

1

2 **Multi-ancestry GWAS of diarrhea during acute**
3 **SARS-CoV2 infection identifies multiple novel loci and**
4 **contrasting etiological roles of irritable bowel**
5 **syndrome subtypes**

6

7 Ninad S. Chaudhary¹, Catherine H. Weldon¹, Priyanka Nandakumar¹,
8 Janie F. Shelton², 23andMe Research Team¹, Michael V. Holmes^{1*}, Stella
9 Aslibekyan^{1*}

10

11

12

13 Affiliations: ¹23andMe, Inc., Sunnyvale, CA; ²Bristol Myers Squibb, Inc.

14

15 * joint senior authors

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35 ABSTRACT

36

37 A substantial proportion of acute SARS-CoV2 infection cases exhibit gastrointestinal
38 symptoms, yet the genetic determinants of these extrapulmonary manifestations are
39 poorly understood. Using survey data from 239,866 individuals who tested positively for
40 SARS-CoV2, we conducted a multi-ancestry GWAS of 80,289 cases of diarrhea
41 occurring during acute COVID-19 infection (33.5%). Six loci (*CYP7A1*, *LZFTL1* -
42 *CCR9*, *TEME182*, *NALCN*, *LFNG*, *GCKR*) met genome-wide significance in a
43 trans-ancestral analysis. The top significant GWAS hit mapped to the *CYP7A1* locus,
44 which plays an etiologic role in bile acid metabolism and is in high LD ($r^2= 0.93$) with the
45 *SDCBP* gene, which was previously implicated in antigen processing and presentation
46 in the COVID-19 context. Another association was observed with variants in the
47 *LZTFL1-CCR9* region, which is a known locus for COVID-19 susceptibility and severity.
48 PheWAS showed a shared association across three of the six SNPs with irritable bowel
49 syndrome (IBS) and its subtypes. Mendelian randomization showed that genetic liability
50 to IBS-diarrhea increased (OR=1.40,95%,CI[1.33-1.47]), and liability to IBS-constipation
51 decreased (OR=0.86, 95%CI[0.79-0.94]) the relative odds of experiencing COVID-19+
52 diarrhea. Our genetic findings provide etiological insights into the extrapulmonary
53 manifestations of acute SARS-CoV2 infection.

54

55

56

57

58

59

60

61

62 INTRODUCTION

63 Diarrhea is a common extrapulmonary gastrointestinal (GI) symptom of acute
64 COVID-19. Retrospective studies estimate that the prevalence of COVID-19 related
65 diarrhea (defined herein as COVID-19+ diarrhea) varies between 7% and 18% for
66 population based and hospital based cohorts, respectively ¹⁻³. There is growing
67 evidence that SARS-CoV2 tropism (the ability of a virus to infect multiple cell types) in
68 gastrointestinal tissues may lead to alteration of gut microbiota ^{4,5} and persistence of
69 virus in the gastrointestinal system, which is associated with a higher risk of post-acute
70 sequelae such as long COVID⁶. SARS-CoV2 RNA can be detected in fecal samples of
71 COVID-19+ patients up to 4 months after acute infection, with fecal viral RNA detectable
72 for longer durations among individuals with COVID-19+ diarrhea as compared to those
73 without diarrhea⁷. Yet, the gastrointestinal symptoms of SARS-CoV2 positive individuals
74 have not been studied in GWAS.

75 The SARS-CoV virus enters the body via angiotensin-converting enzyme 2
76 (ACE2) and co-receptor transmembrane protease serine 2 (TMPRSS2)⁸ receptors that
77 are predominantly expressed in the lung. ACE2 and TMPRSS2 receptors are also
78 expressed in gastric mucosal cells, enterocytes, and colonocytes⁹, making the GI tract a
79 potential extrapulmonary site of SARS-CoV2 infection. SARS-CoV2 RNA and viral
80 proteins were detected in epithelial cells from intestinal biopsies of COVID-19 patients
81 with acute infection¹⁰. In postmortem data from 13 patients, SARS-CoV2 subgenomic
82 RNA was detected in the small intestinal tissues from eight patients, indicative of viral
83 replication¹¹. It is not clear whether GI symptoms are the direct consequence of
84 SARS-CoV2 infection of intestinal cells or due to a systemic immune-inflammatory

85 response. GWAS can provide insights into the genetic architecture of disease and shed
86 light on potential etiological mechanisms. To date, no GWAS studies of COVID-19+
87 diarrhea have been conducted. We sought to bridge this gap by conducting a GWAS of
88 self-reported diarrhea symptoms among 239,866 COVID-19 test positive 23ndMe
89 research participants.

90

91 RESULTS

92 Of the individuals who self-reported testing positive for SARS-CoV2, 33%
93 reported diarrhea as a symptom (N=80,289 cases of COVID-19+ diarrhea / 239,866
94 acute COVID-19 cases who tested positive). Among 80,289 COVID-19+ diarrhea cases,
95 the vast majority reported diarrhea symptoms at baseline (N=79,361 cases) with a few
96 individuals reporting diarrhea at 1 month (N=313 cases), 2 months (N=346 cases), or 3
97 months of follow up (N=269 cases). Women accounted for the majority of the survey
98 participants (71.1% of COVID-19+ diarrhea cases and 65.4% of controls), and were
99 more likely to experience COVID-19+ diarrhea than men (adjusted OR[95%CI]
100 =1.29[1.27,1.31]). Compared to those with COVID-19+ who didn't experience diarrhea,
101 individuals with COVID-19+ diarrhea were younger, more likely to be of non-European
102 ancestry, and to have diabetes, depression, and high triglyceride levels (**Table 1**).
103 Individuals of East Asian (adjusted OR[95%CI]=0.85[0.79,0.94]) or African American
104 (adjusted OR [95%CI]=0.92[0.88,96]) ancestry had lower relative odds of reporting
105 diarrhea compared to those of European ancestry, whereas Latinx individuals were at
106 higher relative odds of reporting diarrhea (adjusted OR[95%CI]=1.21[1.09,1.14]). Those
107 reporting COVID-19+ diarrhea were more than twice as likely to be hospitalized during

108 the acute infection (6.6% in COVID-19+ diarrhea cases vs 3.0% in COVID-19+ controls;
109 corresponding to an adjusted relative odds of 2.56; 95%CI: 2.46, 2.57) and were more
110 likely to have severe respiratory disease (10.2% in COVID-19+ diarrhea cases vs 4.8%
111 in COVID-19+ diarrhea controls; adjusted relative odds of 2.38; 95%CI: 2.30, 2.46) after
112 adjusting for age, sex, and ancestry. Furthermore, those with COVID-19+ diarrhea had
113 twice the odds of subsequently experiencing long COVID (43% in COVID-19+ diarrhea
114 cases vs 23% in COVID-19+ diarrhea controls; adjusted relative odds of 2.53; 95%CI:
115 2.46, 2.61) and of long COVID impacting daily living activities (16% of COVID-19+
116 diarrhea cases vs 9.6% of COVID-19+ diarrhea controls; adjusted relative odds of 2.27;
117 95%CI: 2.21, 2.34).

118

119 **GWAS:** We conducted GWAS within each of European, African American, Latinx,
120 East Asian, and South Asian ancestry groups, and then meta-analyzed across
121 ancestries using fixed effect modeling. Individuals of European ancestry contributed
122 75% of the sample size to multi-ancestry analysis, with African-Americans contributing
123 4%, Latinx 18%, East Asians 2%, and South Asians 1%. We identified six distinct loci
124 associated with COVID-19+ diarrhea in multi-ancestry analysis (**Table 2, Figure 1**). The
125 associations of index variants were statistically significant at baseline but not at the
126 following time points (1 and 3 months after onset of acute infection) likely owing to lack
127 of statistical power (**Supplementary Table 1**). The top statistically significant index
128 variant at chr8q12.1 (rs10504255) is situated in the intergenic region of genes *UBXN2B*
129 and *CYP7A1*, located 4kb upstream of the *CYP7A1* gene. rs10504255 (A/G with G
130 being the effect allele) was associated with COVID-19+ diarrhea with OR[95% CI] =

131 0.94[0.92,0.95], $p = 2.6 \times 10^{-16}$. The credible set from multi-ancestry analysis contained 16
132 variants covering a 100.4-kilobase(kb) region (**Supplementary Figure 1**). The
133 association was primarily driven by individuals of European ancestry ($p = 1.08 \times 10^{-14}$)
134 as compared to other ancestries (**Supplementary Table 2**). In the analysis of
135 expression quantitative trait loci (eQTL), rs10504255 was found to be in high LD with a
136 variant (rs9297994) associated with *CYP7A1* expression in the thyroid gland ($p =$
137 8.1×10^{-9}) and a variant (rs8192870) associated with expression of a nearby gene
138 (*SDCBP*, $r^2=0.94$, $p=9.9 \times 10^{-10}$) in left ventricular myocardium (**Supplementary Table**
139 **3**)¹². Previous studies have shown *SDCBP* expression to correlate with *HLA-DPB1*
140 expression in normal lung tissue ¹³. The variant rs35044562 (alleles A/G with G being
141 the effect allele) on chr3p21.31 lies in the intergenic region of *LZTFL1* and *CCR9*
142 (**Supplementary Figure 2**). *LZTFL1* is widely expressed in ciliated epithelial cells in
143 lungs ¹⁴. Additional genes in this regulatory locus include *CCR9*, primarily expressed in
144 immune cells, and *SLC6A20* in the gastrointestinal tract.

145

146 The variants, rs75683620, at chr2q12.1 (*TMEM182*, alleles A/G with G being the
147 effect allele, $p=5.0 \times 10^{-09}$, **Supplementary Figure 3**), rs536843010 at chr13q33.1
148 (*NALCN*, alleles A/C with C being the effect allele, $p=9.8 \times 10^{-09}$, **Supplementary Figure**
149 **4**), and rs13245319 at chr7p22.3 (*GRIFIN--[]-LFNG*, alleles C/T with T being the effect
150 allele, $p = 1.79 \times 10^{-08}$, **Supplementary Figure 5**) were not associated with functional
151 effects in eQTL or pQTL analysis. These variants are relatively rare in the studied
152 populations, except for rs75683620 (MAF= 0.042) among individuals of African

153 American ancestry. These variants were monomorphic among East and South Asians,
154 meaning that data from these populations did not contribute to the meta-analysis.

155

156 We also identified an association at chr2p23.3 (*GCKR*, rs1260326, alleles C/T
157 with T being the effect allele, **Supplementary Figure 6**) that was genome-wide
158 significant in the meta-analysis ($p = 1.98 \times 10^{-08}$) and in the European population ($p =$
159 2.0×10^{-08}). rs1260326 is a missense variant that is also in high LD with eQTLs for other
160 nearby genes across multiple tissues, including kidney tubules (eQTL gene = *NRBP1*,
161 $p = 5.3 \times 10^{-07}$), CD4+ T cells (eQTL gene = *NRBP1*, $p = 2.8 \times 10^{-07}$), and liver (eQTL
162 gene = *C2orf16*, $p = 4.9 \times 10^{-23}$) (**Supplementary Table 3d**)^{15–18}. rs1260326 is also located
163 within 500kb with $r^2 > 0.8$ of multiple pQTLs based on data from blood plasma
164 (**Supplementary Table 3b**)^{19,20}. As previously reported, rs1260326 is a pleiotropic
165 variant associated with multiple traits²¹. Similarly to other variants, most of the support
166 for this association also comes from the European population.

167

168 **Association with COVID-19 measures:** To explore the specificity of these
169 genetic variants with COVID-19+ diarrhea vs COVID-19 severity, we examined the
170 association of the lead variants with COVID-19 susceptibility and COVID-19 severity
171 measures. The variant, rs3504462 (A/G with G as the risk allele), on chr3p21.31 was
172 associated with COVID-19 test positivity ($p = 1.1 \times 10^{-05}$) and acute COVID-19 infection
173 leading to hospitalization ($p = 1.1 \times 10^{-56}$) or severe respiratory disease ($p = 7.1 \times 10^{-69}$). In
174 contrast, the other variants were not associated with any of these COVID-19 measures
175 (**Supplementary Table 4**). Thus, with the exception of the chr3p21.31 locus, the

176 genetic signals of COVID-19+ diarrhea are likely to drive their effects through
177 mechanisms unrelated to COVID-19 severity.

178

179 **PheWAS:** To characterize these loci, we performed a phenome-wide association
180 study (PheWAS) across 1,482 phenotypes available in the 23andMe, Inc. database in
181 the population of European ancestry. PheWAS results for variants are presented in
182 **Supplementary Table 5**. Among phenotypes associated with rs10504255, the
183 strongest associations were with high cholesterol ($p = 7.53 \times 10^{-172}$) and IBS-D ($p = 5.31$
184 $\times 10^{-134}$). Other associated phenotypes were cardiometabolic traits such as type 2
185 diabetes, high blood pressure and statin use (**Supplementary Table 5a**). In addition to
186 the measures of COVID-19 susceptibility and severity, the rs35044562 variant was also
187 associated with autoimmune conditions (Hashimoto's disease, celiac disease)
188 (**Supplementary Table 5b**). PheWAS for rs1260326 identified associations with lipid
189 profile, allergic conditions, and blood glucose levels (**Supplementary Table 5c**). The
190 main associations for rs1324319 were IBS, kidney stones, and obesity (**Supplementary**
191 **Table 5d**). No traits were associated with rs75683620 or rs536843010 at the
192 Bonferroni-corrected p-value threshold.

193

194 Focusing specifically on gastrointestinal disorders among outcomes constituting
195 the PheWAS analysis, and orienting effect alleles to a higher risk of COVID-19+
196 diarrhea, four SNPs (rs10504255, rs1260326, rs35044562 and rs13245319)
197 demonstrated associations with GI traits in addition to COVID19+ diarrhea (**Figure 2**),
198 including IBD (both ulcerative colitis and Crohn's disease), celiac disease, lactose

199 intolerance and IBS. Directionally consistent associations with higher risks of IBS (for
200 rs10504255, rs1260326 and rs13245319) and IBS-D (rs10504255, rs1260326) were
201 identified whereas a directionally opposite association between COVID-19+ diarrhea
202 and IBS-C (rs10504255) was found (**Supplementary Figure 7**).

203

204 **Mendelian randomization:** Given that a common association on PheWAS
205 across three of the six lead variants was an association with IBS and IBS subtypes, and
206 to more fully characterize the translational relevance, we investigated the potential
207 causal relationship between genetic liability to IBS-C and IBS-D and COVID-19+
208 diarrhea through Mendelian randomization (MR). Using a genetic instrument consisting
209 of 180 SNPs for IBS-D, we identified strong evidence of a potential causal effect of
210 liability to IBS-D and risk of COVID-19+ diarrhea (OR=1.40,95%,CI[1.33-1.47] from
211 random-effects IVW modeling). Steiger filtering removed three SNPs with minimal
212 impact to the IVW MR estimate (OR=1.39,95% CI[1.32-1.46]). These causal estimates
213 persisted or strengthened on robust MR approaches: weighted median provided
214 findings that were largely similar to IVW MR (OR=1.37,95%,CI[1.28-1.47]), with the
215 estimate from MR Egger yielding a stronger predicted causal effect
216 (OR=2.02,95%,CI[1.76-2.32]) (**Figure 3, Supplementary Figure 8**). The intercept from
217 MR Egger regression was statistically significantly different from zero [β (SE) =
218 -0.121(0.002); $p= 1.74 \times 10^{-07}$], indicating presence of directional pleiotropy
219 (**Supplementary Table 6**).

220 The MR estimates for IBS-C were directionally opposite (0.86, 95%CI[0.79-0.94];
221 IVW modeling) to those of IBS-D. MR estimates remained unchanged after removing

222 one SNP following Steiger filtering (OR=0.89,95% CI[0.82-0.96]). Robust MR analyses
223 provided consistent findings although the magnitude of the predicted causal estimate
224 was further from the null on MR-Egger; the intercept from MR Egger was $\beta(\text{SE}) =$
225 0.008(0.008) with a p-value of 0.31.

226 Given partial overlap of individuals contributing to GWAS (**Supplementary**
227 **Figure 9**), we conducted a sensitivity analysis using non-overlapping samples which
228 yielded nearly identical findings (**Supplementary Table 7**).

229

230 **DISCUSSION**

231 In this era of widespread documentation of SARS-CoV2 effects on human health,
232 there is an important gap in understanding extrapulmonary symptoms. In this study, we
233 address this gap by utilizing a direct-to-consumer research platform contributing data on
234 80,289 cases of diarrhea during acute SARS-CoV2 infection. We identified three loci
235 (*CYP7A1*, *LZTFL1-CCR9*, and *GCKR*) that have plausible biological mechanisms of
236 action. Except for *LZTFL1-CCR9*, none of these genetic signals was associated with
237 COVID-19 severity, indicating the role of distinct biological pathways. We further
238 observed a consistent association of top variants at the *CYP7A1*, *GCKR*, and *LFNG* loci
239 with IBS and IBS sub-types. Genetic liability towards IBS-D increased the risk of having
240 COVID-19+ diarrhea, while genetic liability for IBS-C reduced this risk. These results
241 highlight the role of genetic predisposition to preexisting comorbidities in the
242 extrapulmonary manifestations of SARS-CoV2 infection.

243

244 The top significant locus in our study, mapped to *CYP7A1*, encodes a member of
245 the cytochrome p450 enzyme family, which has an important role in the bile acid
246 synthesis pathway²². Polymorphisms at this locus are associated with defects in bile
247 acid synthesis, affecting enzymatic activity of cholesterol 7- α hydroxylase and
248 resulting in bile acid diarrhea. On eQTL analysis, we observed rs10504255 to be in high
249 LD with a variant (rs8192870) associated with expression of a nearby gene (*SDCBP*,
250 $r^2=0.93, p=9.9 \times 10^{-10}$) in the left ventricular myocardium¹². *SDCBP* is also expressed in
251 the intestine, as well as in lung cancers including adenocarcinoma and small cell
252 carcinoma¹². A recent single-cell RNA sequencing study observed that *SDCBP* likely
253 plays a role in antigen processing and presentation in bronchial epithelial cells from
254 COVID-19 patients,¹³ with *SDCBP* expression correlating with *HLA-DPB1* expression in
255 normal lung tissues¹³. *HLA-DPB1* represents a critical immune mechanism as it is
256 expressed on antigen presenting cells that participate in eliciting an immune response
257 to foreign viral peptides. Overall, these findings suggest a possible role of bile acid
258 synthesis pathways and/or the immune system in COVID-19+ diarrhea.

259

260 *LZTFL1*, *SLC6A20* and *CCR9* are part of the chemokine receptor gene cluster at
261 chr3p21, previously identified as COVID-19 susceptibility and severity loci^{23–26}. *LZTFL1*
262 at this locus likely regulates viral response pathways by inhibiting the transcription
263 factors that reduce levels of ACE2 and TMPRSS¹⁴. *SLC6A20* is extensively expressed
264 in the gastrointestinal tract, where it forms a complex with ACE2 receptors, facilitating
265 viral entry²⁴. *CCR9* is expressed in T-lymphocytes of the small intestine and colon,
266 where it regulates chemokines and eosinophil recruitment^{23,27}. Collectively, these

267 mechanisms may contribute to pathophysiology of diarrheal disease during acute
268 SARS-CoV2 infection.

269

270 PheWAS characterization of top GWAS loci identified directionally opposite
271 signals with IBS subtypes. Genetic liability to IBS-D increased, and liability to IBS-C
272 decreased the probability of experiencing diarrhea during the acute phase of
273 SARS-CoV2 infection. These effects could be driven by genetic liability to frequency of
274 bowel motions, modifying the manifestation of diarrhea in the context of an acute
275 infection affecting the GI tract. An alternative explanation might be that individuals at
276 greater liability to IBS-D have altered bowel characteristics that make them more
277 susceptible to GI infection and diarrhea in the context of acute COVID-19 infection. For
278 example, studies have demonstrated that individuals with IBS-D have an altered gut
279 microbiome which could play a role in susceptibility to COVID-19+ diarrhea ^{28,29}.

280

281 The remaining significantly associated genetic loci do not have direct evidence of
282 involvement in the gastrointestinal tract or SARS-CoV2 infection. *GCKR* is a pleiotropic
283 locus associated with C-reactive protein, fasting plasma glucose levels, and blood cell
284 traits ^{21,30}. *GCKR*-mediated effects are driven via regulation of glucokinase enzymes that
285 control the first step of glycolysis. The gene *TMEM182* (rs75683620) at chr2q12.1 has
286 been associated with central obesity and systolic blood pressure in Asian populations
287 ^{31,32}. The gene is expressed in heart tissue and regulates its effects via tumor necrosis
288 factor-alpha ³². The *NALCN* gene (rs536843010) at chr13q33.1 contributes to
289 physiological processes in the neuromuscular junctions by maintaining resting

290 membrane potential via voltage independent, nonselective cation channels³³. It
291 accordingly is involved in control of muscular activity, respiration, and circadian rhythms
292^{33,34}. The function of the LFNG gene (rs13245319) is regulated via the notch signaling
293 pathway³⁵. The notch pathway maintains the homeostasis of multiple tissues and thus
294 plays a role in cancerous growth, including colorectal adenocarcinoma^{36,37}. These
295 associations require independent validation and additional functional investigations to
296 identify any biological relevance to acute COVID-19+ diarrhea.

297

298 Our study has several limitations. European participants contributed 75% of our
299 sample size, so the statistical power to detect genetic associations was driven largely by
300 this group. While the large sample size is an important strength for our study, replication
301 of the GWAS and MR findings by other studies would establish reproducibility and
302 strengthen the findings. Diarrhea cases reported 2 to 3 months after acute infection
303 might be misattributed to COVID-19 as opposed to another source. However, the vast
304 majority (99%) of our cases of COVID-19+ diarrhea occurred at the same time as the
305 diagnosis of COVID-19. Furthermore, individuals with genetic liability to IBS-D may
306 have misattributed diarrhea at the time of acute COVID-19 to the infection as opposed
307 to their underlying IBS-D. Refuting this hypothesis is our observation that individuals
308 with IBS-C had a lower relative odds of experiencing COVID-19+ diarrhea, which
309 collectively may be indicative of an acute infection operating on top of a background
310 genetic liability to IBS-D or IBS-C, with liability to either IBS-D or IBS-C influencing the
311 probability of experiencing COVID-19+ diarrhea in a potentially etiological way. Finally,
312 the MR-Egger intercept for IBS-D indicated evidence of directional pleiotropy. However,

313 the pleiotropy-corrected point estimate from MR-Egger was larger (OR=2.02) than that
314 derived from IVW (OR=1.40), making the latter a conservative estimate.

315

316 In conclusion, the first GWAS of COVID-19+ diarrhea in an ancestrally diverse,
317 large-scale population-based cohort identifies biologically relevant signals associated
318 with bile acid synthesis and immune function, including antigen presenting and cytokine
319 signaling. By virtue of limiting our study to COVID-19+ individuals, we infer and
320 substantiate empirically that all but one of the identified loci are specific to COVID-19+
321 diarrhea rather than markers of susceptibility to SARS-CoV2 infection. Finally, we
322 provide causal evidence in support of the etiological role of liability to chronic intestinal
323 disorders and gastrointestinal symptoms during acute infection.

324 **METHODS**

325 Study Population

326 Participants older than 18 years old were recruited for 23andMe COVID-19 study
327 from April 2020 using email-based surveys. The surveys were distributed to 6.7 million
328 individuals who provided informed consent and volunteered to participate in the
329 research online, under a protocol approved by the external AAHRPP-accredited IRB,
330 Ethical & Independent (E&I) Review Services. As of 2022, E&I Review Services is part
331 of Salus IRB. Recruitment was initially geo-targeted to capture cases as the outbreak
332 spread across the U.S and was continued later on to recruit additional participants. We
333 included only those participants who had responded to COVID-19 survey and provided
334 symptom information by August, 2023 for this study. Details on diagnosis, testing, and

335 symptoms of COVID-19, as well as markers of severity and relevant comorbid
336 conditions were collected via one baseline and three follow-up surveys administered
337 three months apart. Full details of the data collection procedures for this study have
338 been described previously³⁸.

339

340 Phenotype

341 Using the survey information, we defined the diarrhea outcome among
342 individuals who self-reported testing positive for COVID-19. Participants were asked the
343 question whether they experienced any of the following symptoms to which they could
344 select as many as needed from the following list of responses: ‘muscle or body aches/
345 fatigue/ dry cough/ sore throat/ coughing up of sputum or phlegm (productive cough)/
346 loss of smell or taste/ chills/ difficulty breathing or shortness of breath/ pressure or
347 tightness in upper chest/ diarrhea/ nausea or vomiting/ sneezing/ loss of appetite/ runny
348 nose/ headache/ intensely red or watery eyes’. Participants who reported experiencing
349 diarrhea symptoms at baseline or follow-up surveys were identified as cases and those
350 who did not experience diarrhea were controls. The analytical dataset included
351 unrelated individuals with non-missing information on COVID-19+ diarrhea and
352 covariates included in GWAS analyses (N=239,866).

353

354 Genotyping

355 DNA extraction and genotyping were performed on saliva samples by Clinical
356 Laboratory Improvement Amendments-certified and College of American
357 Pathologists-accredited clinical laboratories of Laboratory Corporation of America.

358 Samples were genotyped across the five genotyping platforms and imputed using three
359 combined independent reference panels: the publicly available Human Reference
360 Consortium (HRC), and UK BioBank (UKBB) 200K Whole Exome Sequencing (WES)
361 reference panels and the 23andMe reference panel, which was built by 23andMe using
362 internal and external cohorts. Each genotyping platform was imputed and phased
363 separately. The final genotyped variants included 1,469,237 variants and the final
364 imputation panel included a total of 99,675,338 variants (90,582,19 SNPs and
365 9,093,144 indels). The variant quality control statistics were computed independently
366 with each phasing panel and genotyping platform. We removed variants with low
367 imputation quality ($r^2 < 0.5$ averaged across batches or a minimum $r^2 < 0.3$) or
368 evidence of differences in effects across batches. For genotyped variants, we removed
369 variants only present on our V1 or V2 arrays (due to small sample size) that failed a
370 Mendelian transmission test in trios ($P < 10^{-20}$), failed a Hardy–Weinberg test in
371 individuals of European ancestry ($P < 10^{-20}$), failed a batch effect test (ANOVA $P < 10^{-50}$)
372 or had a call rate $< 90\%$. For imputed variants in HRC panel, following filters were used:
373 singletons were excluded, multi-allelics were split into bi-allelic variants (with bcftools),
374 variants with $> 20\%$ missingness were removed, variants with minor allele count $= 0$
375 were removed, and variants with inbreeding coefficient < -0.3 (high heterozygosity) were
376 removed.

377

378 Ancestry Classification

379 Ancestries in the 23andMe database are determined using a classifier algorithm based
380 on analysis of local ancestry³⁹. Phased genotyped data were first partitioned into

381 windows of about 300 SNPs and a support vector machine (SVM) approach was
382 applied within each window to classify individual haplotypes into one of 45 worldwide
383 reference populations. The SVM classifications are then fed into a hidden Markov model
384 (HMM) that accounts for switch errors and incorrect assignments, and gives
385 probabilities for each reference population in each window. Finally, we used simulated
386 admixed individuals to recalibrate the HMM probabilities so that the reported
387 assignments are consistent with the simulated admixture proportions. We aggregated
388 the probabilities of the 45 reference populations into six main ancestries (European,
389 African-American, Latinx, East Asian, South Asian, Middle Eastern) using a
390 predetermined threshold³⁹. African Americans and Latinx were admixed with broadly
391 varying contributions from Europe, Africa and the Americas. No single threshold of
392 genome-wide ancestry could effectively discriminate between African Americans and
393 Latinx. However, the distributions of the length of segments of European, African and
394 American ancestry are very different between African Americans and Latinx, because of
395 distinct admixture timing between the three ancestral populations in the two ethnic
396 groups. Therefore, we trained a logistic classifier that took the participant's length
397 histogram of segments of African, European and American ancestry, and predicted
398 whether the customer is likely African American or Latinx.

399

400

401 GWAS Analysis

402 To obtain unrelated participants for our GWAS analyses, individuals were
403 included such that no two individuals shared more than 700 cM of DNA identical by

404 descent. We excluded approximately 1.80% of the sample to obtain such a set of
405 unrelated individuals. If a case and control were identified as having at least 700cM of
406 DNA IBD, we retained the case from the sample. We conducted association analysis
407 using logistic regression, assuming additive allelic effects and adjusting for age,
408 age-squared, sex, sex-age interaction, genotyping platform variables and ten principal
409 components to account for residual population structure. We combined the GWAS
410 summary statistics from both genotyped and imputed data. When choosing between
411 imputed and genotyped GWAS results, we favored the imputed result, unless the
412 imputed variant was unavailable or failed quality control. The summary statistics were
413 adjusted for inflation using genomic control when the inflation factor was estimated to be
414 greater than one ($\lambda = 1.029, 1.037, 1.024, 1.042$ and 1.071 within the European, Latinx,
415 African American, East Asian and South Asian ancestry GWAS, respectively). We
416 defined the region boundaries by identifying all SNPs with $P < 10^{-5}$ within the vicinity of a
417 genome-wide significance association and then grouping these regions into intervals so
418 that no two regions were separated by less than 250 kb. We considered the SNP with
419 the smallest P value within each interval to be the **index variant**. Within each region,
420 we calculated a credible set of variants using the method of Maller et al 2021⁴⁰.

421

422 We conducted the GWAS analysis separately in five population cohorts
423 (European, Latinx, African-Americans, East Asian, and South Asian ancestry). We then
424 meta-analyzed the GWAS summary statistics of these populations using an inverse
425 variance fixed effect model. For this approach, we included variants that had at least 1%
426 minor allele frequency in the pooled sample and minor allele count > 30 within each

427 subpopulation. The resulting meta-analyses were also adjusted for inflation ($\lambda = 1.001$)
428 using genomic control.

429

430 PheWAS Analysis

431 We conducted phenome-wide associations (PheWAS) on the index variants from loci
432 that were statistically significant on multi-ancestry analysis. The PheWAS analysis was
433 limited to data from participants of European ancestry. We used data on 1,482
434 phenotypes that were available in the 23andMe research database. Data on these
435 phenotypes were collected using online survey-based questionnaires completed by the
436 participants at the time of recruitment in 23andMe genetic database. PheWAS analysis
437 was performed, adjusting the association between each phenotype and variant of
438 interest for age, sex, and the first five principal components of ancestry. We reported the
439 associations that met the statistical threshold of significance after correcting for multiple
440 testing ($p < 0.05 / (1,482 \times 6) = 5.62 \times 10^{-06}$).

441

442 Mendelian Randomization

443 Because PheWAS highlighted a shared association with irritable bowel syndrome
444 (IBS) and IBS subtypes, we explored whether there was a potential causal role of
445 genetic liability to IBS subtypes and COVID-19+ diarrhea by conducting two-sample
446 Mendelian randomization (MR) using the *twosampleMR* R package⁴¹. The initial MR
447 analysis included overlapping samples of IBS subtypes and COVID-19+ diarrhea in the
448 European population. We focused on the IBS subtypes of diarrhea (IBS-D) and
449 constipation (IBS-C). We first obtained separate genetic instruments for IBS-C and

450 IBS-D from genome-wide significant SNPs for the relevant endpoint that had a minor
451 allele frequency of more than 0.01. We then identified variants with the smallest p-value
452 (index variants) by defining the regions of significant association based on genome-wide
453 significant SNPs as described above. These index variants were included in the genetic
454 instrument of the relevant phenotype (IBS-C or IBS-D). The genetic instrument for
455 IBS-D included 180 SNPs (mean F-statistic per SNP=50.2) and for IBS-C included 52
456 SNPs (mean F-statistic per SNP=43.6). Using the summary statistics of SNPs with
457 IBS-D and IBS-C adjusted for age, sex, and principal components of ancestral structure
458 in the European population, we then fitted random effects inverse variance weighted
459 (IVW) models and obtained MR estimates of IBS subtypes on COVID-19+ diarrhea. To
460 test the veracity of the findings, we used Steiger filtering and robust MR approaches
461 including weighted median, and MR-Egger method^{41–43}. As a sensitivity analysis, we
462 repeated MR analysis between IBS subtypes and COVID-19+ diarrhea after excluding
463 the overlapping samples between them.

464

465 Functional Annotation of GWAS Index Variants

466

467 To perform variant-to-gene mapping, hypotheses of functionally relevant genes
468 are generated by annotating the strongest associations (index variants) with nearby
469 functional variants. The mapping is computed by searching functional variants within
470 500 kb of the index variant with a filter of linkage disequilibrium $r^2 > 0.8$. Functional
471 variants for mapping include coding variants (annotated by the Ensembl Variant Effect
472 Predictor (VEP) v109⁴⁴), eQTLs, and pQTLs. The eQTL annotation resources consist of
473 a comprehensive collection of standardized variants impacting gene expression in
474 various tissues obtained from publicly available datasets^{12,15–18,45–48} and datasets

475 processed by the 23andMe eQTL pipeline . The pQTL annotation resources similarly
476 include a collection of curated protein QTLs from relevant public datasets ^{19,20}.
477 *eQTL discovery*. eQTL calling was performed with one of two versions of 23andMe
478 pipelines, depending on the dataset in question (**Supplementary File**). The first pipeline
479 used FastQTL ⁴⁹ in permutation mode, restricting all tests to variants within a window
480 defined to be 1Mbp up- or downstream of a given gene's transcription start site (TSS)
481 (**Supplementary File Table 1**). Variants tested were single nucleotide polymorphisms
482 with an in-sample MAF $\geq 1\%$, to avoid errors in detection or mapping of larger genetic
483 variants in cross-ancestry comparisons, and models were adjusted for age (if available),
484 sex, probabilistic estimation of expression residuals (PEER) factors ⁵⁰, genetic PCs, and
485 per-dataset covariates. For each gene, the index variant was identified by the minimal
486 permutation p-value. eQTLs were called then on lead variants if they passed a 5% FDR
487 filter using Storey's q-value ⁵¹ methodology. Conditional eQTLs were identified via
488 FastQTL's permutation mode, by using each eQTL as an additional covariate in the
489 model for a given gene. The lead conditional eQTL for all genes were again FDR
490 controlled at 5%, and a maximum of 10 conditional steps were run. Finally, for a set of
491 conditional eQTLs for a given gene, a joint model was fit, and the final eQTL callset
492 consisted of those eQTLs whose joint model test passed a 5% FDR filter. eQTLs were
493 only called for genes classified as one of 'protein_coding', 'miRNA', 'IG_C_gene',
494 'IG_D_gene', 'IG_J_gene', 'IG_V_gene', 'TR_C_gene', 'TR_D_gene', 'TR_J_gene',
495 'TR_V_gene' as defined in GENCODE⁵².

496 The second version of the 23andMe pipeline uses strand-aware RNA-seq
497 quantification, and the eQTLs were called using SusieR package ⁵³ instead of FastQTL,

498 with expression PCs (selected with the elbow method) replacing PEER factors in the
499 modeling, and using the GENCODE v43 gene model (**Supplementary File Table 2**).
500 The pipeline natively generated credible sets with a set probability to contain a SNP
501 tagging the causal variant.

502

503

504 **DATA AVAILABILITY**

505

506 The full set of GWAS summary statistics can be made available to qualified
507 investigators upon request and signing agreement with 23andMe to protect participant
508 confidentiality. The information can be accessed at
509 <https://research.23andme.com/covid19-dataset-access/>

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530 REFERENCES

531

- 532 1. Mao, R. *et al.* Manifestations and prognosis of gastrointestinal and liver involvement in
533 patients with COVID-19: a systematic review and meta-analysis. *Lancet Gastroenterol.*
534 *Hepatol.* **5**, 667–678 (2020).
- 535 2. Parasa, S. *et al.* Prevalence of Gastrointestinal Symptoms and Fecal Viral Shedding in
536 Patients With Coronavirus Disease 2019: A Systematic Review and Meta-analysis. *JAMA*
537 *Netw. Open* **3**, e2011335 (2020).
- 538 3. Sultan, S. *et al.* AGA Institute Rapid Review of the Gastrointestinal and Liver Manifestations
539 of COVID-19, Meta-Analysis of International Data, and Recommendations for the
540 Consultative Management of Patients with COVID-19. *Gastroenterology* **159**, 320-334.e27
541 (2020).
- 542 4. Nardo, A. D. *et al.* Pathophysiological mechanisms of liver injury in COVID-19. *Liver Int. Off.*
543 *J. Int. Assoc. Study Liver* **41**, 20–32 (2021).
- 544 5. Lin, L. *et al.* Gastrointestinal symptoms of 95 cases with SARS-CoV-2 infection. *Gut* **69**,
545 997–1001 (2020).
- 546 6. Xu, E., Xie, Y. & Al-Aly, Z. Long-term gastrointestinal outcomes of COVID-19. *Nat. Commun.*
547 **14**, 983 (2023).
- 548 7. Natarajan, A. *et al.* Gastrointestinal symptoms and fecal shedding of SARS-CoV-2 RNA
549 suggest prolonged gastrointestinal infection. *Med N. Y. N* **3**, 371-387.e9 (2022).
- 550 8. Hoffmann, M. *et al.* SARS-CoV-2 Cell Entry Depends on ACE2 and TMPRSS2 and Is
551 Blocked by a Clinically Proven Protease Inhibitor. *Cell* **181**, 271-280.e8 (2020).
- 552 9. Hamming, I. *et al.* Tissue distribution of ACE2 protein, the functional receptor for SARS
553 coronavirus. A first step in understanding SARS pathogenesis. *J. Pathol.* **203**, 631–637
554 (2004).
- 555 10. Zollner, A. *et al.* Postacute COVID-19 is Characterized by Gut Viral Antigen Persistence in

- 556 Inflammatory Bowel Diseases. *Gastroenterology* **163**, 495-506.e8 (2022).
- 557 11. Stein, S. R. *et al.* SARS-CoV-2 infection and persistence in the human body and brain at
558 autopsy. *Nature* **612**, 758–763 (2022).
- 559 12. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human
560 tissues. *Science* **369**, 1318–1330 (2020).
- 561 13. Ma, D. *et al.* Single-cell RNA sequencing identify SDCBP in ACE2-positive bronchial
562 epithelial cells negatively correlates with COVID-19 severity. *J. Cell. Mol. Med.* **25**,
563 7001–7012 (2021).
- 564 14. Downes, D. J. *et al.* Identification of LZTFL1 as a candidate effector gene at a COVID-19
565 risk locus. *Nat. Genet.* **53**, 1606–1615 (2021).
- 566 15. Qiu, C. *et al.* Renal compartment-specific genetic variation analyses identify new pathways
567 in chronic kidney disease. *Nat. Med.* **24**, 1721–1731 (2018).
- 568 16. Kettunen, J. *et al.* Genome-wide association study identifies multiple loci influencing human
569 serum metabolite levels. *Nat. Genet.* **44**, 269–276 (2012).
- 570 17. Franzén, O. *et al.* Cardiometabolic risk loci share downstream cis- and trans-gene
571 regulation across tissues and diseases. *Science* **353**, 827–830 (2016).
- 572 18. Raj, T. *et al.* Polarization of the effects of autoimmune and neurodegenerative risk alleles in
573 leukocytes. *Science* **344**, 519–523 (2014).
- 574 19. Zhang, J. *et al.* Plasma proteome analyses in individuals of European and African ancestry
575 identify cis-pQTLs and models for proteome-wide association studies. *Nat. Genet.* **54**,
576 593–602 (2022).
- 577 20. Ferkingstad, E. *et al.* Large-scale integration of the plasma proteome with genetics and
578 disease. *Nat. Genet.* **53**, 1712–1721 (2021).
- 579 21. Orho-Melander, M. *et al.* Common missense variant in the glucokinase regulatory protein
580 gene is associated with increased plasma triglyceride and C-reactive protein but lower
581 fasting glucose concentrations. *Diabetes* **57**, 3112–3121 (2008).

- 582 22. Chiang, J. Y. L. & Ferrell, J. M. Up to date on cholesterol 7 alpha-hydroxylase (CYP7A1) in
583 bile acid synthesis. *Liver Res.* **4**, 47–63 (2020).
- 584 23. Pathak, M. & Lal, G. The Regulatory Function of CCR9+ Dendritic Cells in Inflammation and
585 Autoimmunity. *Front. Immunol.* **11**, 536326 (2020).
- 586 24. Severe Covid-19 GWAS Group *et al.* Genomewide Association Study of Severe Covid-19
587 with Respiratory Failure. *N. Engl. J. Med.* **383**, 1522–1534 (2020).
- 588 25. Shelton, J. F. *et al.* Trans-ancestry analysis reveals genetic and nongenetic associations
589 with COVID-19 susceptibility and severity. *Nat. Genet.* **53**, 801–808 (2021).
- 590 26. COVID-19 Host Genetics Initiative. Mapping the human genetic architecture of COVID-19.
591 *Nature* **600**, 472–477 (2021).
- 592 27. Uehara, S., Grinberg, A., Farber, J. M. & Love, P. E. A role for CCR9 in T lymphocyte
593 development and migration. *J. Immunol. Baltim. Md 1950* **168**, 2811–2819 (2002).
- 594 28. Menees, S. & Chey, W. The gut microbiome and irritable bowel syndrome. *F1000Research*
595 **7**, F1000 Faculty Rev-1029 (2018).
- 596 29. Wang, B. *et al.* Alterations in microbiota of patients with COVID-19: potential mechanisms
597 and therapeutic interventions. *Signal Transduct. Target. Ther.* **7**, 143 (2022).
- 598 30. Dupuis, J. *et al.* New genetic loci implicated in fasting glucose homeostasis and their impact
599 on type 2 diabetes risk. *Nat. Genet.* **42**, 105–116 (2010).
- 600 31. Kim, Y. K. *et al.* Identification of a genetic variant at 2q12.1 associated with blood pressure
601 in East Asians by genome-wide scan including gene-environment interactions. *BMC Med.*
602 *Genet.* **15**, 65 (2014).
- 603 32. Ma, M., Lee, J. H. & Kim, M. Identification of a TMEM182 rs141764639 polymorphism
604 associated with central obesity by regulating tumor necrosis factor- α in a Korean population.
605 *J. Diabetes Complications* **34**, 107732 (2020).
- 606 33. Cochet-Bissuel, M., Lory, P. & Monteil, A. The sodium leak channel, NALCN, in health and
607 disease. *Front. Cell. Neurosci.* **8**, 132 (2014).

- 608 34. Bourque, D. K. *et al.* Periodic breathing in patients with NALCN mutations. *J. Hum. Genet.*
609 **63**, 1093–1096 (2018).
- 610 35. Moloney, D. J. *et al.* Fringe is a glycosyltransferase that modifies Notch. *Nature* **406**,
611 369–375 (2000).
- 612 36. Zhou, B. *et al.* Notch signaling pathway: architecture, disease, and therapeutics. *Signal*
613 *Transduct. Target. Ther.* **7**, 95 (2022).
- 614 37. Del Castillo Velasco-Herrera, M. *et al.* Comparative genomics reveals that loss of lunatic
615 fringe (LFNG) promotes melanoma metastasis. *Mol. Oncol.* **12**, 239–255 (2018).
- 616 38. Shelton, J. F. *et al.* The UGT2A1/UGT2A2 locus is associated with COVID-19-related loss of
617 smell or taste. *Nat. Genet.* **54**, 121–124 (2022).
- 618 39. Durand, E. Y., Do, C. B., Mountain, J. L. & Macpherson, J. M. *Ancestry Composition: A*
619 *Novel, Efficient Pipeline for Ancestry Deconvolution.*
620 <http://biorxiv.org/lookup/doi/10.1101/010512> (2014) doi:10.1101/010512.
- 621 40. Wellcome Trust Case Control Consortium *et al.* Bayesian refinement of association signals
622 for 14 loci in 3 common diseases. *Nat. Genet.* **44**, 1294–1301 (2012).
- 623 41. Bowden, J., Davey Smith, G., Haycock, P. C. & Burgess, S. Consistent Estimation in
624 Mendelian Randomization with Some Invalid Instruments Using a Weighted Median
625 Estimator. *Genet. Epidemiol.* **40**, 304–314 (2016).
- 626 42. Bowden, J., Davey Smith, G. & Burgess, S. Mendelian randomization with invalid
627 instruments: effect estimation and bias detection through Egger regression. *Int. J.*
628 *Epidemiol.* **44**, 512–525 (2015).
- 629 43. Hemani, G., Tilling, K. & Davey Smith, G. Orienting the causal relationship between
630 imprecisely measured traits using GWAS summary data. *PLoS Genet.* **13**, e1007081
631 (2017).
- 632 44. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122 (2016).
- 633 45. Kerimov, N. *et al.* A compendium of uniformly processed human gene expression and

- 634 splicing quantitative trait loci. *Nat. Genet.* **53**, 1290–1299 (2021).
- 635 46. Lappalainen, T. *et al.* Transcriptome and genome sequencing uncovers functional variation
636 in humans. *Nature* **501**, 506–511 (2013).
- 637 47. Koolpe, G. A. *et al.* Opioid agonists and antagonists. 6-Desoxy-6-substituted lactone,
638 epoxide, and glycidate ester derivatives of naltrexone and oxymorphone. *J. Med. Chem.* **28**,
639 949–957 (1985).
- 640 48. Craig, D. W. *et al.* RNA sequencing of whole blood reveals early alterations in immune cells
641 and gene expression in Parkinson’s disease. *Nat. Aging* **1**, 734–747 (2021).
- 642 49. Ongen, H., Buil, A., Brown, A. A., Dermitzakis, E. T. & Delaneau, O. Fast and efficient QTL
643 mapper for thousands of molecular phenotypes. *Bioinformatics* **32**, 1479–1485 (2016).
- 644 50. Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of
645 expression residuals (PEER) to obtain increased power and interpretability of gene
646 expression analyses. *Nat. Protoc.* **7**, 500–507 (2012).
- 647 51. Storey, J. D. A Direct Approach to False Discovery Rates. *J. R. Stat. Soc. Ser. B Stat.*
648 *Methodol.* **64**, 479–498 (2002).
- 649 52. Harrow, J. *et al.* GENCODE: the reference human genome annotation for The ENCODE
650 Project. *Genome Res.* **22**, 1760–1774 (2012).
- 651 53. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A Simple New Approach to Variable
652 Selection in Regression, with Application to Genetic Fine Mapping. *J. R. Stat. Soc. Ser. B*
653 *Stat. Methodol.* **82**, 1273–1300 (2020).

654

655

656

657

658

659 **ACKNOWLEDGEMENTS**

660

661 We thank the 23andMe research participants and employees who made this study

662 possible. The following members of the 23andMe Research Team contributed to this

663 study:

664 Stella Aslibekyan, Adam Auton, Elizabeth Babalola, Robert K. Bell, Jessica Bielenberg,

665 Jonathan Bowes, Katarzyna Bryc, Ninad S. Chaudhary, Daniella Coker, Sayantan Das,

666 Emily DelloRusso, Sarah L. Elson, Nicholas Eriksson, Teresa Filshtein, Pierre

667 Fontanillas, Will Freyman, Zach Fuller, Chris German, Julie M. Granka, Karl Heilbron,

668 Alejandro Hernandez, Barry Hicks, David A. Hinds, Ethan M. Jewett, Yunxuan Jiang,

669 Katelyn Kukar, Alan Kwong, Yanyu Liang, Keng-Han Lin, Bianca A. Llamas, Matthew H.

670 McIntyre, Steven J. Micheletti, Meghan E. Moreno, Priyanka Nandakumar, Dominique T.

671 Nguyen, Jared O'Connell, Aaron A. Petrakovitz, G. David Poznik, Alexandra Reynoso,

672 Shubham Saini, Morgan Schumacher, Leah Selcer, Anjali J. Shastri, Janie F. Shelton,

673 Jingchunzi Shi, Suyash Shringarpure, Qiaojuan Jane Su, Susana A. Tat, Vinh Tran,

674 Joyce Y. Tung, Xin Wang, Wei Wang, Catherine H. Weldon, Peter Wilton, Corinna D.

675 Wong.

676

677

678

679

680

681

682

683

684

685

686

687 AUTHOR INFORMATION

688

689 Affiliations:

690

691 **23andMe, Inc:**

692

693 Ninad S. Chaudhary, Catherine H. Weldon, Priyanka Nandakumar,

694 23andMe Research Team, Michael V. Holmes, Stella Aslibekyan

695

696 **Bristol Myers Squibb, Inc:**

697

698 Janie F. Shelton

699

700

701 Author Contributions

702

703 The 23andMe COVID-19 Team developed the recruitment and participant engagement

704 strategy and acquired and processed the data. N.S.C analyzed the data. N.S.C., P.A,

705 C.H.W, S.A, and M.V.H interpreted the data. N.S.C., S.A. and M.V.H. wrote the

706 manuscript. All authors participated in the preparation of the manuscript by reading and

707 commenting on the drafts before submission.

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725 ETHICS DECLARATIONS

726

727 Competing Interests

728

729 C.H.W, P.N, S,A, and M.V.H are current employees of 23andMe and hold stock or stock

730 options in 23andMe. J.S is a current employee of Bristol Myers Squibb. N.S.C works as

731 a postdoctoral fellow on the 23andMe Genetic Epidemiology Team.

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763 FIGURE LEGENDS

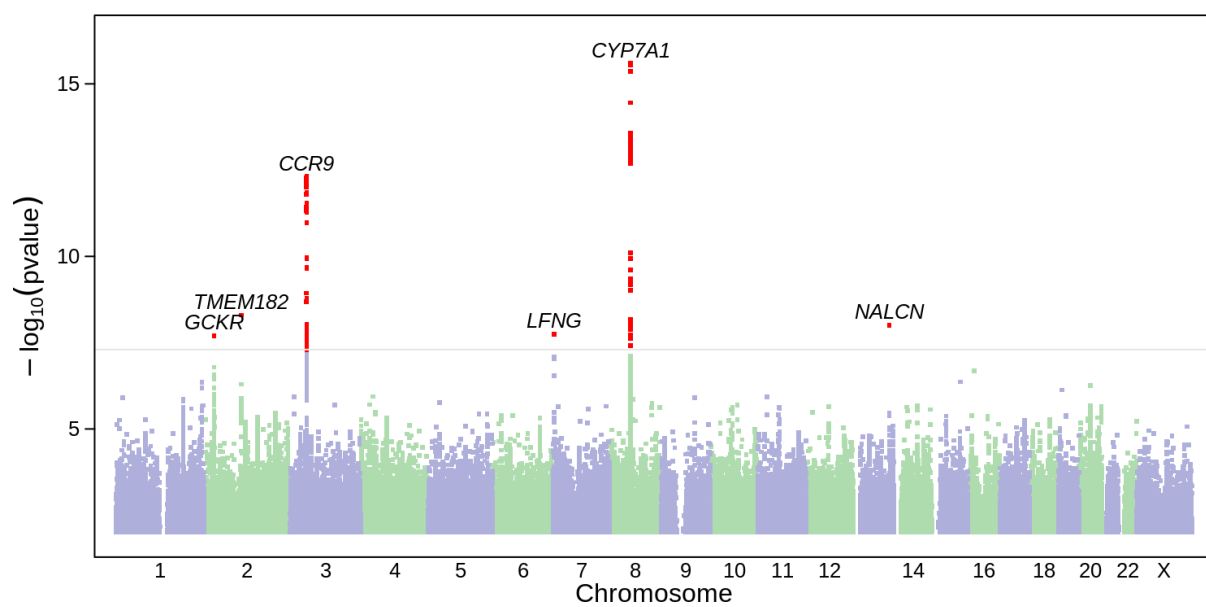
764

765 **Figure 1 : Manhattan plot of COVID-19+ diarrhea among 23andMe participants** 766 **who tested positive for SARS-CoV2**

767 Manhattan plot depicts findings from meta-analysis of five ancestral groups (European,
768 African, Latinx, East Asian, and South Asian). X-axis represents chromosomal position
769 for each SNP. Y-axis represents negative log p values based on logistic regression
770 model under the additive model. Statistically significant variants are highlighted in red.
771 The regions of associations are annotated with index variants.

772

773



774

775

776

777

778

779

780

781

782

783

784

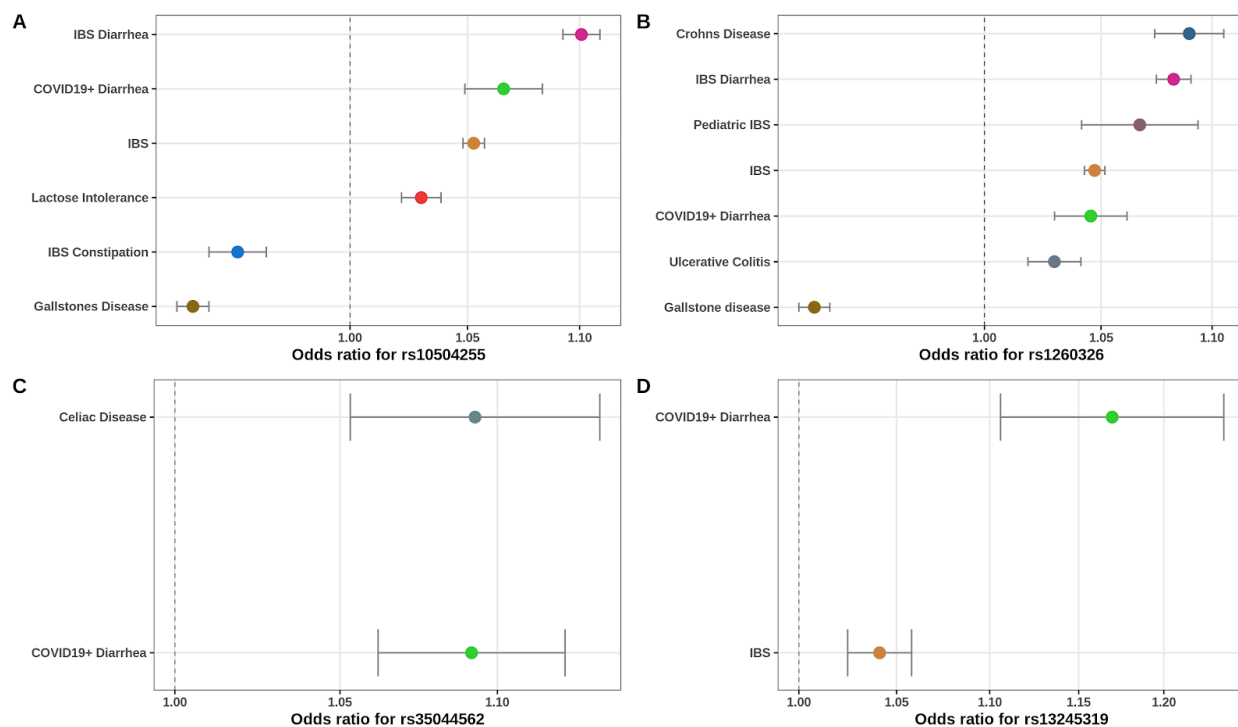
785

786

787

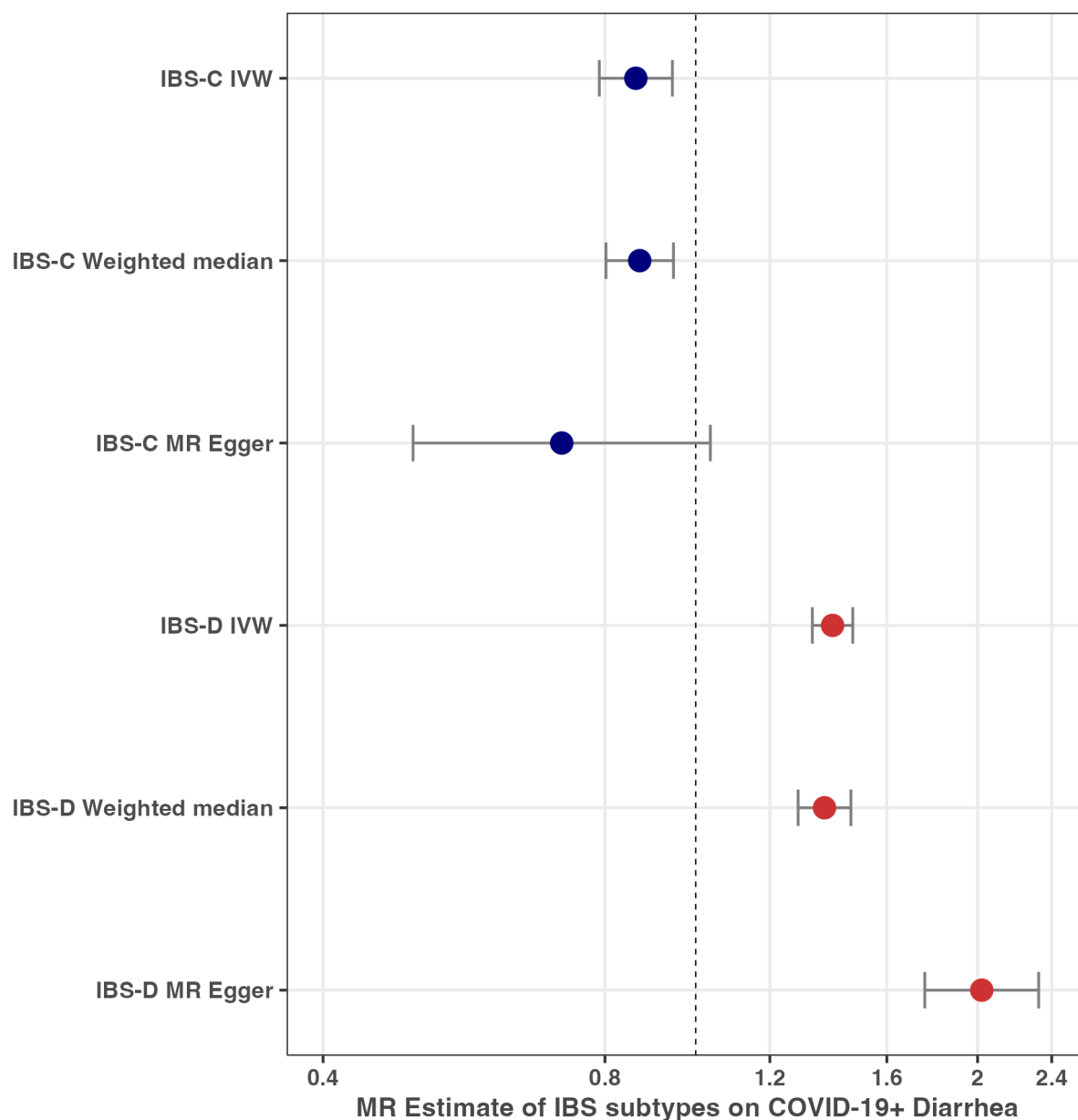
788 **Figure 2: Phenome-wide association of GWAS significant hits of COVID-19+**
789 **diarrhea in 23andMe participants of European ancestry**

790 Forest plot represents statistically significant findings from PheWAS analysis of four loci.
791 The results are presented as odds ratio and 95% confidence intervals under additive
792 model for allele of each loci to represent increased odds of having COVID-19+ diarrhea.
793 X-axis shows estimates on log scale. Y-axis shows phenotypes studied. We have
794 included gastrointestinal phenotypes from PheWAS analysis that met the statistically
795 significant threshold of 5.62×10^{-06} . A complete list of phenotypes that met statistical
796 threshold are included in Supplementary Table 5.
797



798
799
800
801
802
803
804
805
806
807
808
809
810

811 **Figure 3: Forest plot representing genetically predicted effects of irritable bowel**
812 **syndrome subtypes on COVID-19+ diarrhea using Mendelian randomization**
813 The estimates in the plot depict odds ratio and 95% confidence intervals. IBS-C =
814 Irritable bowel syndrome subtype constipation, IBS-D = Irritable bowel syndrome
815 subtype diarrhea, IVW = Inverse variance weighted. The genetic instrument for IBS-C is
816 derived using information from 50 SNPs (meanF statistic = 44.1) and genetic instrument
817 for IBS-D is derived using information from 213 SNPs (meanF statistic = 49.5).
818
819



820
821
822

Table 1. Characteristics of 239,866 research participants in 23andMe COVID-19+ study

Characteristics	COVID-19+ diarrhea	COVID-19+ without diarrhea
No	80,289	159,577
Demographics		
Age mean (SD)	43.23 (14.12)	44.11 (15.10)
Female N (%)	57,065 (71.1)	104,329 (65.4)
Education in years mean (SD)	15.06 (2.37)	15.34 (2.47)
Ancestry N (%)		
<i>African American</i>	3297 (4.1)	7179 (4.5)
<i>East Asian</i>	1282 (1.6)	3047 (1.9)
<i>European</i>	59821 (74.5)	120470 (75.5)
<i>Latinx</i>	15508 (19.3)	27656 (17.3)
<i>South Asian</i>	381 (0.5)	1225 (0.8)
BMI mean (SD)	30.05 (7.55)	28.55 (6.82)
Tobacco Use N (%)	30,091 (38.5)	53,654 (34.5)
Two or more alcohol drinks per week – current N (%)	6,398(23.0)	33,544(23.2)
Health Comorbidities		
High blood pressure N (%)	20,363 (25.8)	34,836 (22.2)
Depression N (%)	32,418 (44.5)	50,773 (34.9)
Diabetes N (%)	10,477 (14.1)	16,083 (10.9)
High total cholesterol N (%)	5679 (19.1)	10371 (16.8)
High triglycerides N (%)	4277 (26.3)	7194 (21.1)
COVID-19+ related characteristics		
COVID-19+ and hospitalized N (%)	5,111 (6.6)	4,570 (3.0)

Table 1. Characteristics of 239,866 research participants in 23andMe COVID-19+ study

Characteristics	COVID-19+ diarrhea	COVID-19+ without diarrhea
COVID-19+ and severe respiratory disease N (%)	8,170 (10.2)	7,702 (4.8)

SD=standard deviation, IQR=Interquartile-range, COPD=Chronic Obstructive Pulmonary Disease, BMI=Body-Mass Index, Severe respiratory disease= pneumonia, difficulty breathing that may have required supplementary oxygen, or ventilatory support

824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844

Table 2: Statistically significant genetic variants associated with COVID-19+ diarrhea on meta-analysis

SNP	Chr	Position	Alleles	Gene	Effect allele	EAF	OR(95%CI)	P value
rs10504255	chr8	58485902	A/G	UBXN2B--[]-CYP7A1	G	0.344	0.94(0.92,0.95)	2.56X10 ⁻¹⁶
rs35044562	chr3	45867532	A/G	LZTFL1--[]-CCR9	G	0.078	1.10(1.07,1.13)	4.88X10 ⁻¹³
rs75683620	chr2	103792340	A/G	TMEM182---[]	G	0.9991	0.67(0.59,0.77)	5.04X10 ⁻⁰⁹
rs536843010	chr13	101100662	A/C	[NALCN]	C	0.991	0.81(0.75,0.87)	9.82X10 ⁻⁰⁹
rs13245319	chr7	2514631	C/T	GRIFIN--[]-LFNG	T	0.021	1.16(1.10,1.22)	1.79X10 ⁻⁰⁸
rs1260326	chr2	27508073	C/T	[GCKR]	T	0.416	1.04(1.03,1.05)	1.98X10 ⁻⁰⁸

Chr = chromosome, EAF = Effect allele frequency, OR = Odds ratio, CI = Confidence Interval

