datasets, and then performing  
44 bidirectional IV analysis to classify these metabolites based on their inferred causal relationships  
45 with BMI. Meta-analysis and pathway analysis of both known and unknown metabolites across  
46 datasets were enabled by our recently developed bioinformatics suite, PAIRUP-MS.

47 **Results:** We identified 10 known metabolites that are more likely to be the causes (e.g.  
48 alpha-hydroxybutyrate) or effects (e.g. valine) of BMI, or may have more complex bidirectional  
49 cause-effect relationships with BMI (e.g. glycine). Importantly, we also identified about 5 times  
50 more unknown than known metabolites in each of these three categories. Pathway analysis  
51 incorporating both known and unknown metabolites prioritized 40 enriched ( $p < 0.05$ ) metabolite  
52 sets for the cause versus effect groups, providing further support that these two metabolite  
53 groups are linked to obesity via distinct biological mechanisms.

54 **Conclusions:** These findings demonstrate the potential utility of our approach to uncover  
55 causal connections with obesity from untargeted metabolomics datasets. Combining genetically  
56 informed causal inference with the ability to map unknown metabolites across datasets provides

57 a path to jointly analyze many untargeted datasets with obesity or other phenotypes. This  
58 approach, applied to larger datasets with genotype and untargeted metabolite data, should  
59 generate sufficient power for robust discovery and replication of causal biological connections  
60 between metabolites and various human diseases.

## 61 INTRODUCTION

62 Abnormal blood metabolite levels are important, frequent, and quantifiable feature of  
63 obesity and its associated phenotypes, which are major health problems globally<sup>1-5</sup>. Recently,  
64 systematic metabolite profiling (metabolomics) studies have described widespread alterations in  
65 the obesity metabolome and identified metabolite markers associated with risk of obesity-related  
66 diseases<sup>6-9</sup>. However, these studies broadly have two key analytic challenges limiting the  
67 biological interpretation and scope of their findings: these correlative studies have not generally  
68 been able to distinguish the cause and effect relationships between metabolites and phenotypes,  
69 and only a portion of the thousands of metabolite signals measured by untargeted profiling  
70 technology could be chemically identified and thereby routinely investigated.

71 Genetic instrumental variable (IV) analysis (for causal inference) and novel  
72 bioinformatics tools (for analysis of untargeted metabolite data) now provide the means to  
73 overcome these limitations and enhance our understanding of the metabolome of any phenotype.  
74 The genetic IV framework, also known as Mendelian randomization, uses genetic variants as  
75 instruments to infer causality from observational data in the presence of unmeasured  
76 confounding, provided certain methodological assumptions are met<sup>10,11</sup>. Bidirectional genetic IV  
77 analysis, using in turn genetic variants affecting metabolite levels and variants affecting a  
78 phenotype such as body mass index (BMI), offers a way to ascribe directionality of causal  
79 relationships and to prioritize potentially causal metabolite-phenotype associations. Previous  
80 genetic IV studies have utilized variants identified in genome-wide association studies (GWAS)  
81 to infer causality between obesity-related phenotypes and curated sets of metabolites (e.g.  
82 branched-chain and aromatic amino acids)<sup>12-16</sup>. However, most studies did not perform  
83 comprehensive bidirectional IV analysis and only focused on the metabolites that could be

84 identified and curated from profiling data, thus likely capturing only a limited slice of obesity  
85 biology and, even within that constraint, not assessing causality.

86 Previously, metabolites of unknown chemical identities – a large portion of untargeted  
87 profiling data – were mostly excluded from analyses (including GWAS) because inter-study  
88 comparison and biological interpretation were technically onerous or intractable<sup>17,18</sup>. To address  
89 these issues, we recently developed a bioinformatics suite, PAIRUP-MS<sup>18</sup>, to match up unknown  
90 metabolites across mass spectrometry-based untargeted profiling datasets, thereby enabling  
91 meta-analysis of multiple datasets and increasing statistical power for detecting biologically  
92 interesting unknowns. In addition, PAIRUP-MS provides a framework for annotating unknown  
93 metabolites using preexisting metabolic pathways and performing pathway analysis  
94 incorporating both known and unknown metabolites.

95 In this study, we demonstrate how the combination of bidirectional genetic IV framework  
96 and PAIRUP-MS can be used to analyze multiple untargeted metabolomics datasets and  
97 characterize causal connections between a phenotype and the metabolome. We identified both  
98 known and unknown BMI-associated metabolites, and then performed GWAS for each  
99 metabolite and for BMI, followed by bidirectional genetic IV analysis to identify metabolites  
100 likely to be causes or effects of obesity. In addition, we highlighted distinct biological pathways  
101 enriched for the cause versus effect metabolites, confirming that the bidirectional IV approach  
102 prioritized two distinct sets of BMI-associated metabolites. This initial work illustrates an  
103 approach that can now be generalized and scaled up to much larger datasets, which will enable  
104 well-powered studies to uncover novel metabolic causes and effects of obesity or any other  
105 phenotype of interest.

## 106 MATERIALS AND METHODS

107 A schematic overview of our analysis plan is shown in **Figure 1** and each step is  
108 described in more detail below.

109

### 110 **Metabolomics datasets and data processing**

111 *Study populations:* The study populations have been described previously<sup>18–20</sup>: (1)

112 Obesity Extremes (OE): N = 300 sampled equally from lean, obese, and the general Estonian

113 Biobank (EB) population, (2) Mexico City Diabetes Study (MCDS): N = 865 in a prospective

114 study, and (3) BioAge Labs Mortality Study (BioAge): N = 583 in a retrospective mortality

115 study nested in EB. All participants provided informed consent. Individual studies were

116 approved by their respective local ethics committees. Boston Children’s Hospital Institutional

117 Review Board approved this research.

118 *Metabolite data processing:* Untargeted liquid chromatography-mass spectrometry (LC-

119 MS) profiling of plasma samples, quality control, and missing value imputation of the data have

120 been described previously<sup>18</sup>. The processed OE dataset contained 298 samples and 13,613

121 metabolite signals (322 known); MCDS contained 821 samples and 7,136 signals (242 known);

122 BioAge contained 583 samples and 14,617 signals (603 known). Within each dataset, we

123 performed rank-based inverse normal transformation on each signal and used the resulting

124 abundance z-scores in downstream analyses. For OE and MCDS data used in BMI and genetic

125 association analyses, we performed covariate adjustment (age, sex, and fasting time for OE; age

126 and sex for MCDS) before the transformation. In this paper, we refer to both known and

127 unknown metabolite signals as “metabolites” for simplicity, recognizing that an unknown signal

128 does not always represent an independent, functional circulating metabolite.

129

## 130 **Mapping and identifying BMI-associated metabolites (Figure 1a)**

131 *Mapping metabolites across datasets:* Using the imputation-based matching algorithm in  
132 PAIRUP-MS<sup>18</sup>, we identified 1,780 metabolite pairs (207 shared known metabolites measured in  
133 both datasets and 1,573 matched unknown or unshared known metabolites) that could be  
134 compared directly across OE and MCDS and restricted subsequent analyses to these metabolites.  
135 For pathway analyses requiring the BioAge-based metabolite set annotations (see below), we  
136 furthered mapped 1,743 (200 shared known and 1,543 matched) of these metabolite pairs to  
137 metabolites measured in BioAge.

138 *Identifying BMI-associated metabolites:* Within each cohort, we adjusted raw BMI  
139 (available for 298 OE and 818 MCDS samples) for age and sex, performed rank-based inverse  
140 normal transformation on the residuals, and used the resulting BMI  $z$ -scores in all further  
141 analyses. (Since the OE obese and lean samples were drawn from the BMI extremes of EB, all  
142 EB samples were used to calculate population-based  $z$ -scores.) To identify BMI-associated  
143 metabolites, we performed linear regression of BMI on each metabolite within each dataset,  
144 followed by inverse variance weighted meta-analysis across the two datasets, and applied a  
145 Bonferroni significance threshold ( $p < 0.05/1,780$ ) in the meta-analysis.

146

## 147 **Bidirectional instrumental variable analyses (Figure 1b)**

148 *Metabolite instrument ( $G_M$ ) selection:* GWAS of the BMI-associated metabolites using  
149 294 OE and 637 MCDS samples (with available genetic data) and subsequent inverse variance  
150 weighted meta-analysis were performed as described previously<sup>18</sup>. To select  $G_M$ , we first  
151 identified the SNP (single nucleotide polymorphism) with the best meta-analyzed  $p$ -value for



152 each metabolite. Next, to avoid using redundant  $G_M$ , we “clumped” the best SNPs for all  
153 metabolites to select independent SNPs that have  $r^2 < 0.5$  or are  $> 250$  kb apart, and only kept the  
154 independent SNPs as  $G_M$  (along with their best-associated metabolites) in further analyses. For  
155 known metabolites in our causality groups (see below), we performed an additional sensitivity  
156 analysis using (where available) genome-wide significant ( $p < 5 \times 10^{-8}$ ) SNPs from published  
157 metabolite GWAS<sup>21-26</sup> as individual  $G_M$ .

158 *BMI instrument ( $G_B$ ) selection:* We used 97 BMI-associated SNPs ( $G_b$ ) previously  
159 identified in GIANT<sup>27</sup> and their effect estimates ( $\beta_b$ ) in our UK Biobank (UKB) GWAS to  
160 calculate a weighted genetic risk score for use as  $G_B$  (i.e.  $G_B = \sum \beta_b \times G_b$ ). We performed BMI  
161 GWAS in UKB using 453,397 European-ancestry samples and sex-combined BMI  $z$ -scores,  
162 using BOLT-LMM<sup>28</sup> to account for relatedness and population structure (**Supplementary Text**  
163 **1**). Analysis of the UKB data was approved by its governing Research Ethics Committee and the  
164 Broad Institute Institutional Review Board. The GIANT, UKB, and metabolomics cohorts have  
165 no known sample overlap. We confirmed that  $G_B$  was significantly associated with BMI in OE  
166 and MCDS and that none of the  $G_b$  are in linkage disequilibrium ( $r^2 > 0.3$ ) with the selected  $G_M$ .

167 *Testing for metabolite-to-BMI causal effect using  $G_M$ :* The association between BMI and  
168 each  $G_M$  was extracted from the UKB GWAS summary statistics and used to calculate the Wald  
169 ratio IV effect estimate of the metabolite (shared known or matched pair) on BMI. The  $p$ -value  
170 for the Wald estimate was calculated using an asymptotic standard error estimate as described  
171 previously<sup>29</sup>. This  $p$ -value – a test of the null hypothesis of no causal effect of the metabolite –  
172 was used to rank metabolites as more or less likely to be causal for BMI.

173 *Testing for BMI-to-metabolite causal effect using  $G_B$ :* We performed linear regression of  
174 each BMI-associated metabolite on  $G_B$  in OE and MCDS separately, followed by inverse

175 variance weighted meta-analysis. The Wald ratio IV effect estimate of BMI on each metabolite  
176 was calculated using the meta-analyzed statistics, and the corresponding  $p$ -value was used to  
177 rank metabolites as more or less likely to be effects of BMI. As a sensitivity analysis, we  
178 performed the MR-PRESSO global test<sup>30</sup> to assess overall horizontal pleiotropy among the  
179 individual SNPs ( $G_b$ ) contained within  $G_B$ , using metabolite- $G_b$  association in the OE-MCDS  
180 meta-analysis and BMI- $G_b$  association in UKB for 96 of 97 BMI SNPs (rs2033529 was excluded  
181 due to absence in our metabolite GWAS).

182

### 183 **Defining cause, effect, and bidirectional metabolite groups (Figure 1c)**

184 To rank BMI-associated metabolites as more or less likely to be the causes or effects of  
185 obesity, we used the  $-\log_{10} p$ -value of the IV effect estimate for either the metabolite ( $G_M$ ) or  
186 BMI ( $G_B$ ) instrument, reasoning that the statistical significance of these estimates is informative.  
187 Metabolites in the top and bottom quartiles of these two  $p$ -value-based rankings were assigned to  
188 three distinct groups corresponding to different types of causal connections with BMI: (1)  
189 “cause”: metabolites that were ranked in the top quartile using  $G_M$  and the bottom using  $G_B$ , and  
190 thus are likely to be upstream causes for BMI; (2) “effect”: metabolites that were ranked in the  
191 bottom quartile using  $G_M$  and the top using  $G_B$ , and thus are likely to be downstream effects of  
192 BMI; (3) “bidirectional”: metabolites that were in the top quartiles of both rankings, suggesting  
193 complex bidirectional cause-effect relationships with BMI.

194

### 195 **Pathway analyses of the defined metabolite groups (Figure 1d)**

196 *Metabolite set annotations:* BioAge data and PAIRUP-MS were used to generate  
197 metabolite set annotations as described previously<sup>18</sup>. Briefly, pathway annotations from

198 ConsensusPathDB<sup>31</sup> were consolidated into 690 metabolite sets with unique metabolite  
199 combinations (i.e. one metabolite set may correspond to multiple pathways containing identical  
200 sets of metabolites). We then used metabolite correlations in BioAge to expand the metabolite  
201 sets to include both known and unknown metabolites, calculating a membership score for each  
202 metabolite in each set.

203 *Pathway analyses:* We applied the pathway analysis framework in PAIRUP-MS to  
204 identify enriched metabolite sets for the cause, effect, and bidirectional metabolite groups we  
205 defined. We compared each of the three groups individually versus all other BMI-associated  
206 metabolites and, in a fourth analysis, compared the cause versus effect groups. First, for each  
207 metabolite set in each comparison analysis, a two-tailed Wilcoxon rank-sum test was performed  
208 to compare the membership scores of the two groups of metabolites. Next, to account for  
209 correlation structure in our data, iterations of this procedure were performed using “null”  
210 metabolite groups to calculate a permutation-based enrichment *p*-value for each metabolite set  
211 **(Supplementary Figure 1).**

212

### 213 **Performing *m/z* query for unknown metabolites**

214 To assess if the unknown metabolites captured information redundant to the known  
215 metabolites in our dataset (and to look up potential identities of unknowns classified in the three  
216 causality groups), we performed *m/z* query as described previously<sup>18</sup>, using the “LC-MS Search”  
217 tool in the Human Metabolome Database (HMDB)<sup>32</sup>. The unknowns were annotated as an *m/z*-  
218 matched adduct of a known metabolite in our data, an *m/z*-matched adduct of an HMDB  
219 metabolite not identified in our data, or as a metabolite without a match in HMDB.

## 220 **RESULTS**

### 221 **Identifying known and unknown metabolites associated with BMI**

222 We used untargeted metabolomics data from OE and MCDS to identify metabolites  
223 associated with BMI. First, we identified 207 pairs of shared known metabolites measured in  
224 both cohorts, and used PAIRUP-MS to match 1,573 additional pairs of unknown or unshared  
225 known metabolites likely to represent identical or highly correlated metabolites. Then, by  
226 performing meta-analysis of both the shared known and matched pairs across the cohorts, we  
227 identified 577 BMI-associated metabolites at Bonferroni significance ( $p < 0.05/1,780$ ), the  
228 majority of which were unknown metabolites: 418 (72.4%) consisted of two paired unknown  
229 metabolites, 59 (10.2%) consisted of a known metabolite matched to an unknown metabolite,  
230 and only 100 (17.3%) consisted of shared known metabolites. When we clustered these  
231 metabolites, we observed metabolite clusters that consisted mostly or entirely of matched pairs of  
232 unknown chemical identities (**Supplementary Figure 2**). Therefore, including these unknown  
233 metabolites in downstream analyses increased the number of candidate metabolites by nearly  
234 five-fold, and allowed us to investigate aspects of obesity biology not represented by the curated,  
235 known metabolites.

236

### 237 **Identifying metabolites more likely to be causal for BMI**

238 Before we could determine whether the BMI-associated metabolites are likely to be  
239 causal for BMI, we first needed to identify the SNP best-associated with each metabolite to use  
240 as genetic instrument ( $G_M$  in **Figure 1**). We therefore performed GWAS of metabolite levels in  
241 both OE and MCDS, followed by meta-analysis. We identified genome-wide significant ( $p < 5 \times$   
242  $10^{-8}$ ) SNPs for 204 (35 shared known and 169 matched) of the BMI-associated metabolites

243 **(Figure 2)**; 66 (14 shared known and 52 matched) of these were also significant after correction  
244 for multiple hypothesis testing ( $p < 5 \times 10^{-8}/577$ ). Overall, the matched, unknown metabolites  
245 showed comparable degree of genetic associations as the shared known metabolites, even in loci  
246 not associated with any of the knowns. Analyzing the unknowns thus greatly improved our  
247 ability to obtain significant and novel genetic instruments for metabolite signals, despite a  
248 relatively small GWAS sample size.

249 We observed that all 577 BMI-associated metabolites had best-associated SNPs with at  
250 least suggestive significance (maximum  $p = 2.5 \times 10^{-6}$ ) and therefore considered the best-  
251 associated SNP for each metabolite as potential instrument. To avoid analyzing metabolites  
252 sharing the same instruments, we included only genetically independent  $G_M$  ( $r^2 < 0.5$  or  $> 250\text{kb}$   
253 apart) and the 324 (40 shared known and 284 matched) metabolites best-associated with these  
254 instruments in subsequent IV analyses (**Supplementary Table 1**). For each metabolite, we  
255 estimated the association between  $G_M$  and BMI using a large independent cohort, UKB, in a two-  
256 sample design to calculate the metabolite-to-BMI IV effect estimate. We identified 50 (11 shared  
257 known and 39 matched) metabolites with nominally significant ( $p < 0.05$ ) metabolite-to-BMI IV  
258  $p$ -values, which indicates that they are more likely to be upstream causes for BMI  
259 (**Supplementary Table 1**).

260

## 261 **Identifying metabolites more likely to be effects of BMI**

262 Next, to determine if the BMI-associated metabolites are likely to be effects of BMI, we  
263 combined 97 BMI SNPs previously identified in GIANT into a weighted genetic risk score using  
264 UKB effect estimates as weights. As expected, the score is a valid genetic instrument for BMI  
265 ( $G_B$  in **Figure 1**) in OE and MCDS (meta-analyzed BMI- $G_B$  association  $p = 5.9 \times 10^{-7}$ ). For each

266 of the 324 BMI-associated metabolites, we estimated the association between  $G_B$  and the  
267 metabolite using OE and MCDS data to calculate the BMI-to-metabolite IV effect estimate. A  
268 total of 56 (8 shared known and 48 matched) metabolites had nominally significant ( $p < 0.05$ )  
269 BMI-to-metabolite IV  $p$ -values and thus are more likely to be downstream effects of BMI  
270 (**Supplementary Table 1**).

271

### 272 **Defining cause, effect, and bidirectional metabolite groups**

273 In order to further characterize the causal relationships between BMI and its associated  
274 metabolites, we ranked the metabolites based on the significance of their  $G_M$  and  $G_B$  IV effect  
275 estimate  $p$ -values (i.e. metabolite-to-BMI or BMI-to-metabolite IV  $p$ -values, respectively), and  
276 classified a subset of them into “cause”, “effect”, or “bidirectional” groups using quartile cutoffs  
277 of the rankings (**Figure 3**). We defined 25 metabolites as more likely to be cause (5 shared  
278 known and 20 matched), 26 as more likely to be effect (3 shared known and 23 matched), and 19  
279 as more likely to be bidirectional (2 shared known and 17 matched) with respect to BMI. The  
280 shared known metabolites in each group are listed in **Table 1**; the top cause, effect, and  
281 bidirectional metabolites are alpha-hydroxybutyrate, valine, and glycine, respectively. Details for  
282 all metabolites in each group are shown in **Supplementary Table 1**. We also performed  $m/z$   
283 query in HMDB to obtain potential identities for the unknowns in the matched metabolite pairs  
284 (**Supplementary Table 2**) and found only 6 out of the 60 matched pairs to be potentially  
285 redundant with the known metabolites curated in our data. Hence, we identified about 5 times  
286 more matched, unknown metabolites in the three causality categories compared to only  
287 analyzing the known metabolites. In addition, we performed sensitivity analyses to assess how  
288 our genetic IV and classification scheme would be influenced by weak instrument or pleiotropy

289 bias (**Supplementary Text 2** and **Supplementary Table 3**); we obtained results that generally  
290 support the robustness of our approach.

291

### 292 **Prioritizing enriched pathways for cause, effect, and bidirectional metabolites**

293 We identified many more matched, unknown metabolite pairs in the cause, effect, and  
294 bidirectional groups compared to the shared known metabolites, but it is difficult to hypothesize  
295 on their roles in obesity biology without knowing their chemical identities. Therefore, to extract  
296 useful information from the unknowns and to gain clues about the biology broadly captured by  
297 the three causality groups, we performed PAIRUP-MS pathway analyses encompassing both  
298 known and unknown metabolites, using metabolite set annotations generated from a separate  
299 cohort, BioAge. First, we carried out three separate analyses to identify pathways with nominally  
300 significant ( $p < 0.05$ ) enrichment for metabolites in the cause, effect, or bidirectional groups,  
301 respectively, when compared against all other BMI-associated metabolites (**Supplementary**  
302 **Table 4**). While the most enriched metabolite sets in each analysis are associated with different  
303 pathways, several metabolite sets were enriched in multiple analyses (e.g. “NAD *de novo*  
304 biosynthesis” was enriched for both cause and effect metabolites).

305 Hence, in order to identify pathways that are the most distinct between the defined  
306 metabolite groups, we next performed a pathway analysis directly comparing the cause versus  
307 effect metabolites, prioritizing 40 metabolite sets at nominal significance ( $p < 0.05$ ;  
308 **Supplementary Table 4**). The 13 cause metabolite sets (in which cause metabolites have higher  
309 membership scores than effect metabolites) are associated with various pathways, such as those  
310 connected to inflammation (e.g. nitric oxide signaling), redox metabolism (e.g.  
311 cysteine/methionine metabolism), and appetite regulation (e.g. endocannabinoid signaling). The

312 27 effect metabolite sets also contain varied pathways including those related to lysine  
313 catabolism, neurobiology (e.g. addiction and catecholamine biosynthesis), and stress response  
314 (e.g. FoxO signaling). While the known metabolites in our analysis have been linked to some of  
315 the enriched metabolite sets in literature, the unknown metabolites contributed most of the data  
316 used to prioritize these sets.

317 Finally, to better visualize the distinguishing features between the cause versus effect  
318 metabolites in terms of their roles in biological pathways, we constructed a heat map of the  
319 metabolites' membership scores in the enriched metabolite sets using unsupervised clustering  
320 (**Figure 4**). The metabolites formed two major clusters consisting of metabolites that are mostly  
321 in the cause or effect groups, with a handful of metabolites clustering with the contrasting group  
322 (i.e. cause metabolite "misclassified" in the effect cluster or vice versa). Even more strikingly,  
323 the cause and effect metabolite sets formed two pure clusters consisting of all cause or all effect  
324 sets. This clustering pattern provides further evidence that the cause and effect metabolites we  
325 defined are involved in distinct biological processes and thus may be associated with BMI  
326 through different mechanisms.



327 **DISCUSSION**

328           The study of comprehensive metabolite profiles defines an exciting frontier in human  
329 pathophysiology. However, metabolite-phenotype associations discovered in metabolomics  
330 studies are often correlative in nature and additional causal inference approaches, such as genetic  
331 IV analysis, are required to help assess causality between metabolites and phenotypes.  
332 Furthermore, unknown metabolite signals are often filtered out prior to analysis of untargeted  
333 metabolomics data, greatly limiting investigation to *a priori* candidate metabolites, reducing the  
334 search space, and hindering downstream analyses such as pathway enrichment. Here we present  
335 a paradigm for combining untargeted metabolomics, genomics, and our recently described  
336 bioinformatics suite, PAIRUP-MS, to overcome these challenges. Using obesity as an exemplar  
337 state of metabolic dysregulation, we illustrate the potential utility of this approach to advance our  
338 understanding of causal connections in metabolic diseases.

339           In this study, we meta-analyzed hundreds of unknown metabolites from two cohorts  
340 using PAIRUP-MS, identifying novel associations between the unknowns, BMI, genetic  
341 variants, and biological pathways. Indeed, using bidirectional genetic IV analysis, we discovered  
342 about 5 times as many unknown than known metabolites with potential causal connections to  
343 BMI. While these unknowns are likely not all fully independent and functional circulating  
344 molecules, their associations with genetic variants and BMI, distinct from those with known  
345 metabolites, suggest that a sizable number of unknown metabolites reflect aspects of BMI  
346 biology not captured by known metabolites. Furthermore, the much larger number of candidate  
347 metabolites allowed us to perform PAIRUP-MS pathway analyses that account for potential  
348 redundancy, prioritizing biological pathways specific to the metabolites with cause or effect  
349 relationships to BMI. Because of the relatively small sample sizes of our cohorts, some of our

350 results did not meet stringent multiple hypothesis testing significance thresholds; nevertheless,  
351 they demonstrate a useful and generalizable analytic framework to probe the metabolome of  
352 obesity and other diseases as larger datasets become available.

353 We identified novel metabolites that may be causes of obesity, as well as replicating two  
354 known metabolites, valine and tyrosine, that may be the effects of BMI<sup>14</sup>. The strongest causal  
355 evidence among known metabolites was for alpha-hydroxybutyrate, which has been linked to  
356 insulin resistance, oxidative stress, glutathione biosynthesis, and mitochondrial dysfunction<sup>6,33,34</sup>.  
357 The oxidative stress and glutathione connections are especially intriguing since “glutathione-  
358 mediated detoxification” emerged as a significant causal pathway when we compared the cause  
359 and effect metabolite groups in pathway analysis. It is also notable that the IV effect estimate of  
360 alpha-hydroxybutyrate on BMI is protective while the observational association suggests this  
361 metabolite is obesogenic. We postulate that a mitochondrial dysfunction/altered redox state  
362 linked to high alpha-hydroxybutyrate level could lead to decreased weight gain, while shared  
363 common causes, such as an obesogenic diet, may lead to increases in both alpha-hydroxybutyrate  
364 level and BMI. This example highlights the advantage of genetic IV analyses over observational  
365 studies alone to explore the potential impact of a theoretical intervention targeted to obesity-  
366 associated metabolites that have yet to be fully characterized<sup>35,36</sup>.

367 The validity of genetic IV analysis rests upon several key assumptions. Specifically, the  
368 genetic instrument must explain variation in the exposure variable and the instrument must not  
369 be associated with the outcome variable except through its relationship with the exposure (no  
370 genetic pleiotropy). Weak instrument bias towards the null and pleiotropy bias away from the  
371 null may lead to misclassification of metabolites in our three causality groups. To address weak  
372 instrument bias for our known metabolite instruments, we performed sensitivity analysis using

373 stronger instruments from published metabolite GWAS, showing that our results are generally  
374 robust against weak instrument bias, although some misclassification is possible due to limited  
375 power of our internal instruments. However, we could not conduct similar analysis for the  
376 unknown metabolite instruments since there is not yet a straightforward way to obtain external  
377 instruments for comparison. To address pleiotropy bias for our BMI instrument, we used a  
378 recently developed method, MR-PRESSO, to show that our BMI IV estimates are likely robust  
379 against extreme cases of pleiotropy bias. We could not examine pleiotropy in the metabolite  
380 instruments due to the lack of multiple instruments for each metabolite (especially for the  
381 unknowns where additional instruments could not be obtained from published GWAS).

382         Larger GWAS of both known and unknown metabolites, conducted across multiple  
383 datasets, will make it possible to extend our paradigm to understand causal biological  
384 mechanisms for various metabolic diseases and alleviate the limitations described above. With  
385 more candidate metabolites and genetic instruments emerging from better-powered studies, our  
386 approach can be expanded to mediation analyses<sup>37</sup>, to pathway Mendelian randomization<sup>38</sup>, or to  
387 metabolite IV subsetting according to predicted biological pathway memberships<sup>39</sup>. In  
388 conclusion, this study showcases the benefit of combining untargeted metabolomics with a  
389 bidirectional genetic IV approach to define the metabolome of a major human disease state,  
390 obesity. We therefore advocate for broader sharing of untargeted metabolomics and genetic  
391 datasets, similar to the approach taken by international efforts to optimize GWAS of many other  
392 phenotypes. Broader sharing would improve power and reliability of methodological frameworks  
393 such as the one presented here, and would enable a fuller realization of the potential of  
394 metabolomics to generate important insights into human diseases.

395 **ACKNOWLEDGEMENTS**

396 We thank the Broad Metabolomics Platform and SIGMA T2D Consortium for sharing  
397 data resources. This research has been conducted using the UK Biobank Resource under  
398 Application Number 11898. This work was supported by the National Heart, Lung, and Blood  
399 Institute grant F31HL126581 (Y.H.H.), National Institute of Diabetes and Digestive and Kidney  
400 Diseases grants T32DK110919 (Y.H.H.), K12DK094721 (C.M.A), and R01DK075787 (J.N.H.),  
401 Endocrine Scholars Award (C.M.A.), Doris Duke Charitable Foundation grant 215205 (J.N.H.),  
402 Estonian Research Council grants IUT20-60 (A.M.), PUT1665 (K.Fischer), and PUT1660 (T.E.),  
403 European Union through Horizon 2020 grant 692145 (A.M.), and European Union through the  
404 European Regional Development Fund Project No. 2014-2020.4.01.15-0012 (A.M.).

405

406 **COMPETING INTERESTS**

407 K.Fortney and E.K.M. are affiliated with BioAge Labs, Inc.; J.N.H. serves on the  
408 Scientific Advisory Board of Camp4 Therapeutics.

409 **REFERENCES**

- 410 1 Ng M, Fleming T, Robinson M, Thomson B, Graetz N, Margono C *et al.* Global, regional,  
411 and national prevalence of overweight and obesity in children and adults during 1980-  
412 2013: a systematic analysis for the Global Burden of Disease Study 2013. *Lancet* 2014;  
413 **384**: 766–781.
- 414 2 Kahn SE, Hull RL, Utzschneider KM. Mechanisms linking obesity to insulin resistance  
415 and type 2 diabetes. *Nature* 2006; **444**: 840–846.
- 416 3 Poirier P, Giles TD, Bray GA, Hong Y, Stern JS, Pi-Sunyer FX *et al.* Obesity and  
417 cardiovascular disease: pathophysiology, evaluation, and effect of weight loss. *Arter*  
418 *Thromb Vasc Biol* 2006; **26**: 968–976.
- 419 4 Renehan AG, Tyson M, Egger M, Heller RF, Zwahlen M. Body-mass index and incidence  
420 of cancer: a systematic review and meta-analysis of prospective observational studies.  
421 *Lancet* 2008; **371**: 569–578.
- 422 5 Flegal KM, Kit BK, Orpana H, Graubard BI. Association of all-cause mortality with  
423 overweight and obesity using standard body mass index categories: a systematic review  
424 and meta-analysis. *JAMA* 2013; **309**: 71–82.
- 425 6 Newgard CB, An J, Bain JR, Muehlbauer MJ, Stevens RD, Lien LF *et al.* A Branched-  
426 Chain Amino Acid-Related Metabolic Signature that Differentiates Obese and Lean  
427 Humans and Contributes to Insulin Resistance. *Cell Metab* 2009; **9**: 311–326.
- 428 7 Ho JE, Larson MG, Ghorbani A, Cheng S, Chen MH, Keyes M *et al.* Metabolomic  
429 Profiles of Body Mass Index in the Framingham Heart Study Reveal Distinct  
430 Cardiometabolic Phenotypes. *PLoS One* 2016; **11**: e0148361.
- 431 8 Cheng S, Rhee EP, Larson MG, Lewis GD, McCabe EL, Shen D *et al.* Metabolite

- 432 profiling identifies pathways associated with metabolic risk in humans. *Circulation* 2012;  
433 **125**: 2222–2231.
- 434 9 Wang TJ, Larson MG, Vasan RS, Cheng S, Rhee EP, McCabe E *et al.* Metabolite profiles  
435 and the risk of developing diabetes. *Nat Med* 2011; **17**: 448–453.
- 436 10 Smith GD, Ebrahim S. Mendelian randomization: prospects, potentials, and limitations.  
437 *Int J Epidemiol* 2004; **33**: 30–42.
- 438 11 Davey Smith G, Hemani G. Mendelian randomization: genetic anchors for causal  
439 inference in epidemiological studies. *Hum Mol Genet* 2014; **23**: R89–98.
- 440 12 Fall T, Hagg S, Magi R, Ploner A, Fischer K, Horikoshi M *et al.* The role of adiposity in  
441 cardiometabolic traits: a Mendelian randomization analysis. *PLoS Med* 2013; **10**:  
442 e1001474.
- 443 13 Holmes M V, Lange LA, Palmer T, Lanktree MB, North KE, Almgueira B *et al.* Causal  
444 effects of body mass index on cardiometabolic traits and events: a Mendelian  
445 randomization analysis. *Am J Hum Genet* 2014; **94**: 198–208.
- 446 14 Wurtz P, Wang Q, Kangas AJ, Richmond RC, Skarp J, Tiainen M *et al.* Metabolic  
447 Signatures of Adiposity in Young Adults: Mendelian Randomization Analysis and Effects  
448 of Weight Change. *PLoS Med* 2014; **11**: e1001765.
- 449 15 Liu J, van Klinken JB, Semiz S, van Dijk KW, Verhoeven A, Hankemeier T *et al.* A  
450 Mendelian Randomization Study of Metabolite Profiles, Fasting Glucose, and Type 2  
451 Diabetes. *Diabetes* 2017; **66**: 2915–2926.
- 452 16 Haase CL, Tybjaerg-Hansen A, Qayyum AA, Schou J, Nordestgaard BG, Frikke-Schmidt  
453 R. LCAT, HDL cholesterol and ischemic cardiovascular disease: a Mendelian  
454 randomization study of HDL cholesterol in 54,500 individuals. *J Clin Endocrinol Metab*

- 455 2012; **97**: E248-56.
- 456 17 Patti GJ, Tautenhahn R, Siuzdak G. Meta-analysis of untargeted metabolomic data from  
457 multiple profiling experiments. *Nat Protoc* 2012; **7**: 508–516.
- 458 18 Hsu Y-HH, Churchhouse C, Pers TH, Mercader JM, Metspalu A, Fischer K *et al.*  
459 PAIRUP-MS: Pathway analysis and imputation to relate unknowns in profiles from mass  
460 spectrometry-based metabolite data. *PLoS Comput Biol* 2019; **15**: 1–26.
- 461 19 Leitsalu L, Haller T, Esko T, Tammesoo ML, Alavere H, Snieder H *et al.* Cohort Profile:  
462 Estonian Biobank of the Estonian Genome Center, University of Tartu. *Int J Epidemiol*  
463 2015; **44**: 1137–1147.
- 464 20 Williams Amy AL, Jacobs Suzanne SBR, Moreno-Macías H, Huerta-Chagoya A,  
465 Churchhouse C, Márquez-Luna C *et al.* Sequence variants in SLC16A11 are a common  
466 risk factor for type 2 diabetes in Mexico. *Nature* 2014; **506**: 97–101.
- 467 21 Rhee EP, Ho JE, Chen MH, Shen D, Cheng S, Larson MG *et al.* A genome-wide  
468 association study of the human metabolome in a community-based cohort. *Cell Metab*  
469 2013; **18**: 130–143.
- 470 22 Draisma HHM, Pool R, Kobl M, Jansen R, Petersen AK, Vaarhorst AAM *et al.* Genome-  
471 wide association study identifies novel genetic variants contributing to variation in blood  
472 metabolite levels. *Nat Commun* 2015; **6**: 7208.
- 473 23 Shin S-YY, Fauman EB, Petersen A-KK, Krumsiek J, Santos R, Huang J *et al.* An atlas of  
474 genetic influences on human blood metabolites. *Nat Genet* 2014; **46**: 543–50.
- 475 24 Long T, Hicks M, Yu HC, Biggs WH, Kirkness EF, Menni C *et al.* Whole-genome  
476 sequencing identifies common-to-rare variants associated with human blood metabolites.  
477 *Nat Genet* 2017; **49**: 568–578.

- 478 25 Burkhardt R, Kirsten H, Beutner F, Holdt LM, Gross A, Teren A *et al.* Integration of  
479 Genome-Wide SNP Data and Gene-Expression Profiles Reveals Six Novel Loci and  
480 Regulatory Mechanisms for Amino Acids and Acylcarnitines in Whole Blood. *PLoS*  
481 *Genet* 2015; **11**: e1005510.
- 482 26 Kettunen J, Demirkan A, Wurtz P, Draisma HH, Haller T, Rawal R *et al.* Genome-wide  
483 study for circulating metabolites identifies 62 loci and reveals novel systemic effects of  
484 LPA. *Nat Commun* 2016; **7**: 11122.
- 485 27 Locke AE, Kahali B, Berndt SI, Justice AE, Pers TH, Day FR *et al.* Genetic studies of  
486 body mass index yield new insights for obesity biology. *Nature* 2015; **518**: 197–206.
- 487 28 Loh PR, Tucker G, Bulik-Sullivan BK, Vilhjalmsdottir BJ, Finucane HK, Salem RM *et al.*  
488 Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat*  
489 *Genet* 2015; **47**: 284–290.
- 490 29 Burgess S, Small DS, Thompson SG. A review of instrumental variable estimators for  
491 Mendelian randomization. *Stat Methods Med Res* 2017; **26**: 2333–2355.
- 492 30 Verbanck M, Chen CY, Neale B, Do R. Detection of widespread horizontal pleiotropy in  
493 causal relationships inferred from Mendelian randomization between complex traits and  
494 diseases. *Nat Genet* 2018; **50**: 693–698.
- 495 31 Kamburov A, Stelzl U, Lehrach H, Herwig R. The ConsensusPathDB interaction  
496 database: 2013 update. *Nucleic Acids Res* 2013; **41**: D793-800.
- 497 32 Wishart DS, Feunang YD, Marcu A, Guo AC, Liang K, Vazquez-Fresno R *et al.* HMDB  
498 4.0: the human metabolome database for 2018. *Nucleic Acids Res* 2018; **46**: D608–D617.
- 499 33 Gall WE, Beebe K, Lawton KA, Adam KP, Mitchell MW, Nakhle PJ *et al.* alpha-  
500 hydroxybutyrate is an early biomarker of insulin resistance and glucose intolerance in a



- 501 nondiabetic population. *PLoS One* 2010; **5**: e10883.
- 502 34 Thompson Legault J, Strittmatter L, Tardif J, Sharma R, Tremblay-Vaillancourt V, Aubut  
503 C *et al.* A Metabolic Signature of Mitochondrial Dysfunction Revealed through a  
504 Monogenic Form of Leigh Syndrome. *Cell Rep* 2015; **13**: 981–989.
- 505 35 Burgess S, Harshfield E. Mendelian randomization to assess causal effects of blood lipids  
506 on coronary heart disease: Lessons from the past and applications to the future. *Curr.*  
507 *Opin. Endocrinol. Diabetes Obes.* 2016. doi:10.1097/MED.0000000000000230.
- 508 36 Voight BF, Peloso GM, Orho-Melander M, Frikke-Schmidt R, Barbalic M, Jensen MK *et*  
509 *al.* Plasma HDL cholesterol and risk of myocardial infarction: A mendelian randomisation  
510 study. *Lancet* 2012. doi:10.1016/S0140-6736(12)60312-2.
- 511 37 Van Der weele T, Vansteelandt S. Mediation analysis with multiple mediators. *Epidemiol*  
512 *Method* 2013. doi:10.1515/em-2012-0010.
- 513 38 Burgess S, Thompson SG. Multivariable Mendelian randomization: The use of pleiotropic  
514 genetic variants to estimate causal effects. *Am J Epidemiol* 2015; **181**: 251–260.
- 515 39 Wittemans LBL, Lotta LA, Oliver-Williams C, Stewart ID, Surendran P, Karthikeyan S *et*  
516 *al.* Assessing the causal association of glycine with risk of cardio-metabolic diseases. *Nat*  
517 *Commun* 2019; **10**: 1–13.

518

519 **FIGURE LEGENDS**

520 **Figure 1. Overview for identifying and characterizing causal connections in the obesity**

521 **metabolome. (a)** OE and MCDS metabolomics datasets, matched using PAIRUP-MS, were used  
522 to identify known and unknown metabolites associated with BMI. **(b)** Independent genetic  
523 instruments ( $G_M$ ) for the BMI-associated metabolites were selected using OE and MCDS data,  
524 and then used to test for a metabolite-to-BMI ( $M \rightarrow B$ ) causal effect in UKB; in parallel, BMI  
525 genetic instrument ( $G_B$ ), a polygenic risk score built using GIANT BMI-associated SNPs ( $G_b$ )  
526 and UKB effect estimate weights ( $\beta_b$ ), was used to test for a BMI-to-metabolite ( $B \rightarrow M$ ) causal  
527 effect in OE and MCDS. **(c)** A subset of metabolites was categorized into “cause”, “effect”, and  
528 “bidirectional” groups based on the significance of the  $G_M$  and  $G_B$  IV effect estimate  $p$ -values,  
529 reflecting different types of causal connections between the metabolites and BMI. **(d)** Pathway  
530 analyses of the three metabolite groups were performed using metabolite set annotations  
531 generated using PAIRUP-MS and BioAge data.  $Y \sim X$ , regression of  $Y$  on  $X$ ;  $U$ , unmeasured  
532 confounder;  $n$ , number of samples;  $m$ , number of metabolites;  $k$ , number of known (or shared  
533 known) metabolites;  $s$ , number of metabolite sets.

534

535 **Figure 2. Joint Manhattan plots summarizing GWAS of BMI-associated metabolites in OE**

536 **and MCDS.** Genetic associations for 100 shared known (top) or 477 matched (bottom) BMI-  
537 associated metabolites were consolidated to plot the best  $p$ -value for each SNP (i.e. only the  $p$ -  
538 value for the best associated metabolite was plotted for each SNP). Genome-wide significance  
539 threshold ( $p < 5 \times 10^{-8}$ ) is marked by the orange lines. Genome-wide significant SNPs are plotted  
540 in red or blue, for shared known or matched metabolites, respectively. Lead SNPs of the most

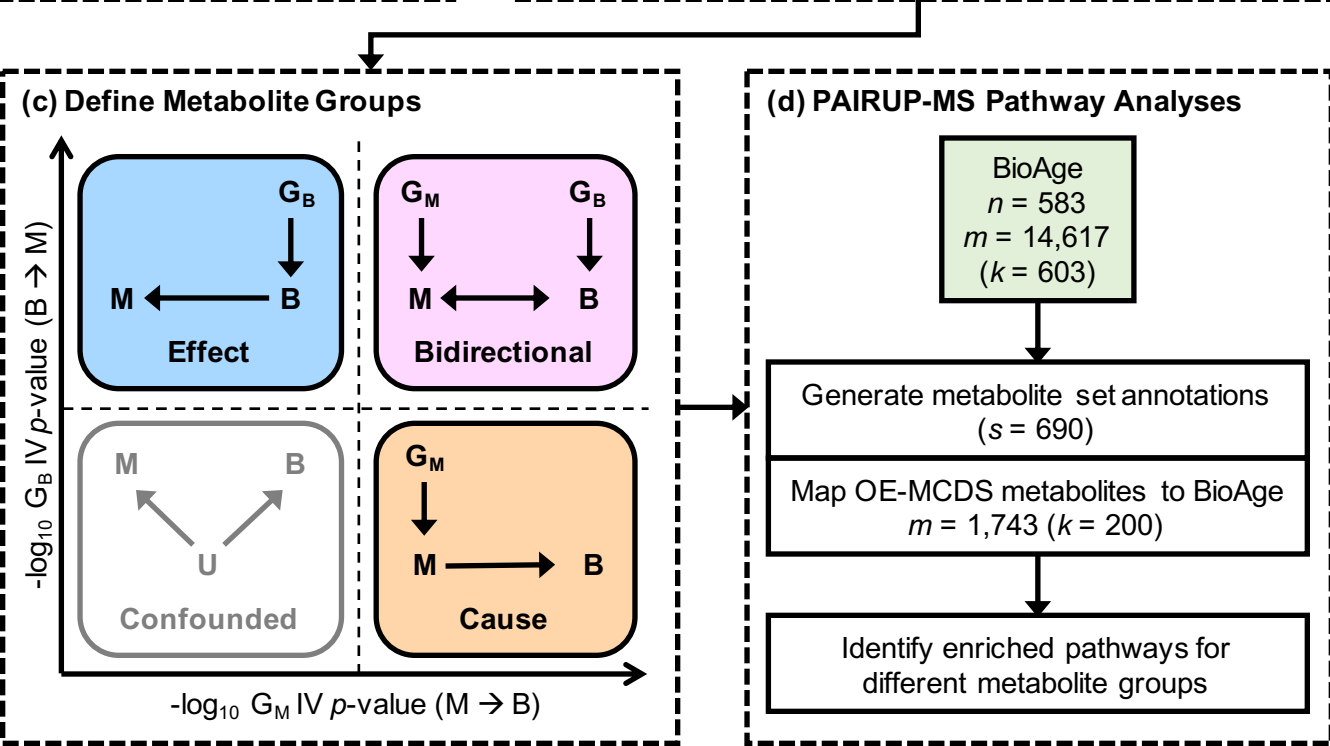
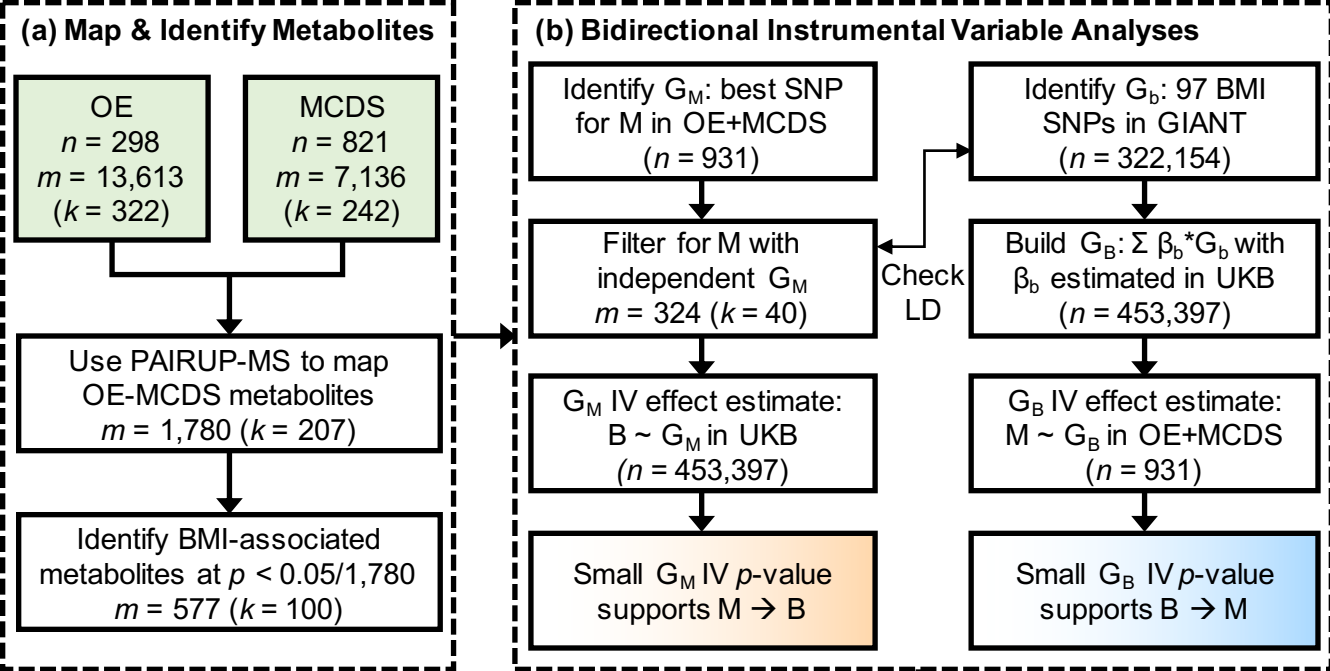
541 significant loci ( $p < 1 \times 10^{-15}$ ) are annotated with nearest genes (within 5kb), along with the best  
542 associated known metabolites if applicable.

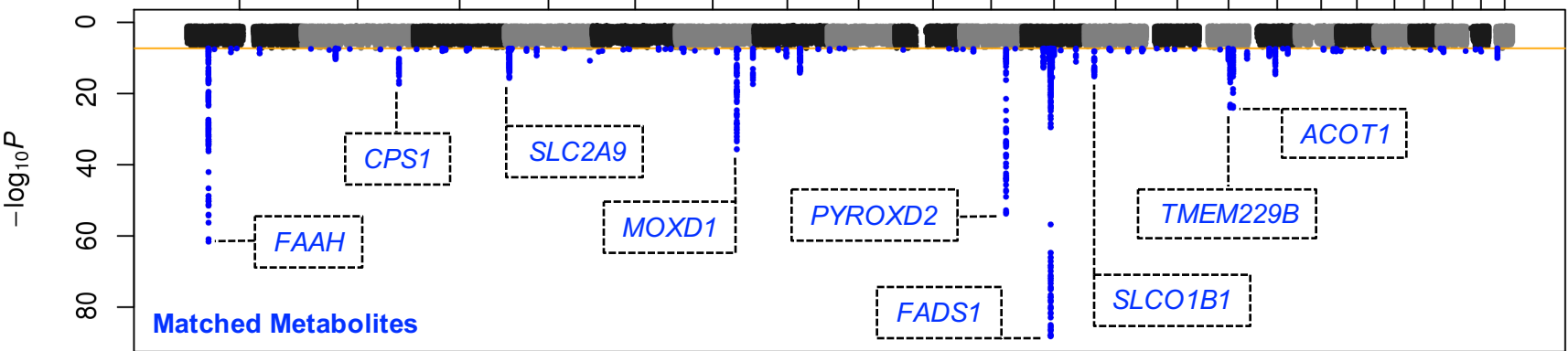
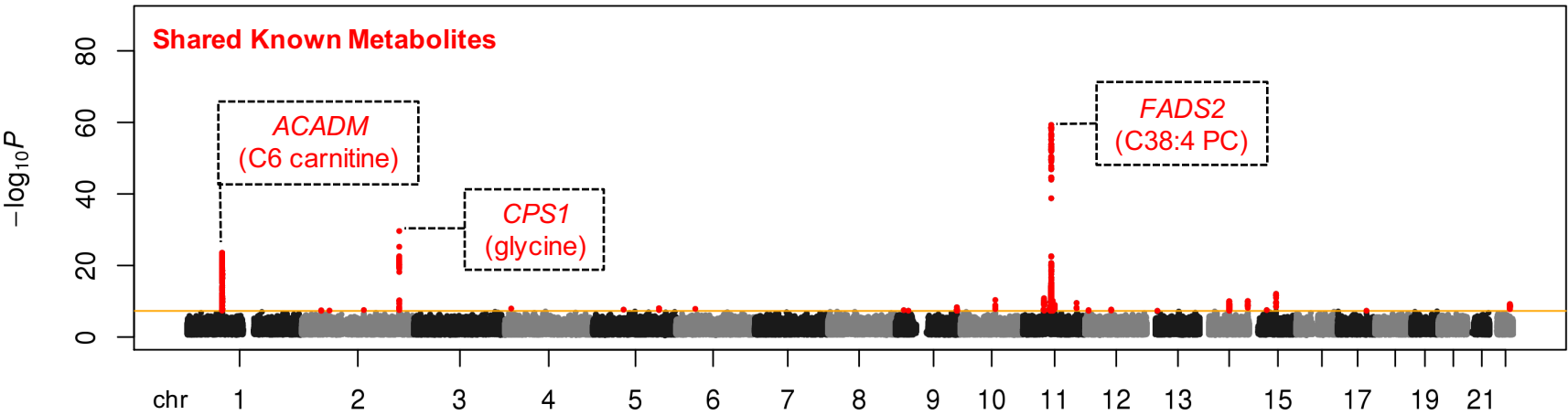
543

544 **Figure 3. Classifying BMI-associated metabolites using IV effect estimate  $p$ -values for  $G_M$**   
545 **(metabolite-to-BMI direction, x-axis) and  $G_B$  (BMI-to-metabolite direction, y-axis).** Top and  
546 bottom quartile cutoffs along each axis are shown as dashed lines. Shared known metabolites in  
547 “cause” (orange), “effect” (blue), and “bidirectional” (pink) regions are labeled with their names.

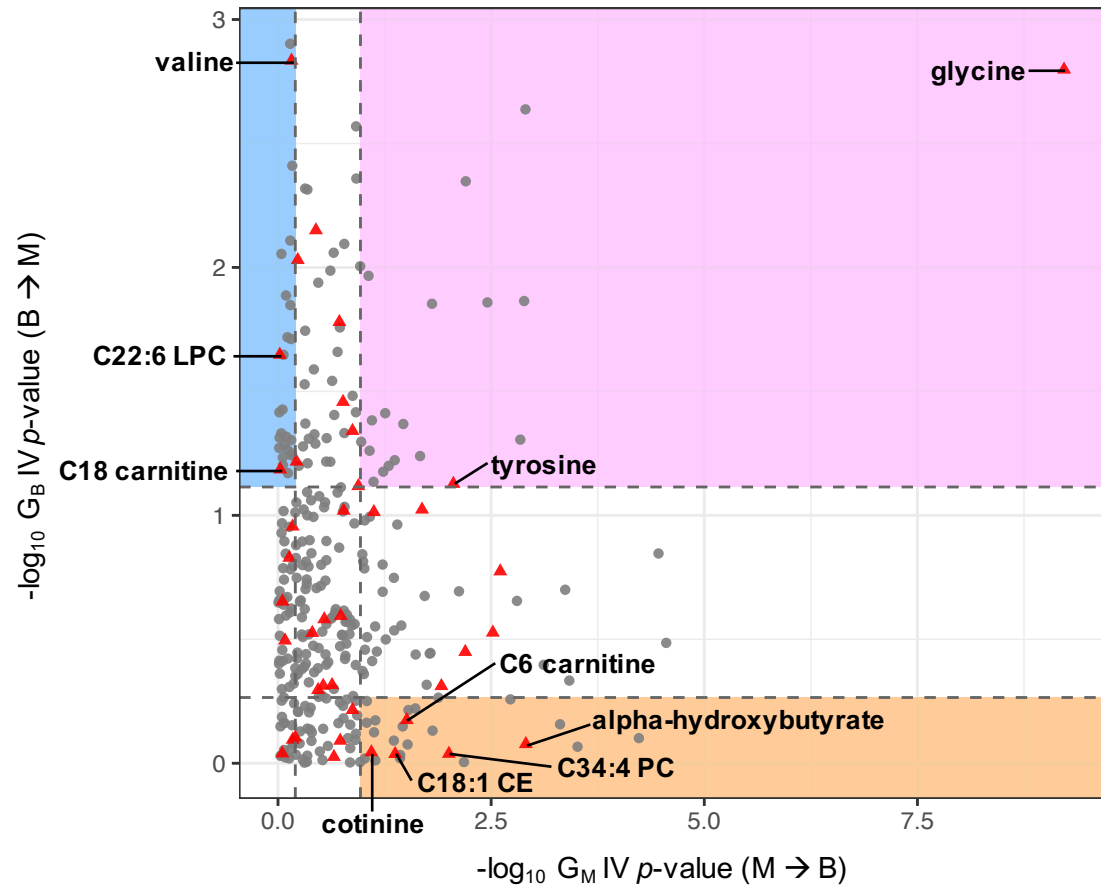
548

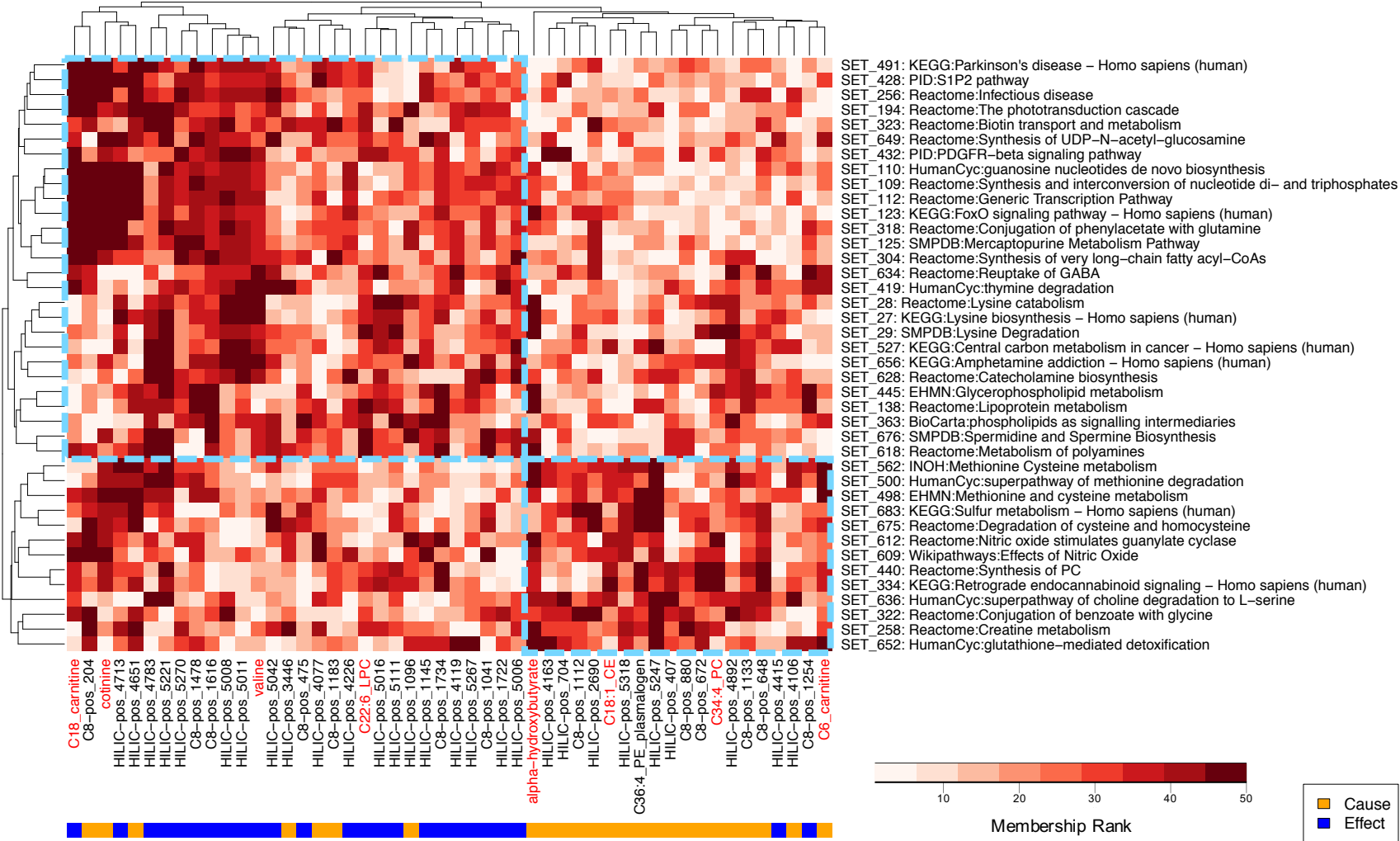
549 **Figure 4. Clustered heat map of cause and effect metabolites’ memberships in metabolite**  
550 **sets prioritized by pathway analysis.** Hierarchical clustering was performed using membership  
551 scores in the BioAge-based metabolite set annotations. Each column is a shared known (red  
552 label) or matched (black label) metabolite from the cause (yellow bar) or effect (blue bar)  
553 metabolite group. Each row is a significant ( $p < 0.05$ ) metabolite set in pathway analysis in either  
554 the cause (yellow bar) or effect (blue bar) direction (with representative pathway name shown in  
555 label; see **Supplementary Table 4** for full pathway list). Larger number in membership rank  
556 (darker red) indicates higher membership score. Dashed light blue boxes highlight the two major  
557 cause and effect clusters according to the clustering dendrograms.





Metabolite Type    ▲ Shared Known    ● Matched





**Table 1. BMI-associated known metabolites classified into cause, effect, or bidirectional groups using the bidirectional IV effect estimate  $p$ -values.**  $Y \sim X$ , regression of  $Y$  on  $X$ ; B, BMI; M, metabolite; covariate adjustment for B and M as described in Methods; SNP, hg19 chromosome:position is shown;  $\beta$ , effect size estimate; EA, effect allele (i.e. metabolite level-increasing allele). Nominally significant  $p$ -values ( $< 0.05$ ) for IV effect estimates are in bold italic.

Group	Metabolite	Observational Association		Metabolite Instrument			$G_M$ IV Estimate ( $M \rightarrow B$ )		$G_B$ IV Estimate ( $B \rightarrow M$ )	
		$\beta$ B ~ M	$P$ B ~ M	SNP	EA	$P$ M ~ $G_M$	$\beta$	$P$	$\beta$	$P$
Cause	alpha-hydroxybutyrate	0.235	7.09E-13	11:119745598	T	4.82E-07	-0.040	<b><i>1.24E-03</i></b>	-0.036	8.36E-01
	C34:4 PC	0.153	3.45E-06	3:182171263	A	1.42E-07	-0.027	<b><i>9.90E-03</i></b>	0.019	9.15E-01
	C6 carnitine	0.207	2.66E-10	1:76224010	C	2.88E-24	-0.010	<b><i>3.09E-02</i></b>	-0.076	6.69E-01
	C18:1 CE	-0.215	5.54E-11	20:38984849	T	3.44E-07	-0.016	<b><i>4.23E-02</i></b>	-0.018	9.17E-01
	cotinine	-0.169	3.10E-07	10:123918365	C	1.07E-06	0.012	8.03E-02	-0.022	9.02E-01
Effect	valine	0.445	1.20E-47	14:32724292	A	7.01E-08	-0.003	6.94E-01	0.708	<b><i>1.46E-03</i></b>
	C22:6 LPC	-0.180	4.38E-08	18:71068347	A	3.26E-07	-0.001	9.52E-01	-0.450	<b><i>2.25E-02</i></b>
	C18 carnitine	-0.193	4.65E-09	6:110760008	A	1.27E-07	0.001	9.41E-01	-0.351	<b><i>6.53E-02</i></b>
Bidirectional	glycine	-0.308	7.55E-22	2:211540507	A	2.35E-30	0.030	<b><i>6.03E-10</i></b>	-0.712	<b><i>1.59E-03</i></b>
	tyrosine	0.376	6.22E-33	6:111477887	C	1.12E-07	-0.022	<b><i>8.77E-03</i></b>	0.334	7.46E-02