

# 1 AncestryDNA COVID-19 Host Genetic Study Identifies Three Novel Loci

2 Genevieve H.L. Roberts\*<sup>1</sup>, Danny S. Park\*<sup>2</sup>, Marie V. Coignet<sup>2</sup>, Shannon R. McCurdy<sup>2</sup>, Spencer C.  
3 Knight<sup>2</sup>, Raghavendran Partha<sup>2</sup>, Brooke Rhead<sup>2</sup>, Miao Zhang<sup>2</sup>, Nathan Berkowitz<sup>2</sup>, AncestryDNA  
4 Science Team<sup>1,2</sup>, Asher K. Haug Baltzell<sup>1</sup>, Harendra Guturu<sup>2</sup>, Ahna R. Girshick<sup>2</sup>, Kristin A. Rand<sup>2</sup>,  
5 Eurie L. Hong<sup>2</sup>, Catherine A. Ball<sup>2</sup>

6 \*These authors contributed equally to this work.

7 1. AncestryDNA, 1300 West Traverse Parkway, Lehi, UT 84043

8 2. AncestryDNA, 153 Townsend St, Suite 800, San Francisco, CA, 94107

9 **Corresponding Author:** Catherine Ball, [cball@ancestry.com](mailto:cball@ancestry.com)

10

## 11 Abstract

12 Human infection with SARS-CoV-2, the causative agent of COVID-19, leads to a remarkably diverse  
13 spectrum of outcomes, ranging from asymptomatic to fatal. Recent reports suggest that both clinical and  
14 genetic risk factors may contribute to COVID-19 susceptibility and severity. To investigate genetic risk  
15 factors, we collected over 500,000 COVID-19 survey responses between April and May 2020 with  
16 accompanying genetic data from the AncestryDNA database. We conducted sex-stratified and meta-  
17 analyzed genome-wide association studies (GWAS) for COVID-19 susceptibility (positive  
18 nasopharyngeal swab test,  $n_{cases}=2,407$ ) and severity (hospitalization,  $n_{cases}=250$ ). The severity GWAS  
19 replicated associations with severe COVID-19 near *ABO* and *SLC6A20* ( $P<0.05$ ). Furthermore, we  
20 identified three novel loci with  $P<5\times 10^{-8}$ . The strongest association was near *IVNSIABP*, a gene  
21 involved in influenza virus replication<sup>1</sup>, and was associated only in males. The other two novel loci  
22 harbor genes with established roles in viral replication or immunity: *SRRMI* and the immunoglobulin  
23 lambda locus. We thus present new evidence that host genetic variation likely contributes to COVID-19  
24 outcomes and demonstrate the value of large-scale, self-reported data as a mechanism to rapidly address  
25 a health crisis.

## 26 Introduction

27 Novel severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), the causative agent of  
28 COVID-19, precipitated a pandemic with >21 million cases and >760,000 deaths worldwide as of  
29 August 2020.<sup>2</sup> Outcomes of SARS-CoV-2 infection in the United States are diverse; most infections  
30 result in mild illness that can be managed at home, yet ~14% of cases are hospitalized and ~5% are  
31 fatal.<sup>3</sup> Epidemiological studies have identified clinical risk factors for severe COVID-19 that include  
32 common health conditions such as hypertension, diabetes, obesity, older age, and male sex.<sup>4,5</sup> Reports of  
33 higher susceptibility to<sup>6,7,8</sup> and severity of<sup>9</sup> SARS-CoV infections in men could suggest important  
34 biological differences in immune response to SARS-CoV-2 in men relative to women.<sup>10</sup>

35  
36 In addition to clinical risk factors, emerging evidence suggests that host genetic variation may contribute  
37 to COVID-19 susceptibility and severity. Ellinghaus *et al.* conducted a genome-wide association study  
38 (GWAS) of COVID-19 cases with respiratory failure and identified two loci that achieved genome-wide  
39 significance: one signal on chromosome (chr) 9 near the *ABO* gene, which determines blood type, and  
40 one signal on chr 3 near a cluster of genes with known immune function including *SLC6A20*, *CXCR6*,  
41 *CCR1*, *CCR2*, and *CCR9*.<sup>11</sup> Additionally, a small whole-exome sequencing study identified *TLR7*, an  
42 X-chromosome gene involved in interferon signal induction, in four male patients with severe  
43 COVID-19.<sup>12</sup> To validate rapidly emerging results from new studies such as the Ellinghaus study,  
44 further investigation in independent datasets is needed. Furthermore, investigation in larger datasets with  
45 increased statistical power may detect additional, novel host genetic variation relevant to COVID-19  
46 susceptibility and severity.

47  
48 To replicate and discover non-genetic<sup>6</sup> and genetic associations with COVID-19 outcomes, we engaged  
49 AncestryDNA members who have consented to research in the United States, with 18 million total

50 individuals in the global network. On April 22, 2020, we released a 54-question COVID-19 survey  
51 intended to assess exposure, risk factors, symptomatology, and demographic information that had been  
52 previously identified as associated with COVID-19 susceptibility and severity in the evolving pandemic.  
53 In under two months, over 500,000 COVID-19 survey responses were collected with a 95% survey  
54 completion rate.

55  
56 From these self-reported outcomes, we constructed two phenotypes: one intended to assess  
57 susceptibility, in which individuals who reported a positive COVID-19 test were compared to those who  
58 reported a negative test (referred to throughout as “susceptibility”) and one intended to assess severity,  
59 in which who were hospitalized with COVID-19 were compared to COVID-19 positive individuals who  
60 were not hospitalized (referred to throughout as “hospitalization”). To identify novel genetic  
61 determinants of these outcomes, we conducted a GWAS for each phenotype in a cohort of European-  
62 descent individuals. Sex-stratified GWAS were performed to investigate potential biological sex-driven  
63 differences in immune response to SARS-CoV-2 infection, and the GWAS results were meta-analyzed  
64 to maximize statistical power for variant discovery.

65

## 66 Results

### 67 COVID-19 survey results and cohort demographics

68 To validate the COVID-19 survey, we examined how representative our cohort is with respect to the  
69 United States population both in terms of COVID-19 infection status and demographics. The COVID-19  
70 survey presented to AncestryDNA customers who consented to research has between 39 and 54  
71 questions, depending on reported COVID-19 test result (**Supplementary Figure 1 and Supplementary**  
72 **Table 1**). The first survey question assesses testing status, and in total, 3,733, or 13.4% of respondents  
73 who were tested, reported a positive test result (**Supplementary Table 2**) – a proportion comparable to

74 the national cumulative positivity rate of 12% during a similar collection period (March 1 - May 30,  
75 2020).<sup>13</sup> Of respondents that reported a positive COVID-19 test, 375 (11%) reported hospitalization,  
76 comparable to a U.S. Centers for Disease Control and Prevention (CDC) report of a 14% hospitalization  
77 rate (**Supplementary Table 3**).<sup>3</sup>

78

79 Extended demographic data presented in **Supplementary Table 3** show that the proportion of  
80 COVID-19 survey respondents with European ancestry is largely consistent with the U.S. population:  
81 roughly 75% of respondents were classified as having European ancestry, comparable to 73% of the  
82 U.S. general population that self-identify as “white” in the United States Census Bureau’s 2018  
83 American Community Survey.<sup>14</sup> In the COVID-19 survey, 65% of respondents were female, and the  
84 United States Census Bureau’s 2018 American Community Survey reports a female population of  
85 50.8%. This is consistent with previous literature reporting gender differences in willingness to  
86 participate in survey completion.<sup>15</sup> The median age of COVID-19 survey respondents was 56, notably  
87 higher than the national median age<sup>14</sup> of 38; however, a minimum age of 18 was required to participate  
88 in the study.

89

90 The survey data were filtered to select a final GWAS cohort to analyze (filtering steps summarized in  
91 **Supplementary Figure 2**). **Table 1** presents the final GWAS cohort sample sizes for the susceptibility  
92 and severity phenotypes. In total, 2,417 cases (862 male, 1,555 female) and 14,933 controls (4,472 male,  
93 10,461 female) were used in the susceptibility GWAS and 250 cases (105 male, 145 female) and 1,967  
94 controls (679 male, 1,288 female) were used in the hospitalization GWAS.

95

96 **Replication of *ABO* and *SLC6A20* COVID-19 severity associations**

97 Ellinghaus *et al.* (2020) performed a GWAS consisting of 1,980 Spanish and Italian COVID-19 cases  
98 with respiratory failure and 1,805 controls, most of whom had not been tested for COVID-19. This study  
99 identified two loci that achieved genome-wide significance: one signal on chr 3 represented by  
100 rs11385942 near the *SLC6A20* gene and one signal on chr 9 represented by rs657152 near the *ABO*  
101 gene.<sup>11</sup>

102  
103 We assessed replication of the lead SNPs at the two loci identified by Ellinghaus *et al.*<sup>11</sup> within our  
104 hospitalization GWAS. Although the phenotypes are not the same, both respiratory failure and  
105 hospitalization assess aspects of COVID-19 infection severity. As summarized in **Table 2**, we observed  
106 nominal replication of the lead SNPs at these two loci in our hospitalization analysis at both the *ABO*  
107 locus ( $P=0.022$ ) and the *SLC6A20* locus ( $P=0.020$ ). For both loci, consistent risk alleles and directions  
108 of effect were observed, but with generally smaller odds ratios (ORs) than those reported by Ellinghaus  
109 and colleagues. We did not observe significant associations at either locus in the susceptibility analysis  
110 (**Supplementary Table 4**).

### 111 112 **Novel associations with COVID-19 susceptibility and severity**

113 We conducted a sex-stratified GWAS adjusted for orthogonal age, orthogonal age<sup>2</sup>, array version, and  
114 PC1-12 (**Supplementary Table 5**), to investigate possible differences in genetic associations with  
115 COVID-19 outcomes in males and females and meta-analyzed the resulting summary statistics to  
116 maximize statistical power. For the susceptibility phenotype, the male GWAS, female GWAS, and  
117 meta-analysis had genomic inflation factors of 1.00, 1.01, and 1.00, respectively (**Supplementary**  
118 **Figure 3**). The hospitalization GWAS had genomic inflation factors of 1.04, 1.01, and 0.99 for males,  
119 females, and the meta-analysis, respectively (**Supplementary Figure 4**). In total, three novel loci

120 surpassed the genome-wide significance threshold of  $5 \times 10^{-8}$  in at least one study: two separate loci on  
121 chr 1 and one locus on chr 22 (**Table 3 and Figure 1, Supplementary Figures 5 and 6**).

122  
123 In the susceptibility analysis, the most significant association was represented by lead SNP rs6668622,  
124 with the association present in males only ( $P=3.28 \times 10^{-9}$ ; OR=0.69) (**Figure 1c-d**) and absent in females  
125 ( $P=0.37$ ; OR=0.96) (**Figure 1e-f**). In the sex-combined meta-analysis, the rs6668622 association did not  
126 reach genome-wide significance ( $P=3.83 \times 10^{-5}$ ; OR=0.87). Consistent with the differential association  
127 observed in men and women, there is significant heterogeneity of effect ( $I^2$ ) for rs6668622 between the  
128 male and the female studies ( $I^2=94$ ;  $P=1.6 \times 10^{-5}$ ; **Supplementary Table 6**). This signal is intergenic and  
129 the nearest protein-coding genes to rs6668622 are *IVNSIABP* (~150Kb) and *SWT1* (~288Kb)  
130 (**Supplementary Figure 7**).

131  
132 We identified an additional locus associated with susceptibility on chr 22, represented by lead SNP  
133 rs73166864. This locus achieved genome-wide significance in the meta-analysis ( $P=1.56 \times 10^{-8}$ ;  
134 OR=1.70); the effect sizes were not significantly different ( $I^2=1.6$ ;  $P=0.20$ ; **Supplementary Table 6**) in  
135 the male ( $P=1.19 \times 10^{-5}$ ; OR=1.97) and female ( $P=2.21 \times 10^{-4}$ ; OR=1.55) analyses. The variant  
136 rs73166864 is intergenic and is within 500kb of the immunoglobulin lambda locus<sup>16</sup>, which encodes  
137 proteins used to construct the light chains of antibodies. The nearest protein-coding genes to rs73166864  
138 are *IGLL5*, *GNAZ*, *RSPH14*, *RAB36* and *BCR* (**Supplementary Figure 8**).

139  
140 In the hospitalization meta-analysis, a locus on chr 1 surpassed genome-wide significance. The lead  
141 SNP, rs111972040, is uncommon with a MAF of approximately 1%, but with a large estimated effect  
142 size in the hospitalization meta-analysis ( $P=8.38 \times 10^{-9}$ ; OR=8.29). Although rs111972040 did not  
143 achieve genome-wide significance in either sex-stratified GWAS, the estimated ORs in the male

144 ( $P=3.46 \times 10^{-3}$ ; OR=6.50) and female ( $P=8.01 \times 10^{-7}$ ; OR=9.37) studies were large and there was no  
145 significant heterogeneity of effect between males and females ( $I^2=0$ ;  $P=0.63$ ; **Supplementary Table 6**).  
146 The variant rs111972040 is a non-coding transcript variant in the gene *SRRMI*, and *NCMAP*, *CLIC4*,  
147 *RCAN3*, *NIPAL3*, and *RUNX3* are all within 500kb (**Supplementary Figure 9**).  
148  
149 To assess whether clinical risk factors other than age and sex had an effect on these associated loci, we  
150 additionally adjusted for other associated risk factors including body mass index (BMI) and having one  
151 or more pre-existing health conditions (asthma, bone marrow transplant, cancer, cardiovascular disease,  
152 kidney disease, chronic obstructive pulmonary disease (COPD), diabetes, hypertension, organ  
153 transplant, autoimmune disease, immunodeficiency, or lung conditions). For all three novel lead SNPs,  
154 the estimated effect sizes remained relatively consistent, though the  $P$ -value dropped below  
155 genome-wide significance in the hospitalization analysis, likely in part due to the small sample size and  
156 the additional decrease in sample size for these extended analyses (**Supplementary Table 7**). Both loci  
157 associated with susceptibility remained genome-wide significant even after adjusting for these additional  
158 clinical risk factors.

159

### 160 **Heritability of COVID-19 susceptibility and severity**

161 We estimated narrow-sense heritability ( $h^2$ ) – defined as the proportion of phenotypic variation due to  
162 additive genetic factors<sup>17</sup> – using the linear mixed model approach (GCTA-GREML)<sup>18</sup> and autosomal  
163 genome-wide imputed variants. The underlying assumption of this linear mixed model approach is that  
164 the included common variants all contribute equally small effects to heritability.<sup>19</sup>

165

166 To our knowledge, these are the first estimates of heritability for COVID-19 susceptibility and severity  
167 (hospitalization) based on empirical genetic similarity, though one twin-based study estimated

168 heritability of COVID-19 symptoms.<sup>20</sup> As shown in **Supplementary Table 8**, for susceptibility, the  
169 estimated liability heritability was  $h^2=0.00$  (standard error [SE]=0.07). The estimated liability  
170 heritability for hospitalization was  $h^2=0.14$  (SE=0.58); the large standard error is partly a result of the  
171 small hospitalization sample size.

172

## 173 Discussion

174 To identify genetic determinants of COVID-19 susceptibility and severity, we conducted GWAS of self-  
175 reported COVID-19 outcomes in a population of survey respondents with European ancestry. To explore  
176 possible differences in biological response to SARS-CoV-2 infection, we analyzed both susceptibility  
177 and severity outcomes via sex-stratified GWAS and sex-combined meta-analyses. In total, three novel  
178 loci surpassed genome-wide significance in one or more study, with lead SNPs rs6668622 (male  
179 susceptibility  $P=3.2\times 10^{-9}$ ), rs73166864 (susceptibility  $P=1.56\times 10^{-8}$ ), and rs111972040 (hospitalization  
180  $P=8.38\times 10^{-9}$ ) near *IVNSIABP*, the immunoglobulin lambda locus, and *SRRMI*, respectively.

181

182 The most significantly associated SNP in any study was rs6668622 with the susceptibility outcome. The  
183 nearest gene to rs6668622 is *IVNSIABP*, which encodes Influenza Virus NS1A Binding Protein. This  
184 protein is known to bind with influenza virulence factor NS1 and this interaction appears to promote  
185 influenza viral gene expression.<sup>21</sup> The variant rs6668622 is a known, strong expression quantitative trait  
186 locus (eQTL) in lung tissue for *IVNSIABP*<sup>22,23</sup>, suggesting that risk variation might impact mRNA  
187 abundance of *IVNSIABP*. Strikingly, haploinsufficiency of *IVNSIABP* appears to associate with primary  
188 immunodeficiency<sup>24</sup>, suggesting *IVNSIABP* may play a role in cellular response to other pathogens  
189 besides influenza. It is unclear why this association is only present only in males, though it may provide  
190 a clue as to why males appear to be at higher risk of COVID-19 infection, hospitalization, and  
191 mortality.<sup>7-12</sup> We speculate that sex hormones or behavioral differences might trigger generally different

192 cellular responses to SARS-Cov-2 infection in men and in women<sup>10</sup>, and one such difference may  
193 involve differential expression of *IVNSIABP*.

194

195 Another locus, represented by the intergenic lead SNP rs73166864, was associated with the  
196 susceptibility outcome with similar effects in men and women. This signal is ~75Kb away from the  
197 immunoglobulin lambda locus, a region that undergoes somatic recombination in B-cells and encodes  
198 proteins used to construct the antigen-binding light chain region of antibodies.<sup>16</sup> It is unclear what the  
199 functional consequence of this intergenic variation might be, but proximity to such an important region  
200 for antibody generation is intriguing.

201

202 The final locus, represented by the lead SNP rs111972040, was the only genome-wide significant  
203 association with the hospitalization outcome. The variant rs111972040 is a non-coding transcript variant  
204 in the gene *SRRM1*, which encodes Serine and Arginine Repetitive Matrix 1. A related gene, *SRRM2*,  
205 was implicated in human immunodeficiency virus (HIV) splicing and replication.<sup>25</sup> Based on this, we  
206 hypothesize that *SRRM1* may likewise play an important role in splicing SARS-Cov-2 viral genes.

207

208 In addition to identifying novel associations, the severity GWAS replicated findings from a previous  
209 COVID-19 respiratory failure GWAS<sup>11</sup> that identified two loci: the blood type *ABO* gene and a cluster  
210 of immune genes near *SLC6A20*. We observed consistent directions of effect at both loci, but the  
211 replication *P*-values were nominal, and we observed smaller estimated ORs; however, these  
212 observations are not surprising. The original study considered a more severe outcome (respiratory  
213 failure) than our analogous severity study (hospitalization). Furthermore, we included only 250 cases  
214 that were hospitalized in our study relative to the 1,980 cases with respiratory failure in the Ellinghaus  
215 study<sup>11</sup>, thus the nominal *P*-values in this study may simply reflect lower statistical power for our

216 severity analysis. Lastly, the winner's curse suggests that overestimated effect sizes are expected in the  
217 discovery study.<sup>26</sup>

218

219 We estimated heritability for both outcomes to assess the contribution of common genetic variation to  
220 COVID-19 susceptibility and severity. The liability heritability estimates were generally small:  $h^2=0.00$   
221 (SE=0.07) for susceptibility and  $h^2=0.14$  (SE=0.58) for hospitalization. However, for hospitalization, the  
222 sample sizes were small for this type of analysis and standard errors were correspondingly high. A key  
223 assumption underlying the heritability estimates is that many small effects were distributed across the  
224 entire genome<sup>19</sup>. Thus, the low heritability estimates might not simply reflect low genetic contribution,  
225 but rather could suggest a genetic architecture involving a limited number of high-effect variants or a  
226 large contribution by uncommon and rare variants.

227

228 A key limitation of our data is that COVID-19 cases who suffered very severe or fatal infections are less  
229 likely to participate in our survey, which consequently results in undersampling cases with severe  
230 outcomes. We also restricted to individuals of European ancestry due to small sample sizes in other  
231 genetic ancestry cohorts for the susceptibility and severity outcomes in this early phase of COVID-19  
232 survey data collection. As the COVID-19 survey cohort grows, future analyses will focus on increased  
233 ancestral diversity to increase generalizability. Finally, we lack an independent replication cohort for our  
234 novel findings and will rely on future ascertainment of additional survey respondents and COVID-19  
235 GWAS consortia<sup>27</sup> efforts to determine whether our findings are reproducible.

236

237 In summary, we collected over 500,000 self-reported COVID-19 outcomes in under two months and  
238 conducted one of the largest genetic studies of infection susceptibility and severity to date, thus  
239 demonstrating the value of large-scale self-reported data as a mechanism to rapidly address a serious

240 health crisis. We identified three novel loci, all of which harbor genes with established roles in viral  
241 replication or immunity, and one of which may provide insight into why men appear to be differently  
242 affected by COVID-19 than women. We thus add to growing evidence that host genetic variation  
243 contributes to COVID-19 susceptibility and severity and suggest identification of such genetic risk  
244 factors may provide profound insight into pathogenesis of the novel coronavirus.

## 245 Methods

246 **Ethics statement.** All data for this research project was from subjects who provided prior informed  
247 consent to participate in AncestryDNA’s Human Diversity Project, as reviewed and approved by our  
248 external institutional review board, Advarra (formerly Quorum). All data was de-identified prior to use.

249  
250 **Study population.** Self-reported COVID-19 outcomes were collected through the Personal Discoveries  
251 Project®, a survey platform available to AncestryDNA customers via the web and mobile applications.  
252 The COVID-19 survey ranged from 39-54 questions, depending on the initial COVID-19 test result  
253 reported. **Supplementary Figure 1** describes the flow of the topics assessed in each section of the  
254 survey. The full questions and possible responses for the two questions used as primary outcomes in this  
255 study are presented in **Supplementary Table 1**. Analyses presented here were performed with data  
256 collected between April 22-May 28, 2020.

257  
258 To participate in the COVID-19 survey, participants must meet the following criteria: they must be 18  
259 years of age or older, a resident of the United States, be an existing AncestryDNA customer who has  
260 consented to participate in research<sup>28</sup>, and be able to complete a short survey. The survey is designed to  
261 assess self-reported COVID-19 positivity and severity, as well as susceptibility and known risk factors  
262 including community exposure and known contacts with individuals diagnosed with COVID-19.

263  
264 **Phenotype definitions.** Two phenotypes were assessed: one for susceptibility and one for severity of  
265 COVID-19. Cases for the COVID-19 susceptibility phenotype were defined as individuals who  
266 responded to the question, “Have you been swab tested for COVID-19, commonly referred to as  
267 coronavirus?” as “Yes, and was positive”. Responders who answered “Yes, and was negative” were  
268 defined as controls for the susceptibility study. Cases for the severity phenotype were defined as

269 individuals who reported testing positive for COVID-19 and responded to the question, “Were you  
270 hospitalized due to these symptoms?” as “Yes”. Severity controls reported testing positive for  
271 COVID-19, but reported no hospitalization related to COVID-19.

272

273 **Genotyping.** Genotyping and quality control procedures have been previously described elsewhere.<sup>28</sup>  
274 Briefly, customer genotype data for this study were generated using an Illumina genotyping array with  
275 approximately 730,000 SNPs and processed either with Illumina or with Quest/Athena Diagnostics. To  
276 ensure quality of each dataset, a sample passes a number of quality control (QC) checks, which includes  
277 identifying duplicate samples, removing individuals with a per-sample call rate <98%, and identifying  
278 discrepancies between reported sex and genetically inferred sex. Samples that pass all quality-control  
279 tests proceed to the analysis pipeline; samples that fail one or more tests must be recollected or manually  
280 cleared for analysis by lab technicians. Array markers with per-variant call rate <0.98 and array markers  
281 that had overall allele frequency differences of >0.10 between any two array versions were additionally  
282 removed prior to downstream analyses.

283

284 **Selecting a European ancestry association cohort.** As a measure to control population stratification,  
285 we selected individuals with estimated European ancestry for inclusion in the GWAS. To determine  
286 these definitions, a proprietary algorithm was used to estimate continental admixture proportions for all  
287 COVID-19 survey respondents.<sup>29</sup> Briefly, this algorithm uses a hidden Markov model to estimate  
288 unphased diploid ancestry across the genome by comparing haplotype structure to a reference panel. The  
289 reference panel consists of a combination of AncestryDNA customers and publicly available datasets  
290 and is designed to reflect global diversity. From our total cohort of 506,743 individuals who participated  
291 in the COVID-19 survey, 379,383 (74.9%) individuals with estimated European ancestry were retained  
292 **(Supplementary Table 3 and Supplementary Figure 10).**

293

294 **Removal of related individuals.** AncestryDNA's identity-by-descent inference algorithm<sup>30</sup> was used to  
295 estimate the relationship between pairs of individuals. Pairs with estimated separation of fewer than four  
296 meioses were considered close relatives. For all close relative pairs, one individual was randomly  
297 selected for exclusion from our study. We excluded 1,741 individuals from the susceptibility analysis  
298 and 223 from the severity analysis due to relatedness.

299

300 **Calculation of principal components to control residual population structure.** After selecting  
301 unrelated individuals with European ancestry as described above, genetic PCs were calculated to include  
302 in the association studies to control residual population structure and were computed using FlashPCA  
303 2.0.<sup>31</sup> Input genotypes were linkage disequilibrium (LD)-pruned using PLINK 1.9 command `--indep-`  
304 `pairwise 100 5 0.2 --maf 0.05 --geno 0.001`.

305

306 **Imputation.** Samples were imputed to the Haplotype Reference Consortium (HRC) reference panel<sup>32</sup>  
307 version 1.1, which consists of 27,165 total individuals and 36 million variants. The HRC reference panel  
308 does not include indels; consequently, indels are not present in the results of our analyses. We  
309 determined best-guess haplotypes with Eagle version 2.4.1<sup>33</sup> and performed imputation with Minimac4  
310 version 1.0.1.<sup>34</sup> We used 1,117,080 unique variants as input and 8,049,082 imputed variants were  
311 retained in the final data set. For these variants, we conservatively restricted our analyses to variants  
312 minor allele frequency (MAF) $>0.01$  and Minimac4  $R^2 > 0.30$  using genotype dosage probabilities for all  
313 variants regardless of whether they were originally genotyped.

314

315 **Statistical analyses.** We used the COVID-19 Host Genetics Initiative (HGI) analysis plan version 1 to  
316 guide our analyses.<sup>35</sup> A key recommendation in this plan is to analyze males and females separately

317 when possible; therefore, we conducted four separate GWAS in total: susceptibility in males,  
318 susceptibility in females, hospitalization in males, and hospitalization in females. Sex was determined  
319 from genotype data. For each GWAS, a fixed-effects logistic regression model was implemented with  
320 PLINK2.0 with either the susceptibility or severity phenotype as the primary outcome and imputed  
321 genotype dosage value as the primary predictor. The following were included as fixed-effect covariates:  
322 PCs 1-12 (described above), array platform, orthogonal age, and orthogonal age<sup>2</sup>. Orthogonal  
323 polynomials were used to eliminate collinearity between age and age<sup>2</sup> and were calculated in R version  
324 3.6.0 with base function `poly(age, degree=2)`. We additionally used PLINK2.0 to remove  
325 variants with Minimac4 imputation quality  $R^2 < 0.3$  or with  $MAF < 0.01$ . The following PLINK2.0 flags  
326 were used for each analysis:

```
327     --vcf [input imputed VCF] dosage=DS
328     --psam [file that provides sex information]
329     --covar [covariates file]
330     --covar-name PC1, PC2, PC3, PC4, PC5, PC6, PC7, PC8, PC9, PC10, PC11, PC12,
331     orthogonal_age, orthogonal_age2, platform
332     --covar-variance-standardize
333     --extract-if-info R2 >= 0.3
334     --freq
335     --glm hide-covar
336     --keep [list of unrelated Europeans]
337     --keep-females OR keep-males
338     --maf 0.01
339     --pheno [phenotype file]
340     --pheno-name [phenotype column name]
```

341 For each phenotype, summary statistics for males and females were combined using a fixed-effects  
342 inverse-variance weighted meta-analysis, implemented with METAL<sup>36</sup> (version released 25 March  
343 2011). Unless otherwise noted, all variant ORs are adjusted for the 15 covariates described above.

345

346 For all individual GWAS and meta-analyses, we considered the European genome-wide significance  
347 threshold<sup>37</sup> of two-tailed  $P < 5 \times 10^{-8}$  to represent a significant association. For lead SNPs at loci  
348 surpassing  $P < 5 \times 10^{-8}$ , we performed extended analyses in which we additionally adjusted for BMI and  
349 self-reported affliction with one or more of any of the following health conditions: asthma, bone marrow  
350 transplant, cancer, cardiovascular disease, kidney disease, COPD, diabetes, hypertension, organ failure  
351 requiring transplant, autoimmune disease, immunodeficiency, and/or “other” lung condition.

352

353 **Replication.** To our knowledge, the only published GWAS of COVID-19 outcomes to date was  
354 performed by Ellinghaus *et al.*<sup>11</sup> We compared  $P$ -values and OR estimates from our analyses to the two  
355 lead variants reported by Ellinghaus: rs657152, representing a region on chr 9 near *ABO*; and  
356 rs11385942, representing a region on chr 3 near *SLC6A20*. The variant rs11385942 is a small indel and  
357 indels are not present in the HRC reference panel; we therefore examined the association with  
358 rs11385942 using rs17713054, a SNP in perfect LD ( $R^2=1$ ,  $D'=1$ ) in a European population from  
359 LDpair.<sup>38</sup> The allele rs17713054-A corresponds to rs11385942-AG.

360

361 **Heritability.** To estimate heritability of COVID-19 phenotypes, we calculated phenotypic variance  
362 explained by a genetic relatedness matrix (GRM) in the European GWAS cohort. To calculate the GRM,  
363 we first performed LD-pruning of autosomal, imputed, high-quality (Minimac4  $R^2 > 0.3$ ) best-guess  
364 genotypes with PLINK-1.9 command `--indep-pairwise 100 5 0.2 --maf 0.05 --geno`  
365 `0.001 --chr 1-22`, resulting in a total of 224,096 variants. From these variants, GCTA version  
366 1.91.5 beta2 was used to estimate a GRM.<sup>18</sup> We subset the resulting matrix to unrelated individuals (`--`  
367 `grm-cutoff 0.025`) and individuals with diagonal elements between 0.95-1.05, resulting in a final  
368 GRM consisting of 18,415 subjects. These 18,415 individuals are candidates for the two heritability  
369 analyses.

370

371 Next, we performed single component heritability estimation for each phenotype by fitting a linear  
372 mixed model (GCTA-GREML)<sup>18</sup>. Each model included PCs 1-12, array platform, sex inferred from  
373 genotype data, orthogonal age, and orthogonal age<sup>2</sup> as fixed effects. For each phenotype, the model also  
374 included a random genetic effect with the covariance given by a GRM calculated with the intersection of  
375 the 18,415 heritability individuals and those with non-missing phenotypes using the same 224,096  
376 variants described above.

377

378 The observed heritability was estimated using the GCTA command `--grm [phenotype-  
379 specific GRM] --reml --pheno [phenotype file] --covar [genetic sex,  
380 platform] --qcovar [PCs 1-12, orthogonal_age, orthogonal_age2]. The  
381 liability heritability was estimated from the observed heritability using the standard transformation39 and  
382 the sample prevalence using the GCTA command above with --prevalence [sample  
383 prevalence].`

384 Tables

385 **Table 1:** Total number of cases and controls for the susceptibility and hospitalization GWAS

	Susceptibility		Hospitalization	
	Cases COVID-19 Positive	Controls COVID-19 Negative	Cases COVID-19 Cases Reporting Hospitalization	Controls COVID-19 Cases Reporting NO Hospitalization
All*	2,417	14,933	250	1,967
Male*	862	4,472	105	679
Female*	1,555	10,461	145	1,288
Age (median)*	52	55	58	50

386 \*All individuals in the GWAS cohort are of European ancestry. A summary of exclusion criteria is presented in  
387 Supplementary Figure 2.

388 **Table 2:** Replication of severe respiratory failure loci from Ellinghaus *et. al.* (2020)<sup>11</sup> with  
 389 hospitalization phenotype  
 390

Chr	Pos	RSID	Genes	Ref	Alt	Ellinghaus <i>P</i> -value	Ellinghaus OR (95% CI)	Alt Freq.	Alt OR (95% CI)	<i>P</i> -value
3	45876460	rs11385942	<i>SLC6A20</i> , <i>CCR9</i> , <i>FYCO1</i> , <i>CXCR6</i> , <i>XCRI</i>	G	GA	1.15×10 <sup>-10</sup>	1.77 (1.48-2.11)	0.08*	1.43* (1.06-1.93)	0.020*
9	136139265	rs657152	<i>ABO</i>	C	A	4.95×10 <sup>-8</sup>	1.32 (1.20-1.47)	0.36	1.26 (1.03-1.54)	0.022

391 *Chr*, chromosome; *Pos*, position (hg19 genome build); *RSID*, dbSNP identifier; *Ref*, reference allele; *Alt*, alternate allele;  
 392 *OR*, odds ratio of alternate allele; *Alt Freq.*, allele frequency of alternate allele; *CI*, confidence interval.

393 \*A SNP (rs17713054) in perfect LD with rs11385942 ( $r^2=1$ ) in European population from 1000 Genomes phase 3 was used  
 394 to assess replication of the chr 3 locus. See Methods for details.

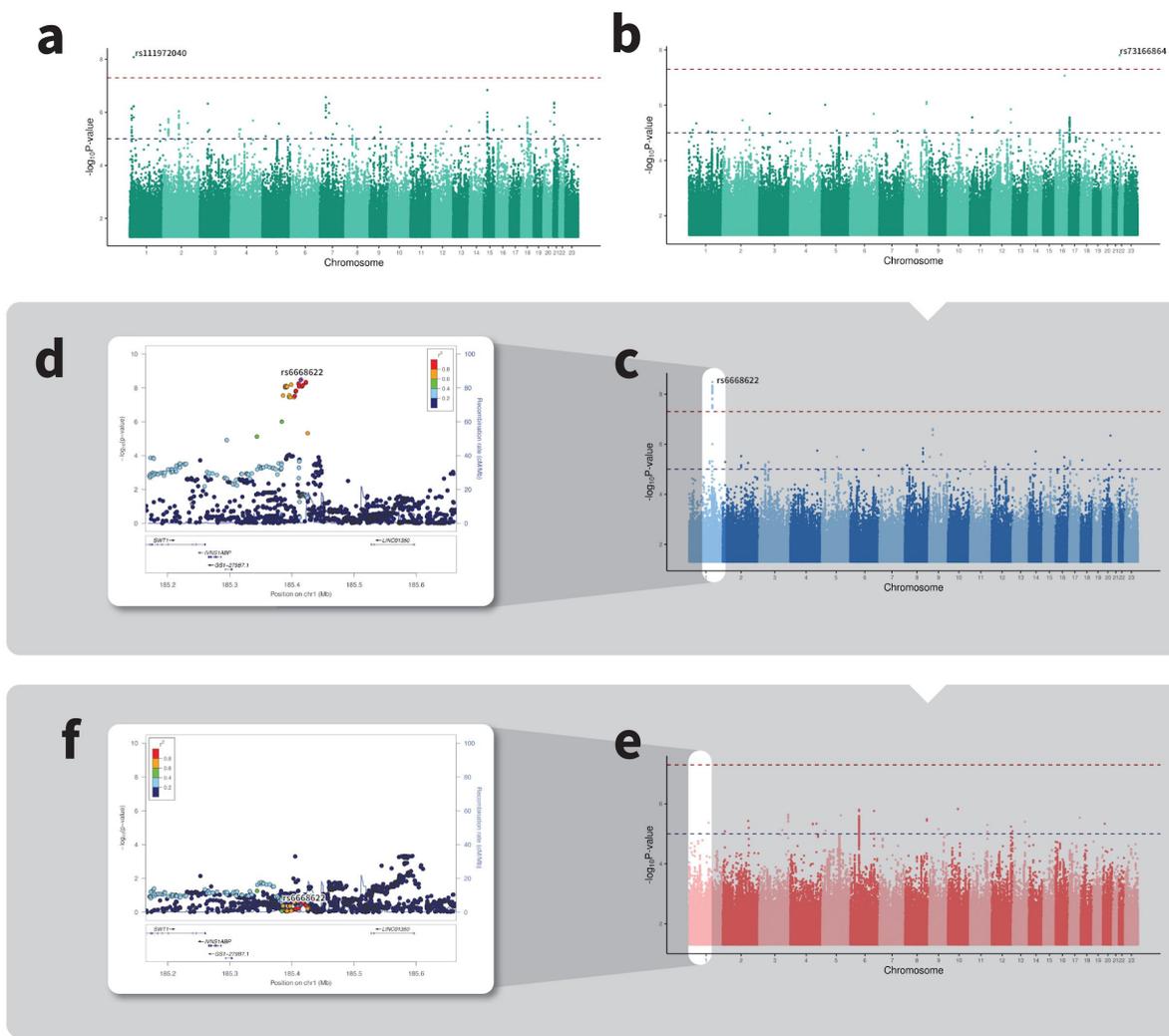
395 **Table 3:** Novel loci with  $P < 5 \times 10^{-8}$  for one or more phenotypes

Chr	Pos (hg19)	RSID	Nearest Genes (within 500 Kb)	Ref	Alt	Phenotype	Alt Freq.	Alt OR (95 % CI)	P-value*
1	185414582	rs6668622	<i>IVNSIABP</i> , <i>SWT1</i>	T	C	<b>Male Susceptibility</b>	<b>0.31</b>	<b>0.69 (0.61-0.78)</b>	<b>3.28x10<sup>-9</sup></b>
						Meta Susceptibility	0.30	0.87 (0.81-0.93)	3.83x10 <sup>-5</sup>
						Female Susceptibility	0.30	0.96 (0.89-1.05)	0.37
						Male Hospitalization	0.31	0.93 (0.64-1.36)	0.71
						Meta Hospitalization	0.30	0.97 (0.78-1.20)	0.75
						Female Hospitalization	0.30	0.97 (0.74-1.28)	0.84
1	24999361	rs111972040	<i>SRRM1</i> , <i>NCMAP</i> , <i>CLIC4</i> , <i>RCAN3</i> , <i>NIPAL3</i> , <i>RUNX3</i>	A	G	<b>Meta Hospitalization</b>	<b>0.01</b>	<b>8.29 (4.04-17.02)</b>	<b>8.38x10<sup>-9</sup></b>
						Female Hospitalization	0.01	9.37 (3.85-22.8)	8.01x10 <sup>-7</sup>
						Male Hospitalization	0.01	6.50 (1.85-22.77)	3.46x10 <sup>-3</sup>
						Female Susceptibility	0.01	0.76 (0.49-1.17)	0.21
						Meta Susceptibility	0.01	0.80 (0.56-1.14)	0.22
						Male Susceptibility	0.01	0.89 (0.49-1.62)	0.70
22	23340580	rs73166864	Immunoglobulin Lambda Locus, <i>RSPH14</i> , <i>GNAZ</i> , <i>RAB36</i> , <i>BCR</i>	T	C	<b>Meta Susceptibility</b>	<b>0.02</b>	<b>1.70 (1.42-2.05)</b>	<b>1.56x10<sup>-8</sup></b>
						Male Susceptibility	0.03	1.97 (1.45-2.66)	1.19x10 <sup>-5</sup>
						Female Susceptibility	0.02	1.55 (1.23-1.96)	2.21x10 <sup>-4</sup>
						Female Hospitalization	0.02	0.73 (0.32-1.70)	0.47
						Male Hospitalization	0.03	1.24 (0.53-2.89)	0.63
						Meta Hospitalization	0.02	0.94 (0.52-1.70)	0.84

396 *Chr*, chromosome; *Pos*, position; *RSID*, dbSNP identifier; *Kb*, kilobase; *Ref*, reference allele; *Alt*, alternate allele; *OR*, odds ratio of alternate allele; *Alt*  
397 *Freq.*, allele frequency of alternate allele; *CI*, confidence interval.

398 \*Results for each locus are sorted by P-value. The study with the most significant association is indicated with bold text.

399 Figures



400  
401 **Figure 1:** Manhattan plots with horizontal red dashed line representing genome-wide significance  
402 ( $P=5\times 10^{-8}$ ) and horizontal blue dashed line representing suggestive significance ( $P=1\times 10^{-5}$ ) for the (a)  
403 hospitalization meta-analysis (b) susceptibility meta-analysis (c) susceptibility GWAS in males only  
404 (blue), highlighting the association peak in the *IVNSIABP* region represented by rs6668622 (d)  
405 LocusZoom plot indexed on rs6668622 in male-only susceptibility GWAS (e) susceptibility GWAS in  
406 females only (red), highlighting the absence of an association peak in the *IVNSIABP* region represented  
407 by rs6668622 (f) LocusZoom plot indexed on rs6668622 in female-only susceptibility GWAS.

## 408 Footnotes

409 **Acknowledgements** We thank our AncestryDNA customers who made this study possible by  
410 contributing information about their experience with COVID-19 through our survey. Without them, this  
411 work would not be possible. We would like to thank Zach Bass, Robert Dowling, Disha Akarte, Swapnil  
412 Sneham, Sean Enright and the entire Cyborg team for their tireless work in the release and continued  
413 support of the COVID-19 survey. We would like to acknowledge Chris Trainor and David Serventi for  
414 their work on the Figures. We additionally thank Mark Daly and the COVID-19 Host Genetics  
415 Consortium for advice and preliminary feedback.

416  
417 **Author contributions** GHLR and DSP contributed equally to the manuscript. GHLR and DSP wrote  
418 the manuscript with support from MVC and SRM. GHLR and DSP designed and conducted  
419 genome-wide association studies, GHLR and DSP interpreted results with support from MVC. SCK,  
420 RP, and BR helped with additional analyses and interpretation. MZ, HG, SCK, NB, AHB, and DSP  
421 performed genotype imputation and data management. MVC and KAR designed the COVID-19 survey  
422 questionnaire and MVC assessed concordance of phenotype prevalence with the U.S. population. SRM  
423 conducted heritability analyses. ARG, AHB, and HG facilitated forward progression of the manuscript  
424 and provided input and guidance. The AncestryDNA Science Team contributed to additional work,  
425 allowing for the completion of the COVID-19 research and manuscript. KAR led the COVID-19  
426 research and data teams. KAR, ELH, and CAB provided project guidance. All authors have contributed  
427 to and reviewed the final manuscript.

428  
429 **Competing interests** The authors declare competing financial interests: authors affiliated with  
430 AncestryDNA are employed by Ancestry and may have equity in Ancestry.

431

432 **AncestryDNA Science Team** Yambazi Banda, Ke Bi, Robert Burton, Marjan Champine, Ross Curtis,  
433 Karen Delgado, Abby Drokhlyansky, Ashley Elrick, Cat Foo, Michael Gaddis, Jialiang Gu, Heather  
434 Harris, Shannon Hateley, Shea King, Christine Maldonado, Evan McCartney-Melstad, Alexandra  
435 McFarland, Patty Miller, Luong Nguyen, Keith Noto, Milos Pavlovic, Jingwen Pei, Jenna Petersen,  
436 Scott Pew, Chodon Sass, Josh Schraiber, Alisa Sedghifar, Andrey Smelter, Sarah South, Barry Starr,  
437 David Turissini, Cecily Vaughn, Yong Wang  
438

## 439 References

- 
- <sup>1</sup> Othumpangat, S., Noti, J.D., Blachere, F.M. & Beezhold, D.H. Expression of non-structural-1A binding protein in lung epithelial cells is modulated by miRNA-548an on exposure to influenza A virus. *Virology*. **447**, 84-94 (2013).
- <sup>2</sup> World Health Organization. Coronavirus disease (COVID-19) Situation Report – 209. [https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200816-covid-19-sitrep-209.pdf?sfvrsn=5dde1ca2\\_2](https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200816-covid-19-sitrep-209.pdf?sfvrsn=5dde1ca2_2) (2020).
- <sup>3</sup> Stokes, E.K., *et al.* Coronavirus disease 2019 case surveillance – United States, January 22-May 30, 2020. *MMWR Morb Mort Weekly Rep.* **69**, 759-765 (2020).
- <sup>4</sup> Grasselli, G., *et al.* Baseline characteristics and outcomes of 1591 patients infected with SARS-CoV-2 admitted to ICUs of the Lombardy region, Italy. *JAMA*. **323**, 1574-1581 (2020).
- <sup>5</sup> Richardson, S., *et al.* Presenting characteristics, comorbidities, and outcomes among 5700 patients hospitalized with COVID-19 in the New York City area. *JAMA*. **323**, 2052-2059 (2020).
- <sup>6</sup> Knight, S., *et al.* COVID-19 Susceptibility and Severity Risks in Over 500,000 Individuals. **[In Preparation]**.
- <sup>7</sup> Channappanavar, R., *et al.* Sex-based differences in susceptibility to severe acute respiratory syndrome coronavirus infection. *J Immunol*, **10**, 198, 4046-4053 (2017).
- <sup>8</sup> Allen, William E., *et al.* Population-scale longitudinal mapping of COVID-19 symptoms, behaviour and testing. *Nat. Hum. Behav.*, 1-11 (2020).
- <sup>9</sup> Womersley, K., Ripullone, K. & Peters, S.A.E., Woodward, M. Covid-19: male disadvantage highlights the importance of sex disaggregated data. *BMJ*. **370**, m2870 (2020).
- <sup>10</sup> Scully, E.P., Haverfield, J. & Ursin, R.L. Considering how biological sex impacts immune responses and COVID-19 outcomes. *Nat Rev Immunol*. **20**, 442-447 (2020).
- <sup>11</sup> Ellinghaus, D., *et al.* Genomewide association study of severe Covid-19 with respiratory failure. *N Engl J Med.* (2020).

- 
- <sup>12</sup> van der Made, C.I., *et al.* (2020). Presence of genetic variants among young men with severe COVID-19. *JAMA*, 324(7), 663-673.
- <sup>13</sup> U.S. Centers for Disease Control and Prevention (CDC). COVIDView: a weekly surveillance summary of U.S. COVID-19 activity – key updates for Week 22, ending May 30, 2020. <https://www.cdc.gov/coronavirus/2019-ncov/covid-data/pdf/covidview-06-05-2020.pdf> (2020).
- <sup>14</sup> United States Census Bureau. American Community Survey (ACS) 1-year estimates data profiles, Table DP05. [https://data.census.gov/cedsci/table?q=United%20States&table=DP05&tid=ACSDP1Y2018.DP05&g=0100000US&lastDisplayedRow=29&vintage=2017&layer=state&cid=DP05\\_0001E](https://data.census.gov/cedsci/table?q=United%20States&table=DP05&tid=ACSDP1Y2018.DP05&g=0100000US&lastDisplayedRow=29&vintage=2017&layer=state&cid=DP05_0001E) (2018).
- <sup>15</sup> Underwood, D., Kim, H., & Matier, M. (2000). To Mail or To Web: Comparisons of Survey Response Rates and Respondent Characteristics. AIR 2000 Annual Forum Paper.
- <sup>16</sup> Collins, A.M. & Watson, C.T. Immunoglobulin light chain gene rearrangements, receptor editing and the development of a self-tolerant antibody repertoire. *Front Immunol.* **9**, 2249 (2018).
- <sup>17</sup> Visscher, P.M., Hill, W.G. & Wray, N.R. Heritability in the genomics era – concepts and misconceptions. *Nat Rev Genet.* **9**, 255-266 (2008).
- <sup>18</sup> Yang, J., Lee, H., Goddard, M.E. & Visscher, P.M. GCTA: a tool for genome-wide complex trait analysis. *Am J Hum Genet.* **88**, 76-82 (2011).
- <sup>19</sup> Yang, J., *et al.* Common SNPs explain a large proportion of heritability for human height. *Nat Genet.* **42**, 565-569 (2010).
- <sup>20</sup> Williams, F.M.K., *et al.* Self-reported symptoms of covid-19 including symptoms most predictive of SARS-CoV-2 infection, are heritable. *MedRxiv*. doi: <https://doi.org/10.1101/2020.04.22.20072124> (2020).
- <sup>21</sup> Zhang, K., *et al.* Structural-functional interactions of NS1-BP protein with the splicing and mRNA export machineries for viral and host gene expression. *Proc Nat Acad Sci.* **115**, E12218-E12227 (2018).
- <sup>22</sup> GTEx Consortium. Genetic effects on gene expression across human tissues. *Nature.* **550**, 204-213. (2017)
- <sup>23</sup> UCSC Genome Browser GTEx Track: Combined Expression QTLs from 44 Tissues from GTEx (midpoint release, V6) (IVNS1ABP). <https://genome.ucsc.edu/cgi->

---

[bin/hgc?hgsid=863508327\\_p8QQtWdOtziU6pA2hnivGQqAXIAK&c=chr1&l=185414331&r=185414832&o=185414581&t=185414582&g=gtxEqtlCluster&i=IVNS1ABP](https://www.ncbi.nlm.nih.gov/bioproject/528132) (2017).

<sup>24</sup> Thaventhiran, J.E.D., *et al.* Whole-genome sequencing of a sporadic primary immunodeficiency cohort. *Nature*. **583**, 90-95 (2020).

<sup>25</sup> Wojcechowskyj, J.A., *et al.* Quantitative phosphoproteomics reveals extensive cellular reprogramming during HIV-1 entry. *Cell Host Microbe*. **13**, 613-623 (2013).

<sup>26</sup> Xiao, R. & Boehnke, M. Quantifying and correcting for the winner's curse in genetic association studies. *Genet Epidemiol*. **33**, 453-462. (2009).

<sup>27</sup> COVID-19 Host Genetics Initiative. The COVID-19 Host Genetics Initiative, a global initiative to elucidate the role of host genetic factors in susceptibility and severity of the SARS-CoV-2 virus pandemic. *Eur J Hum Genet*. **28**, 715-718. (2020).

<sup>28</sup> Wright KM, *et al.* A prospective analysis of genetic variants associated with human lifespan. *G3 (Bethesda)*. **9**, 2863-2878. (2019).

<sup>29</sup> Ball, C.A., *et al.* AncestryDNA Ethnicity Estimate 2020 White Paper. [https://www.ancestrycdn.com/dna/static/pdf/whitepapers/Ethnicity2020\\_white%20paper.pdf](https://www.ancestrycdn.com/dna/static/pdf/whitepapers/Ethnicity2020_white%20paper.pdf) (2020).

<sup>30</sup> Ball, C.A., *et al.* AncestryDNA Matching White Paper. <https://www.ancestrycdn.com/support/us/2020/08/matchingwhitepaper.pdf> (2020).

<sup>31</sup> Abraham, G., Qiu, Y. & Inouye, M. FlashPCA2: principal component analysis of biobank-scale genotype datasets. *Bioinformatics*. **33**, 2776-2778 (2017)

<sup>32</sup> Loh, P.R., *et al.* Reference-based phasing using the Haplotype Reference Consortium panel. *Nat Genet*. **48**, 1443-1448 (2016).

<sup>33</sup> Loh, P.R. Eagle v2.4.1 user manual. <https://alkesgroup.broadinstitute.org/Eagle/> (2018).

<sup>34</sup> Center for Statistical Genetics. Minimac4. <https://genome.sph.umich.edu/wiki/Minimac4> (2019).

<sup>35</sup> COVID-19 Host Genetics Initiative. The COVID-19 Host Genetics Initiative, a global initiative to elucidate the role of host genetic factors in susceptibility and severity of the SARS-CoV-2 virus pandemic. *COVID-19 Host Genetics Pilot Analysis Plan v1.0* (2020).

<sup>36</sup> Willer, C. J., Li, Y., & Abecasis, G. R. (2010). METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics*, **26**, 2190-2191 (2010)

<sup>37</sup> Pe'er, I., Yelensky, R., Altshuler, D. & Daly, M.J. Estimation of the multiple testing burden for genomewide association studies of nearly all common variants. *Genet Epidemiol.* **32**, 381-385 (2008).

<sup>38</sup> LDlink. National Institutes of Health. U.S. Department of Health and Human Services. <https://ldlink.nci.nih.gov/?var1=rs17713054> (2020).

<sup>39</sup> Lee, S.R., Wray, N.R., Goddard, M.E. & Visscher, P.M. Estimating missing heritability for disease from genome-wide association studies. *Am J Hum Genet.* **88**, 294-305 (2011).