

# 1 Whole genome sequencing identifies multiple loci for critical 2 illness caused by COVID-19

3  
4 Athanasios Kousathanas<sup>‡,1</sup>, Erola Pairo-Castineira<sup>‡,2,3</sup>, Konrad Rawlik<sup>2</sup>, Alex Stuckey<sup>1</sup>, Christopher  
5 A Odhams<sup>1</sup>, Susan Walker<sup>1</sup>, Clark D Russell<sup>2,4</sup>, Tomas Malinauskas<sup>5</sup>, Jonathan Millar<sup>2</sup>, Katherine  
6 S Elliott<sup>5</sup>, Fiona Griffiths<sup>2</sup>, Wilna Oosthuyzen<sup>2</sup>, Kirstie Morrice<sup>6</sup>, Sean Keating<sup>7</sup>, Bo Wang<sup>2</sup>, Daniel  
7 Rhodes<sup>1</sup>, Lucija Klaric<sup>3</sup>, Marie Zechner<sup>2</sup>, Nick Parkinson<sup>2</sup>, Andrew D. Bretherick<sup>3</sup>, Afshan Siddiq<sup>1</sup>,  
8 Peter Goddard<sup>1</sup>, Sally Donovan<sup>1</sup>, David Maslove<sup>8</sup>, Alistair Nichol<sup>9</sup>, Malcolm G Semple<sup>10,11</sup>, Tala  
9 Zainy<sup>1</sup>, Fiona Maleady-Crowe<sup>1</sup>, Linda Todd<sup>1</sup>, Shahla Salehi<sup>1</sup>, Julian Knight<sup>5</sup>, Greg Elgar<sup>1</sup>, Georgia  
10 Chan<sup>1</sup>, Prabhu Arumugam<sup>1</sup>, Tom A Fowler<sup>1,12</sup>, Augusto Rendon<sup>1</sup>, Manu Shankar-Hari<sup>13</sup>, Charlotte  
11 Summers<sup>14</sup>, Charles Hinds<sup>15</sup>, Peter Horby<sup>16</sup>, Danny McAuley<sup>17,18</sup>, Hugh Montgomery<sup>19</sup>, Peter J.M.  
12 Openshaw<sup>20,21</sup>, Yang Wu, Jian Yang<sup>22</sup>, Paul Elliott<sup>23</sup>, Timothy Walsh<sup>7</sup>, GenOMICC Investigators,  
13 23andMe, Covid-19 Human Genetics Initiative, Angie Fawkes<sup>6</sup>, Lee Murphy<sup>6</sup>, Kathy Rowan<sup>24</sup>,  
14 Chris P Ponting<sup>3</sup>, Veronique Vitart<sup>3</sup>, James F Wilson<sup>3,25</sup>, Richard H Scott<sup>1,26</sup>, Sara Clohisey<sup>\*,2</sup>,  
15 Loukas Moutsianas<sup>\*,1</sup>, Andy Law<sup>\*,2</sup>, Mark J Caulfield<sup>\*,1,27</sup>, J. Kenneth Baillie<sup>\*,2,3,4,7</sup>.

16 ‡ - joint first authors

17 \* - joint last authors

18 <sup>1</sup>Genomics England, London UK

19 <sup>2</sup>Roslin Institute, University of Edinburgh, Easter Bush, Edinburgh, EH25 9RG, UK

20 <sup>3</sup>MRC Human Genetics Unit, Institute of Genetics and Molecular Medicine, University of Edinburgh,  
21 Western General Hospital, Crewe Road, Edinburgh, EH4 2XU, UK

22 <sup>4</sup>Centre for Inflammation Research, The Queen's Medical Research Institute, University of Edinburgh,  
23 47 Little France Crescent, Edinburgh, UK

24 <sup>5</sup>Wellcome Centre for Human Genetics, University of Oxford, Roosevelt Drive, Oxford, OX3 7BN,  
25 UK

26 <sup>6</sup>Edinburgh Clinical Research Facility, Western General Hospital, University of Edinburgh, EH4  
27 2XU, UK

28 <sup>7</sup>Intensive Care Unit, Royal Infirmary of Edinburgh, 54 Little France Drive, Edinburgh, EH16 5SA,  
29 UK

30 <sup>8</sup>Department of Critical Care Medicine, Queen's University and Kingston Health Sciences Centre,  
31 Kingston, ON, Canada

32 <sup>9</sup>Clinical Research Centre at St Vincent's University Hospital, University College Dublin, Dublin,  
33 Ireland

34 <sup>10</sup>NIHR Health Protection Research Unit for Emerging and Zoonotic Infections, Institute of Infection,  
35 Veterinary and Ecological Sciences University of Liverpool, Liverpool, L69 7BE, UK

36 <sup>11</sup>Respiratory Medicine, Alder Hey Children's Hospital, Institute in The Park, University of Liverpool,

37 Alder Hey Children's Hospital, Liverpool, UK  
38 <sup>12</sup>Test and Trace, the Health Security Agency, Department of Health and Social Care, Victoria St,  
39 London, UK  
40 <sup>13</sup>Department of Intensive Care Medicine, Guy's and St. Thomas NHS Foundation Trust, London,  
41 UK  
42 <sup>14</sup>Department of Medicine, University of Cambridge, Cambridge, UK  
43 <sup>15</sup>William Harvey Research Institute, Barts and the London School of Medicine and Dentistry,  
44 Queen Mary University of London, London EC1M 6BQ, UK  
45 <sup>16</sup>Centre for Tropical Medicine and Global Health, Nuffield Department of Medicine, University of  
46 Oxford, Old Road Campus, Roosevelt Drive, Oxford, OX3 7FZ, UK  
47 <sup>17</sup>Wellcome-Wolfson Institute for Experimental Medicine, Queen's University Belfast, Belfast,  
48 Northern Ireland, UK  
49 <sup>18</sup>Department of Intensive Care Medicine, Royal Victoria Hospital, Belfast, Northern Ireland, UK  
50 <sup>19</sup>UCL Centre for Human Health and Performance, London, W1T 7HA, UK  
51 <sup>20</sup>National Heart and Lung Institute, Imperial College London, London, UK  
52 <sup>21</sup>Imperial College Healthcare NHS Trust:London,London,UK  
53 <sup>22</sup>Westlake Laboratory of Life Sciences and Biomedicine, Hangzhou, Zhejiang 310024, China  
54 <sup>23</sup>Imperial College, London  
55 <sup>24</sup>Intensive Care National Audit & Research Centre, London, UK  
56 <sup>25</sup>Centre for Global Health Research, Usher Institute of Population Health Sciences and Informatics,  
57 Teviot Place, Edinburgh EH8 9AG, UK  
58 <sup>26</sup>Great Ormond Street Hospital, London UK  
59 <sup>27</sup>William Harvey Research Institute, Queen Mary University of London, Charterhouse Square,  
60 London EC1 6BQ

## 61 Abstract

62 Critical illness in COVID-19 is caused by inflammatory lung injury, mediated by the host immune  
63 system. We and others have shown that host genetic variation influences the development of illness  
64 requiring critical care<sup>1</sup> or hospitalisation<sup>2;3;4</sup> following SARS-Co-V2 infection. The GenOMICC  
65 (Genetics of Mortality in Critical Care) study recruits critically-ill cases and compares their genomes  
66 with population controls in order to find underlying disease mechanisms.

67 Here, we use whole genome sequencing and statistical fine mapping in 7,491 critically-ill cases  
68 compared with 48,400 population controls to discover and replicate 22 independent variants that  
69 significantly predispose to life-threatening COVID-19. We identify 15 new independent associations  
70 with critical COVID-19, including variants within genes involved in interferon signalling (*IL10RB*,  
71 *PLSCR1*), leucocyte differentiation (*BCL11A*), and blood type antigen secretor status (*FUT2*).  
72 Using transcriptome-wide association and colocalisation to infer the effect of gene expression  
73 on disease severity, we find evidence implicating expression of multiple genes, including reduced  
74 expression of a membrane flippase (*ATP11A*), and increased mucin expression (*MUC1*), in critical  
75 disease.

76 We show that comparison between critically-ill cases and population controls is highly efficient for  
77 genetic association analysis and enables detection of therapeutically-relevant mechanisms of disease.  
78 Therapeutic predictions arising from these findings require testing in clinical trials.

## 79 Introduction

80 Critical illness in COVID-19 is both an extreme disease phenotype, and a relatively homogeneous  
81 clinical definition including patients with hypoxaemic respiratory failure<sup>5</sup> with acute lung injury,<sup>6</sup>  
82 and excluding many patients with non-pulmonary clinical presentations<sup>7</sup> who are known to have  
83 divergent responses to therapy.<sup>8</sup> In the UK, the critically-ill patient group is younger, less likely  
84 to have significant comorbidity, and more severely affected than a general hospitalised cohort,<sup>5</sup>  
85 characteristics which may amplify observed genetic effects. In addition, since development of critical  
86 illness is in itself a key clinical endpoint for therapeutic trials,<sup>8</sup> using critical illness as a phenotype  
87 in genetic studies enables detection of directly therapeutically-relevant genetic effects.<sup>1</sup>

88 Using microarray genotyping in 2,244 cases, we previously reported that critical COVID-19 is  
89 associated with genetic variation in the host immune response to viral infection (*OAS1*, *IFNAR2*,  
90 *TYK2*) and the inflammasome regulator *DPP9*.<sup>1</sup> In collaboration with international groups, we  
91 recently extended these findings to include a variant near *TAC4* (rs77534576).<sup>2</sup> Several variants  
92 have been associated with milder phenotypes, such as the need for hospitalisation or management  
93 in the community, including the ABO blood type locus,<sup>4</sup> a pleiotropic inversion in chr17q21.31,<sup>9</sup>  
94 and associations in 5 additional loci including the T lymphocyte-associated transcription factor,  
95 *FOXP4*.<sup>2</sup> An enrichment of rare loss-of-function variants in candidate interferon signalling genes has  
96 been reported,<sup>3</sup> but this has yet to be replicated at genome-wide significance thresholds.<sup>10;11</sup>

97 We established a partnership between the GenOMICC Study and Genomics England to perform  
98 whole genome sequencing (WGS) to improve resolution and deepen fine-mapping of significant  
99 signals to enhance the biological insights into critical COVID-19. Here, we present results from a  
100 cohort of 7,491 critically-ill patients from 224 intensive care units, compared with 48,400 population  
101 controls, describing discovery and validation of 22 gene loci for susceptibility to life-threatening  
102 COVID-19.

## 103 Results

### 104 Study design

105 Cases were defined by the presence of COVID-19 critical illness in the view of the treating clinician -  
106 specifically, the need for continuous cardio-respiratory monitoring. Patients were recruited from  
107 224 intensive care units across the UK in the GenOMICC (Genetics Of Mortality In Critical Care)  
108 study. As a control population, unrelated participants recruited to the 100,000 Genomes Project  
109 were selected, excluding those with a known positive COVID-19 test, as severity information was  
110 not available. The 100,000 Genomes Project cohort (100k cohort) is comprised of UK individuals  
111 with a broad range of rare diseases or cancer and their family members. We included an additional  
112 prospectively-recruited cohort of volunteers (mild cohort) who self-reported testing positive for  
113 SARS-CoV-2 infection, and experienced mild or asymptomatic disease.

### 114 GWAS analysis

115 Whole genome sequencing and subsequent alignment and variant calling was performed for all  
116 subjects as described below (Methods). Following quality control procedures, we used a logistic  
117 mixed model regression, implemented in SAIGE,<sup>12</sup> to perform association analyses with unrelated

chr:pos (hg38)	rsid	REF	ALT	RAF	pop	OR	OR <sub>CI</sub>	Pval	HetPVal	Consequence	Gene	Expression
1:155066988	rs114301457	C	T*	0.0058	EUR	2.40	1.82-3.16	$6.8 \times 10^{-10}$	1	synonymous	EFNA4	-
1:155175305	rs7528026	G	A*	0.032	META	1.39	1.24-1.55	$7.16 \times 10^{-9}$	0.96	intron	TRIM46	-
1:155197995	rs41264915	A*	G	0.89	EUR	1.28	1.19-1.37	$1.02 \times 10^{-12}$	0.29	intron	THBS3	MUC1
2:00480453	rs1123573	A*	G	0.61	META	1.13	1.09-1.18	$9.85 \times 10^{-10}$	0.29	intron	BCL11A	-
3:45796521	rs2271616	G	T*	0.14	EUR	1.29	1.21-1.37	$9.9 \times 10^{-17}$	0.0011	5' UTR	SLC6A20	SLC6A20, CCR5
3:45859597	rs73064425	C	T*	0.077	EUR	2.71	2.51-2.94	$1.97 \times 10^{-133}$	0.010	intron	LZTFL1	LZTFL1, CCR9
3:146517122	rs343320	G	A*	0.081	EUR	1.25	1.16-1.35	$4.94 \times 10^{-9}$	0.53	missense	PLSCR1	-
5:131995059	rs56162149	C	T*	0.17	EUR	1.20	1.13-1.26	$7.65 \times 10^{-11}$	0.17	intron	ACSL6	ACSL6, FNIP1
6:32623820	rs9271609	T*	C	0.65	EUR	1.14	1.09-1.19	$3.26 \times 10^{-9}$	0.24	upstream	HLA-DQA1	HLA-DQA1, HLA-DQA2
6:41515007	rs2496644	A*	C	0.015	META	1.45	1.32-1.60	$7.59 \times 10^{-15}$	0.49	intron	LINC01276	-
9:21206606	rs28368148	C	G*	0.013	EUR	1.74	1.45-2.09	$1.93 \times 10^{-9}$	1	missense	IFNA10	-
11:34482745	rs61882275	G*	A	0.62	EUR	1.15	1.10-1.20	$1.61 \times 10^{-10}$	0.29	intron	ELF5	-
12:132489230	rs56106917	GC	G*	0.49	EUR	1.13	1.09-1.18	$2.08 \times 10^{-9}$	0.90	upstream	FBRSL1	-
13:112889041	rs9577175	C	T*	0.23	EUR	1.18	1.12-1.24	$3.71 \times 10^{-11}$	0.10	downstream	ATP11A	ATP11A
15:93046840	rs4424872	T*	A	0.0079	EUR	2.37	1.87-3.01	$8.61 \times 10^{-13}$	$1.82 \times 10^{-7}$	intron	RGMA	-
16:89196249	rs117169628	G	A*	0.15	EUR	1.19	1.12-1.26	$4.9 \times 10^{-9}$	0.80	missense	SLC22A31	SLC22A31, CDH15
17:46152620	rs2532300	T*	C	0.77	EUR	1.16	1.10-1.22	$4.19 \times 10^{-9}$	0.32	intron	KANSL1	ARHGAP27
17:49863260	rs3848456	C	A*	0.029	EUR	1.50	1.33-1.70	$4.19 \times 10^{-11}$	0.14	regulatory	.	-
19:4717660	rs12610495	A	G*	0.31	EUR	1.32	1.27-1.38	$3.91 \times 10^{-36}$	0.069	missense	DPP9	-
19:10305768	rs73510898	G	A*	0.093	EUR	1.28	1.19-1.37	$1.57 \times 10^{-11}$	0.011	intron	ZGLP1	-
19:10352442	rs34536443	G	C*	0.050	EUR	1.50	1.36-1.65	$6.98 \times 10^{-17}$	0.63	missense	TYK2	TYK2, PDE4A
19:48697960	rs368565	C	T*	0.44	EUR	1.15	1.1-1.2	$3.55 \times 10^{-11}$	0.22	intron	FUT2	FUT2, NTN5, RASIP1
21:33230000	rs17860115	C	A*	0.32	EUR	1.24	1.19-1.3	$9.69 \times 10^{-22}$	0.63	5' UTR	IFNAR2	-
21:33287378	rs8178521	C	T*	0.27	EUR	1.18	1.12-1.23	$3.53 \times 10^{-12}$	0.67	intron	IL1ORB	-
21:33959662	rs35370143	T	TAC*	0.083	EUR	1.26	1.17-1.36	$1.24 \times 10^{-9}$	1	intron	LINC00649	-

Table 1: Lead variants from independent regions in the per-population GWAS and trans-ancestry meta-analysis. Variants and the reference and alternate allele are reported with hg38 build coordinates. Asterisk (\*) indicates the risk allele. For each variant, we report the risk allele frequency in Europeans (RAF), the odds ratio and 95% confidence interval, and the association  $P$ -value. Consequence indicates the worst consequence predicted by VEP99, and Gene indicates the VEP99-predicted gene, but not necessarily the causal mediator. Expression indicates genes where is evidence of gene expression affecting COVID-19 severity, found by TWAS and colocalisation analysis.

118 individuals (critically-ill cases  $n = 7,491$ , controls (100k)  $n = 46,770$ , controls (mild COVID-  
119 19)  $n = 1,630$ ) (Methods, Supplementary Table 2). 1,339 of these cases were included in the  
120 primary analysis for our previous report.<sup>1</sup> Genome wide association studies (GWAS) were performed  
121 separately for genetically predicted ancestry groups (European - EUR, South Asian - SAS, African  
122 - AFR, East Asian - EAS, see Methods). Subsequently, we conducted inverse-variance weighted  
123 fixed effects meta-analysis across the four predicted ancestry cohorts using METAL<sup>13</sup> (Methods).  
124 In order to reduce the risk of spurious associations arising from genotyping or pipeline errors, we  
125 required supporting evidence from variants in linkage disequilibrium for all genome-wide significant  
126 variants: observed z-scores for each variant were compared to imputed z-scores for the same variant,  
127 with discrepant values being excluded (see Methods, Supplementary Figure 12).

128 In population-specific analyses, we discovered 22 independent genome-wide significant associations  
129 in the EUR ancestry group (Figure 1, Supplementary Figure 11 and Table 1) at a  $P$ -value threshold  
130 adjusted for multiple testing for 2,264,479 independent linkage disequilibrium-pruned genetic variants:  
131  $2.2 \times 10^{-08}$  (Supplementary Table 3). The strong association at 3p21.31 also reached genome-wide  
132 significance in the SAS ancestry group (Supplementary Figure 11).

133 In trans-ancestry meta-analysis, we identified an additional three loci with genome-wide significant  
134 associations (Figure 1, Table 1). We tested the meta-analysed set of 25 loci for heterogeneity of  
135 effect size between predicted ancestries and detected significant (at  $P < 1.83 \times 10^{-3}$ ) evidence for  
136 heterogeneity for two variants (Table 1, Supplementary Figure 13).

137 Fine mapping of the association signal revealed putative causal variants for several genes (See

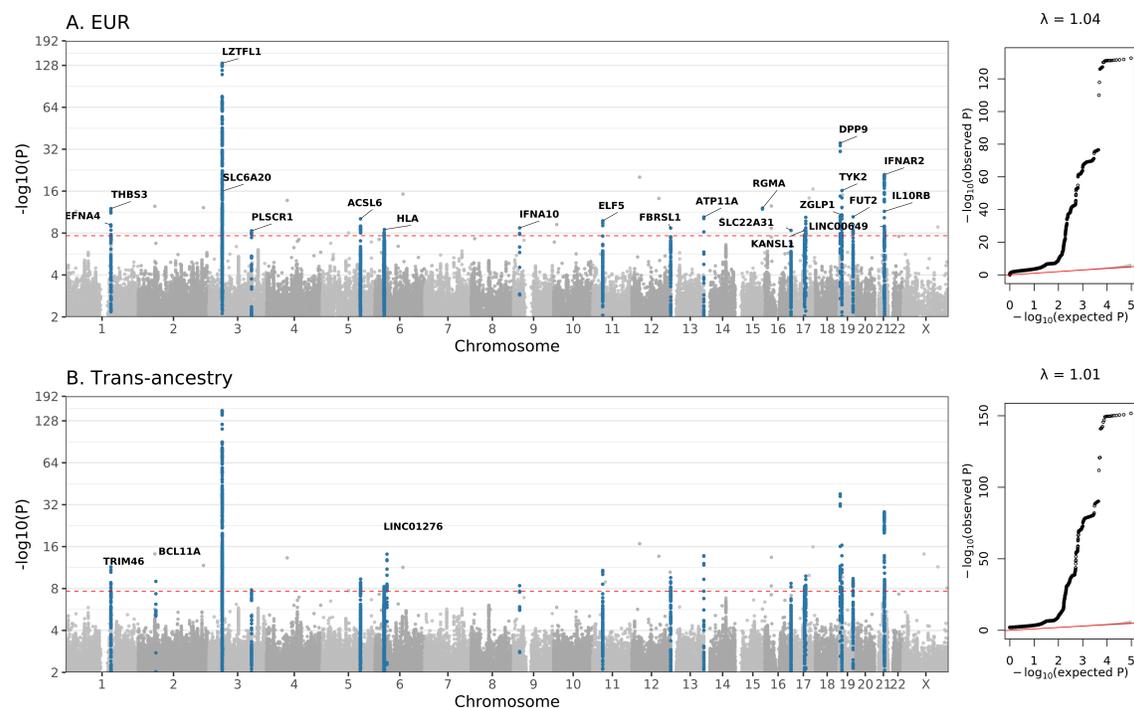


Figure 1: GWAS results for EUR ancestry group, and trans-ancestry meta-analysis. Manhattan plots are shown on the left and quantile–quantile (QQ) plots of observed versus expected  $P$  values are shown on the right, with genomic inflation ( $\lambda$ ) displayed for each analysis. Highlighted results in blue in the Manhattan plots indicate variants that are LD-clumped ( $r^2=0.1$ ,  $P_2=0.01$ , EUR LD) with the lead variants at each locus. Gene name annotation by Variant Effect Predictor (VEP) indicates genes impacted by the predicted consequence type of each lead variant. The dashed line shows the Bonferroni-corrected  $P$ -value= $2.2 \times 10^{-8}$ .

138 Supplementary Information). For example, we detected variants at 3q24 and 9p21.3 predicted to  
139 be missense mutations by Variant Effect Predictor (VEP). These impact *PLSCR1* and *IFNA10*  
140 respectively, and both are predicted to be deleterious by the Combined Annotation Dependent  
141 Depletion (CADD) tool<sup>14</sup> (*PLSCR1* (chr3:146517122:G:A, rs343320, p.His262Tyr, OR:1.24, 95% CIs  
142 [1.15-1.33], CADD:22.6; *IFNA10* (chr9:21206606:C:G, rs28368148, p.Trp164Cys, OR:1.74, 95% CIs  
143 [1.45-2.09], CADD:23.9). Structural predictions for these loci suggest functional effects (Figure 3  
144 and Supplementary Figure 15).

## 145 Replication

146 Replication was performed using summary statistics generously shared by collaborators: data from  
147 the COVID-19 Host Genetics Initiative (HGI) data freeze 6 were combined using meta-analysis  
148 with data shared by 23andMe (Methods). Although the HGI programme included an analysis  
149 intended to mirror the GenOMICC study (analysis "A2"), there are currently insufficient cases  
150 from other sources available to attempt replication, so we used the broader hospitalised phenotype  
151 (analysis "B2") for replication. We removed signals in the HGI data derived from GenOMICC cases  
152 using mathematical subtraction (see Methods) to ensure independence. Using LD clumping to find  
153 variants genotyped in both the discovery and replication studies, we required  $P < 0.002$  (0.05/25)  
154 and concordant direction of effect (Table 2) for replication.

155 We replicated 22 of the 25 significant associations identified in the population specific and/or  
156 trans-ancestry GWAS. Two of the three loci not replicated correspond to rare alleles that may not be  
157 well represented in the replication datasets which are dominated by SNP genotyping data. Although  
158 not replicated, for rs28368148 (9:21206606:C:G, *IFNA10*) we observed both a consistent direction  
159 of effect and odds ratio. The third locus is within the human leukocyte antigen (HLA) locus (see  
160 below).

161 We inferred credible sets of variants using Bayesian fine-mapping with susieR<sup>15</sup>, by analysing the  
162 GWAS summaries of 17 3Mbp regions that were flanking groups of lead signals. We obtained 22  
163 independent credible sets of variants for EUR and one for SAS that each had posterior inclusion  
164 probability  $> 0.95$ .

## 165 Gene burden testing

166 To assess the contribution of rare variants to critical illness, we performed gene-based analysis using  
167 SKAT-O as implemented in SAIGE-GENE<sup>16</sup>, using a subset of 12,982 individuals from our cohort  
168 (7,491 individuals with critical COVID-19 and 5,391 controls) for which the genome sequencing  
169 data were processed with the same alignment and variant calling pipeline. We tested the burden of  
170 rare (MAF<0.5%) variants considering the predicted variant consequence type. We assessed burden  
171 using a strict definition for damaging variants (high-confidence loss-of-function (pLoF) variants as  
172 identified by LOFTEE<sup>17</sup>) and a lenient definition (pLoF plus missense variants with CADD  $\geq 10$ )<sup>14</sup>  
173 , but found no significant associations at a gene-wide significance level. All individual rare variants  
174 included in the tests had  $P$ -values  $> 10^{-5}$ .

175 We then further examined the association with 13 genes involved in the regulation of type I and  
176 III interferon immunity that were implicated in critical COVID-19 pneumonia<sup>3</sup> but, as with other  
177 recent studies<sup>10</sup>, we did not find any significant gene burden test associations (tests for all genes

chr:pos (hg38)	rsid	REF	ALT	OR	OR <sub>CI</sub>	P <sub>val</sub>	OR <sub>HGI, 23m</sub>	OR <sub>CI, HGI, 23m</sub>	P <sub>val, HGI, 23m</sub>	Gene	Citation
1:155066988	rs114301457	C	T	2.40	1.81-3.18	1.51×10 <sup>-9</sup>	1.46	1.21-1.77	0.00011 *	EFNA4	-
1:155175305	rs7528026	G	A	1.39	1.24-1.55	7.16×10 <sup>-9</sup>	1.14	1.07-1.22	0.00012 *	TRIM46	-
1:155197995	rs41264915	A	G	0.80	0.76-0.86	3.79×10 <sup>-12</sup>	0.9	0.87-0.933	1.51×10 <sup>-9</sup> *	THBS3	-
2:60480453	rs1123573	A	G	0.88	0.85-0.92	9.85×10 <sup>-10</sup>	0.95	0.93-0.97	0.000018 *	BCL11A	-
3:45796521	rs2271616	G	T	1.26	1.19-1.34	2.45×10 <sup>-15</sup>	1.11	1.07-1.15	4.95×10 <sup>-9</sup> *	SLC6A20	( <sup>2</sup> )
3:45859597	rs73064425	C	T	2.52	2.35-2.70	2.18×10 <sup>-152</sup>	1.46	1.4-1.51	1.02×10 <sup>-77</sup> *	LZTFL1	<sup>4</sup>
3:146517122	rs343320	G	A	1.24	1.15-1.33	1.52×10 <sup>-8</sup>	1.08	1.04-1.13	0.00028 *	PLSCR1	-
5:132441275	rs10066378	T	C	1.20	1.13-1.27	4.48×10 <sup>-10</sup>	1.05	1.02-1.08	0.00074 *	IRF1-AS1	-
6:32623820	rs9271609	T	C	0.88	0.84-0.92	1.27×10 <sup>-8</sup>	1	0.98-1.03	0.89	HLA-DQA1	-
6:41515007	rs2496644	A	C	0.69	0.63-0.76	7.59×10 <sup>-15</sup>	0.87	0.83-0.92	3.17×10 <sup>-7</sup> *	LINC01276	-
9:21206606	rs28368148	C	G	1.74	1.45-2.1	4.09×10 <sup>-9</sup>	1.21	1.07-1.37	0.0024	IFNA10	-
11:34482745	rs61882275	G	A	0.87	0.84-0.91	1.62×10 <sup>-11</sup>	0.93	0.91-0.95	1.9×10 <sup>-10</sup> *	ELF5	*
12:132479205	rs4883585	G	A	1.13	1.09-1.18	1.12×10 <sup>-9</sup>	1.04	1.02-1.06	0.00047 *	FBRSL1	-
13:112889041	rs9577175	C	T	1.18	1.13-1.23	1.61×10 <sup>-12</sup>	1.07	1.04-1.09	1.29×10 <sup>-6</sup> *	ATP11A	-
15:93046840	rs4424872	T	A	0.64	0.53-0.76	1.99×10 <sup>-6</sup>	-	-	-	RGMA	-
16:89196249	rs117169628	G	A	1.18	1.12-1.25	6.04×10 <sup>-9</sup>	1.1	1.07-1.14	6.57×10 <sup>-9</sup> *	SLC22A31	-
17:46152620	rs2532300	T	C	0.87	0.82-0.91	1.4×10 <sup>-8</sup>	0.92	0.89-0.94	2.49×10 <sup>-9</sup> *	KANSL1	<sup>9</sup>
17:49863260	rs3848456	C	A	1.42	1.27-1.58	1.47×10 <sup>-10</sup>	1.15	1.09-1.21	1.34×10 <sup>-7</sup> *	.	<sup>2</sup>
19:4717660	rs12610495	A	G	1.32	1.27-1.38	6.44×10 <sup>-39</sup>	1.11	1.09-1.14	5.74×10 <sup>-19</sup> *	DPP9	<sup>1</sup>
19:10305768	rs73510898	G	A	1.24	1.16-1.33	1.47×10 <sup>-9</sup>	1.08	1.04-1.12	0.00016 *	ZGLP1	-
19:10352442	rs34536443	G	C	1.50	1.37-1.66	4.22×10 <sup>-17</sup>	1.22	1.15-1.29	4.06×10 <sup>-11</sup> *	TYK2	<sup>1</sup>
19:48697960	rs368565	C	T	1.13	1.09-1.18	3.74×10 <sup>-10</sup>	1.04	1.02-1.06	0.00087 *	FUT2	-
21:33230000	rs17860115	C	A	1.26	1.21-1.31	6.28×10 <sup>-28</sup>	1.11	1.08-1.13	1.77×10 <sup>-18</sup> *	IFNAR2	<sup>1</sup>
21:33287378	rs8178521	C	T	1.17	1.12-1.22	4.23×10 <sup>-12</sup>	1.06	1.03-1.09	8.02×10 <sup>-6</sup> *	IL10RB	-
21:33914436	rs12626438	A	G	1.22	1.14-1.31	1.78×10 <sup>-8</sup>	1.1	1.06-1.14	2.33×10 <sup>-7</sup> *	LINC00649	-

Table 2: Replication in a combined data from external studies - combined meta-analysis of HGI freeze 6 B2 and 23andMe. Odds ratios and *P*-values are shown for variants in LD with the lead variant that were genotyped/imputed in both sources. Chromosome, reference and alternate allele correspond to the build hg38. An asterisk (\*) next to the HGI and 23andme meta-analysis *P*-value indicates that the lead signal is replicated with *P*-value<0.002 with a concordant direction of effect. Citation lists the first publication of confirmed genome-wide associations with critical illness or (in brackets) any COVID-19 phenotype; in this column, (\*) indicates a variant which met genome-wide significance without GenOMICC data in the public latest version of the HGI analysis (B2 v6) but has not been reported yet.

178 had  $P$ -value  $> 0.05$ , Supplementary File AVTsuppinfo.xlsx). We also did not replicate the reported  
 179 association<sup>10</sup> for the toll-like receptor 7 (*TLR7*) gene.

## 180 Transcriptome-wide association study

181 In order to infer the effect of genetically-determined variation in gene expression on disease sus-  
 182 ceptibility, we performed a transcriptome-wide association study (TWAS) using gene expression  
 183 data (GTEXv8) for two disease-relevant tissues, lung and whole blood. We found 14 genes with  
 184 significant association between predicted expression and critical COVID-19 in the lung and 6 in  
 185 whole blood analyses (Supplementary File: TWAS.xlsx). To increase statistical power using eQTLs  
 186 from multiple tissues, we performed a TWAS meta-analysis using all available tissues in GTEXv8,  
 187 revealing 51 transcriptome-wide significant genes. Since TWAS uses a composite signal derived  
 188 from multiple eQTLs, we used colocalisation to find specific eQTLs in whole blood (eqtlGen and  
 189 GTEXv8) and lung (GTEXv8<sup>18</sup>) which share the same signal with GWAS (EUR) associations. We  
 190 found 16 genes which significantly colocalise in at least one of the studied tissues, shown in Figure 2.

191 We repeated the TWAS analysis using models of intron excision rate from GTEXv8 to obtain splicing  
 192 TWAS. We found 40 signals in lung, affecting 16 genes and 20 signals in whole blood which affect  
 193 9 genes. In a meta-analysis of splicing TWAS using all GTEXv8 tissues, we found 91 significant  
 194 introns in a total of 33 genes. Using GTEXv8 lung and whole blood sqtIs to find colocalising  
 195 signals with splicing TWAS significant results, we found 11 genes with colocalising splicing signals  
 196 (Supplementary File: TWAS.xlsx).

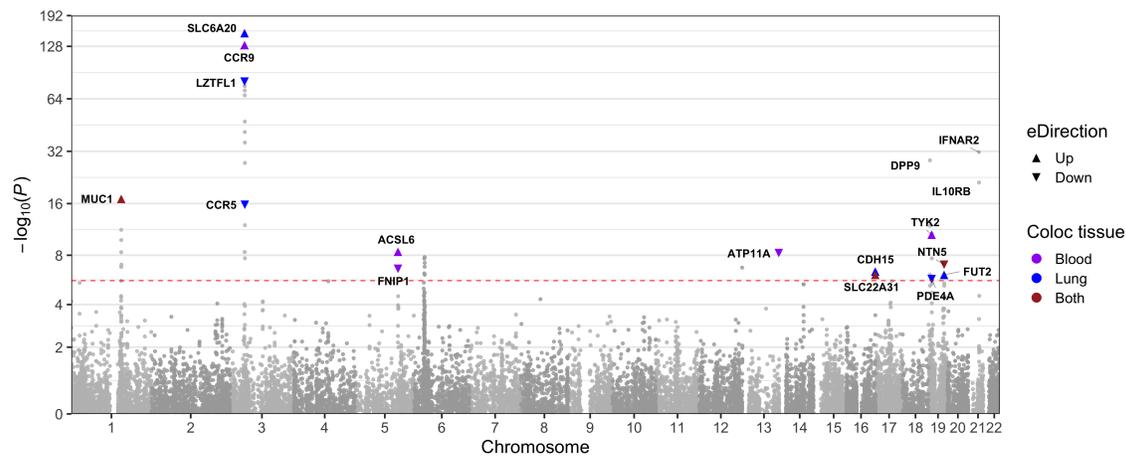


Figure 2: Gene-level Manhattan plot showing results from TWAS meta-analysis and highlighting genes that colocalise with GWAS signals or have strong metaTWAS associations. Highlighting color is different for lung and blood tissue data that were used for colocalisation. Arrows show direction of change in gene expression associated with an increased disease risk. Red dashed line shows significance threshold at  $P < 2.3 \times 10^{-6}$ .

## 197 HLA region

198 To investigate the contribution of specific HLA alleles to the observed association in the HLA region,  
199 we imputed HLA alleles at a four digit (two-field) level using HIBAG<sup>19</sup>. The only allele that reached  
200 genome-wide significance was HLA-DRB1\*04:01 ( $OR = 0.80$ ,  $95\%CI = 0.75 - 0.86$ ,  $P = 1.6 \times 10^{-10}$   
201 in EUR), which has a stronger  $P$ -value than the lead SNP in the region ( $OR : 0.88$ ,  $95\%CIs : 0.84 - 0.92$ ,  
202  $P = 3.3 \times 10^{-9}$  in EUR) and is a better fit to the data ( $AIC_{DRB1*04:01} = 30241.34$ ,  
203  $AIC_{leadSNP} = 30252.93$ ). Results are shown in supplementary figure 25.

## 204 Discussion

205 We report 22 replicated genetic associations with life-threatening COVID-19, and 3 additional loci,  
206 discovered in only 7,491 cases. This demonstrates the efficiency of the design of the GenOMICC  
207 study, which is an open-source international research programme<sup>20</sup> focusing on critically-ill patients  
208 with infectious disease and other critical illness phenotypes (<https://genomicc.org>). By using whole  
209 genome sequencing we were able to detect multiple distinct signals with high confidence for several  
210 of the associated loci, in some cases implicating different biological mechanisms.

211 Several variants associated with life-threatening disease are linked to interferon signalling. A coding  
212 variant in a ligand, *IFNA10A*, and reduced expression of its receptor *IFNAR2* (Figure 2), were  
213 associated with critical COVID-19. The narrow failure of replication for the *IFNA10* variant  
214 (rs28368148, replication  $P = 0.00243$ , significance threshold  $P < 0.002$ ) may be due to limited power  
215 in the replication cohort. The lead variant in *TYK2* in whole genome sequencing is a well-studied  
216 protein-coding variant with reduced phosphorylation activity, consistent with that reported recently,<sup>2</sup>  
217 but associated with significantly increased *TYK2* expression (Figure 2, Methods). Fine mapping  
218 reveals a significant critical illness association with an independent missense variant in *IL10RB*,  
219 a receptor for Type III (lambda) interferons (rs8178521, Trp164Cys, Table 1). Overall, variants  
220 predicted to be associated with reduction in interferon signalling are associated with critical disease.  
221 Importantly, systemic administration of interferon in a large clinical trial, albeit late in disease, did  
222 not reduce mortality.<sup>21</sup>

223 Phospholipid scramblase 1 (*PLSCR1*; chr3:146517122:G:A) functions as a nuclear signal for the  
224 antiviral effect of interferon,<sup>22</sup> and has been shown to control replication of other RNA viruses  
225 including vesicular stomatitis virus, encephalomyocarditis virus and Influenza A virus.<sup>23;22</sup> The risk  
226 allele at the lead variant (chr3:146517122:G:A, rs343320) encodes a substitution, H262Y, which  
227 is predicted to disrupt the non-canonical nuclear localisation signal<sup>24</sup> by eliminating a hydrogen  
228 bond with importin (Figure 3). Deletion of this nuclear localisation signal has been shown to  
229 prevent neutrophil maturation.<sup>25</sup> Although *PLSCR1* is strongly up-regulated when membrane lipid  
230 asymmetry is lost (see below), it may not act directly on this process.<sup>26</sup>

231 We report significant associations in several genes implicated in B-cell lymphopoiesis and differentia-  
232 tion of myeloid cells. *BCL11A* is essential in B- and T-lymphopoiesis<sup>27</sup> and promotes plasmacytoid  
233 dendritic cell differentiation.<sup>28</sup> *TAC4*, reported previously,<sup>2</sup> encodes a regulator of B-cell lymphopoe-  
234 sis<sup>29</sup> and antibody production,<sup>30</sup> and promotes survival of dendritic cells.<sup>31</sup> Finally, although  
235 the strongest fine mapping signal at 5q31.1 (chr5:131995059:C:T, rs56162149) is in an intron of  
236 *ACSL6* (locus, p), the credible set includes a missense variant in *CSF2* of uncertain significance  
237 (chr5:132075767:T:C). *CSF2* encodes granulocyte-macrophage colony stimulating factor, a key

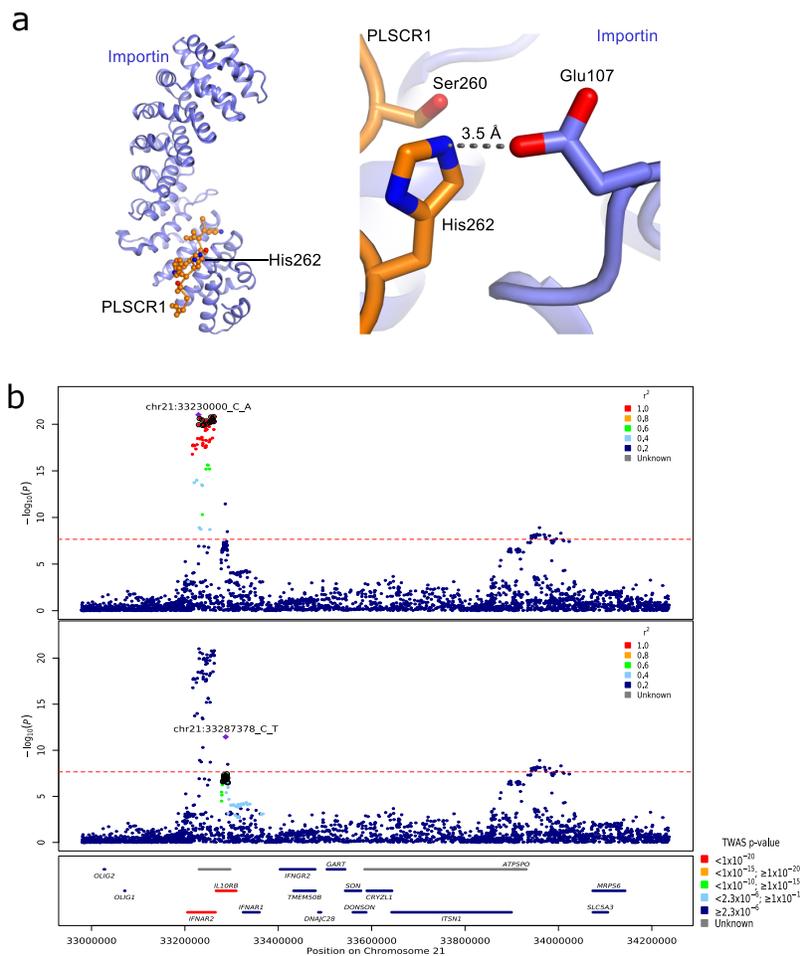


Figure 3: (a) Predicted structural consequences of lead variant at *PLSCR1*. Left panel shows the crystal structure of PLSCR1 nuclear localization signal (orange, Gly257–Ile266, numbering correspond to UniProt entry O15162) in complex with Importin  $\alpha$  (blue), Protein Data Bank (PDB) ID 1Y2A. Side chains of PLSCR1 are shown as connected spheres with carbon atoms coloured in orange, nitrogens in blue and oxygens in red. Hydrogen atoms were not determined at this resolution (2.20 Å) and are not shown. Right panel: a closeup view showing side chains of PLSCR1 Ser260, His262 and Importin Glu107 as sticks. Distance (in Å) between selected atoms (PLSCR1 His262 Ne2 and Importin Glu107 carboxyl O) is indicated. A hydrogen bond between PLSCR1 His262 and Importin Glu107 is indicated with a dashed line. The risk variant is predicted to eliminate this bond, disrupting nuclear import, an essential step for effect on antiviral signalling<sup>22</sup> and neutrophil maturation.<sup>25</sup> (b) Regional detail showing fine-mapping to separate two adjacent independent signals. Top two panels: variants in linkage disequilibrium with the lead variants shown. The loci that are included in two independent credible sets are displayed with black outline circles. Bottom panel: locations of protein-coding genes, coloured by TWAS  $P$ -value.

238 differentiation factor in the mononuclear phagocyte system which is strongly up-regulated in critical  
239 COVID-19,<sup>32</sup> and is already under investigation as a target for therapy.<sup>33</sup>

240 Several new genetic associations implicate genes known to be involved in lung disease. The second  
241 variant in the credible set at 13q14 (chr13:112882313:A:G, rs1278769, in *ATP11A*), has been reported  
242 as a lead variant for idiopathic pulmonary fibrosis.<sup>34</sup> *ATP11A* encodes a flippase which maintains  
243 the asymmetric distribution of phospholipids in cell membranes;<sup>35</sup> disruption of this asymmetry  
244 is a phagocytic signal on apoptotic cells, and is required for platelet activation.<sup>36;37</sup> TWAS and  
245 colocalisation demonstrate that genetic variants predicted to decrease expression of *ATP11A* in lung  
246 are associated with critical illness. A combination of fine mapping, colocalisation with eQTL signals  
247 (GTEx and eQTLgen) and TWAS results provide evidence in support of *MUC1* as the mediator of  
248 the association with rs41264915 (Table 1). This may indicate an important role for mucins in the  
249 development of critical illness in COVID-19. The direction of effect (Figure 2) suggests that agents  
250 that reduce *MUC1* expression, and by extension its abundance, may be a therapeutic option. Finally,  
251 the association on 11p13 (rs61882275) includes GTEx eQTL for the lung fibroblast transcription  
252 factor *ELF5* in lung tissue, and the gene encoding the antioxidant enzyme catalase (*CAT*) in whole  
253 blood with evidence of colocalisation in both signals ( supplementary material: TWAS.xlsx).<sup>18</sup> The  
254 protective allele at this locus is weakly associated with reduced lung function in a previous GWAS.<sup>38</sup>

255 *FUT2* encodes alpha-(1,2)fucosyltransferase, which controls the secretion of ABO blood type glycans  
256 into body fluids and expression on epithelial surfaces. An association with critical COVID-19 was  
257 reported previously in a candidate gene association study by Mankelov et al.<sup>39</sup> The credible set for the  
258 *FUT2* locus includes rs492602 (chr19:48703160:A:G) which is linked to a stop codon gain mutation  
259 (chr19:48703417:G:A), leading to the well-described non-secretor phenotype in homozygotes.<sup>40;41</sup>  
260 We show that the stop-gain, non-secretor allele is protective against life-threatening COVID-19.  
261 The protective variant in our study has been previously reported to protect against other viruses  
262 (rotavirus,<sup>42</sup> mumps and common colds<sup>43</sup>), to enhance antibody responses to polyomavirus BK<sup>44</sup>  
263 and to increase susceptibility to infection with some encapsulated bacteria.<sup>45</sup>

## 264 Limitations

265 In contrast to microarray genotyping, whole genome sequencing is rapidly evolving and a relatively  
266 new technology for genome-wide association studies, with relatively few sources of population  
267 controls. We used selected controls from the 100,000 genomes project, sequenced on a different  
268 platform (illumina HiSeqX) from the cases (illumina NovaSeq6000)(Supplementary Table 1). To  
269 minimise the risk of false positive associations arising due to sequencing or genotyping errors, we  
270 required all significant associations to be supported by local variants in linkage disequilibrium,  
271 which may be excessively stringent (see Methods). Although this approach may remove some true  
272 associations, our priority is to maximise confidence in the reported signals. Of 25 variants meeting  
273 this requirement, 22 are replicated in an independent study, and the remaining 3 may well be true  
274 associations that have failed due to a lack of coverage or power in the replication dataset.

275 The design of our study incorporates genetic signals for every stage in the disease progression  
276 into a single phenotype. This includes exposure, viral replication, inflammatory lung injury and  
277 hypoxaemic respiratory failure. Although we can have considerable confidence that the replicated  
278 associations with critical COVID-19 we report are robust, we cannot determine at which stage in  
279 the disease process, or in which tissue, the relevant biological mechanisms are active, which can have  
280 therapeutic implications.

## 281 **Conclusions**

282 The genetic associations here implicate new biological mechanisms underlying the development of  
283 life-threatening COVID-19, several of which may be amenable to therapeutic targeting. In the  
284 context of the ongoing global pandemic, translation to clinical practice is an urgent priority. As  
285 with our previous work, large-scale randomised trials are essential before translating our findings  
286 into clinical practice.

## 287 Acknowledgements

288 We thank the patients and their loved ones who volunteered to contribute to this study at one of the  
289 most difficult times in their lives, and the research staff in every intensive care unit who recruited  
290 patients at personal risk during the most extreme conditions we have ever witnessed in UK hospitals.

291 GenOMICC was funded by the Department of Health and Social Care (DHSC), LifeArc, the Medical  
292 Research Council, UKRI, Sepsis Research (the Fiona Elizabeth Agnew Trust), the Intensive Care  
293 Society, a Wellcome-Beit Prize award to J. K. Baillie (Wellcome Trust 103258/Z/13/A) and a BBSRC  
294 Institute Program Support Grant to the Roslin Institute (BBS/E/D/20002172, BBS/E/D/10002070  
295 and BBS/E/D/30002275). Whole-genome sequencing was performed by Illumina at the NHS  
296 Genomic Sequencing Centre in partnership and was overseen by Genomics England. We would  
297 like to thank all at Genomics England who have contributed to the supporting the processing of  
298 the sequencing and clinical data. We thank DHSC, the Medical Research Council, UKRI, LifeArc,  
299 Genomics England Ltd and Illumina Inc for funding sequencing. Genomics England and the 100,000  
300 Genomes Project was funded by the National Institute for Health Research, the Wellcome Trust,  
301 the Medical Research Council, Cancer Research UK, the Department of Health and Social Care  
302 and NHS England. We are grateful for the support from Professor Dame Sue Hill and the team  
303 in NHS England and the 13 Genomic Medicine Centres that successfully delivered the 100,000  
304 Genomes Project which provide the control sequences for this study. We thank the participants of  
305 the 100,000 Genomes Project who made this study possible and the Genomics England Participant  
306 Panel for their strategic advice, involvement and engagement. We acknowledge NHS Digital, Public  
307 Health England and the Intensive Care National Audit and Research Centre who provided life course  
308 longitudinal clinical data on the participants. This work forms part of the portfolio of research of  
309 the NIHR Biomedical Research Centre at Barts. Mark Caulfield is an NIHR Senior Investigator.  
310 This study owes a great deal to the National Institute of Healthcare Research Clinical Research  
311 Network (NIHR CRN) and the Chief Scientist Office (Scotland), who facilitate recruitment into  
312 research studies in NHS hospitals, and to the global ISARIC and InFACT consortia.

313 The views expressed are those of the authors and not necessarily those of the DHSC, DID, NIHR,  
314 MRC, Wellcome Trust or PHE.

315 The Genotype-Tissue Expression (GTEx) Project was supported by the Common Fund of the Office  
316 of the Director of the National Institutes of Health, and by NCI, NHGRI, NHLBI, NIDA, NIMH, and  
317 NINDS. The data used for the analyses described in this manuscript were obtained from the GTEx  
318 Portal on August 22nd, 2021 (GTEx Analysis Release V8 (dbGaP Accession phs000424.v8.p2)).

## 319 Data availability

320 Summary statistics will be shared openly with international collaborators to accelerate discovery.  
321 Data can be obtained from [genomicc.org/data](https://genomicc.org/data)

322 Individual-level data will be available in the UK Outbreak Analysis Platform at the University of  
323 Edinburgh and through the Genomics England research environment.

## 324 **Contributions**

325 AK, EP-C, KR, AS, CAO, SW, TM, KSE, BW, DR, LK, MZ, NP, ADB, YW, JY, SC, LMo, AL and  
326 JKB contributed to data analysis. AK, EP-C, KR, AS, CAO, SW, CDR, JM, AR, SC, LMo and AL  
327 contributed to bioinformatics. AK, EP-C, KR, CDR, JM, DM, AN, MGS, SC, LMo, MJC and JKB  
328 contributed to writing and reviewing the manuscript. EP-C, KR, KM, SK, AF, LM, KRo, CPP,  
329 VV, JFW, SC, AL, MJC and JKB contributed to design. SW, FG, WO, PG and SD contributed to  
330 project management. FG, WO, KM, SK, PG, SD, DM, AN, MGS, SS, JK, TAF, MS-H, CS, CH,  
331 PH, DMc, HM, PJO, PE, TW, AF, LM, KRo, CPP, RHS, SC and AL contributed to oversight.  
332 FG, WO, FM-C and JKB contributed to ethics and governance. KM, ASi, AF and LM contributed  
333 to sample handling and sequencing. and ASi contributed to data collection. and TZ contributed  
334 to sample handling. TZ and GE contributed to sequencing. and LT contributed to recruitment of  
335 controls. GC, PA, KRo and AL contributed to clinical data management. KRo, CPP, SC and JKB  
336 contributed to conception. KRo, CPP, VV and JFW contributed to reviewing the manuscript. MJC  
337 and JKB contributed to scientific leadership.

## 338 **Conflict of interest**

339 All authors declare that they have no conflicts of interest relating to this work.

340 Genomics England Ltd is a wholly owned Department of Health and Social Care company created in  
341 2013 to work with the NHS to introduce advanced genomic technologies and analytics into healthcare.  
342 All Genomics England affiliated authors are, or were, salaried by Genomics England during this  
343 programme.

## 344 **Materials and Methods**

### 345 **Ethics**

346 GenOMICC was both approved by the following research ethics committees: Scotland "A" Research  
347 Ethics Committee, 15/SS/0110; Coventry and Warwickshire Research Ethics Committee (England,  
348 Wales and Northern Ireland), 19/WM/0247). Current and previous versions of the study protocol  
349 are available at [genomicc.org/protocol](http://genomicc.org/protocol). All participants gave informed consent.

### 350 **Recruitment of cases**

351 Patients recruited to the GenOMICC study ([genomicc.org](http://genomicc.org)) had confirmed COVID-19 according to  
352 local clinical testing and were deemed, in the view of the treating clinician, to require continuous  
353 cardiorespiratory monitoring. In UK practice this kind of monitoring is undertaken in high-  
354 dependency or intensive care units. This study was approved by research ethics committees in the  
355 recruiting countries (Scotland 15/SS/0110, England, Wales and Northern Ireland: 19/WM/0247).  
356 Current and previous versions of the study protocol are available at [genomicc.org/protocol](http://genomicc.org/protocol). All  
357 participants gave informed consent.

## 358 **Recruitment of controls**

### 359 **Mild/asymptomatic controls**

360 Participants were recruited to the mild COVID-19 cohort on the basis of having experienced mild  
361 (non-hospitalised) or asymptomatic COVID-19. Participants volunteered to take part in the study  
362 via a microsite and were required to self-report the details of a positive COVID-19 test. Volunteers  
363 were prioritised for genome sequencing based on demographic matching with the critical COVID-19  
364 cohort considering self-reported ancestry, sex, age and location within the UK. We refer to this  
365 cohort as the covid-mild cohort.

### 366 **100,000 Genomes project controls**

367 Participants were enrolled in the 100,000 Genomes Project from families with a broad range of  
368 rare diseases, cancers and infection by 13 regional NHS Genomic Medicine Centres across England  
369 and in Northern Ireland, Scotland and Wales. For this analysis, participants for whom a positive  
370 SARS-CoV-2 test had been recorded as of March, 2021 were not included due to uncertainty in the  
371 severity of COVID-19 symptoms. Only participants for whom genome sequencing was performed  
372 from blood derived DNA were included and participants with haematological malignancies were  
373 excluded to avoid potential tumour contamination.

## 374 **DNA extraction**

375 For severe COVID-19 cases and mild cohort controls, DNA was extracted from whole blood either  
376 manually using Nucleon Kit (Cytiva) and re-suspended in 1 ml TE buffer pH 7.5 (10mM Tris-Cl pH  
377 7.5, 1mM EDTA pH 8.0), or automated on the Chemagic 360 platform using Chemagic DNA blood  
378 kit (Perkin Elmer) and re-suspended in 400 $\mu$ L Elution Buffer. The yield of the DNA was measured  
379 using Qubit and normalised to 50ng/ $\mu$ l before sequencing.

## 380 **WGS sequencing**

381 For all three cohorts, DNA was extracted from whole-blood using standard protocols. Sequencing  
382 libraries were generating using the Illumina TruSeq DNA PCR-Free High Throughput Sample  
383 Preparation kit and sequenced with 150bp paired-end reads in a single lane of an Illumina HiSeq  
384 X instrument (for 100,000 Genomes Project samples) or NovaSeq instrument (for the COVID-19  
385 critical and mild cohorts).

## 386 **Sequencing data QC**

387 All genome sequencing data were required to meet minimum quality metrics and quality control  
388 measures were applied for all genomes as part of the bioinformatics pipeline. The minimum data  
389 requirements for all genomes were  $> 85 \times 10^{-9}$  bases with  $Q \geq 30$  and  $\geq 95\%$  of the autosomal  
390 genome covered at  $\geq 15x$  calculated from reads with mapping quality  $> 10$  after removing duplicate  
391 reads and overlapping bases, after adaptor and quality trimming. Assessment of germline cross-  
392 sample contamination was performed using VerifyBamID and samples with  $> 3\%$  contamination  
393 were excluded. Sex checks were performed to confirm that the sex reported for a participant was  
394 concordant with the sex inferred from the genomic data.

## 395 **WGS Alignment and variant calling**

### 396 **COVID-19 cohorts**

397 For the critical and mild COVID-19 cohorts, sequencing data alignment and variant calling was  
398 performed with Genomics England pipeline 2.0 which uses the DRAGEN software (v3.2.22). Align-  
399 ment was performed to genome reference GRCh38 including decoy contigs and alternate haplotypes  
400 (ALT contigs), with ALT-aware mapping and variant calling to improve specificity.

### 401 **100,000 Genome Project cohort (100K-genomes)**

402 All genomes from the 100,000 Genomes Project cohort were analysed with the Illumina North Star  
403 Version 4 Whole Genome Sequencing Workflow (NSV4, version 2.6.53.23); which is comprised of  
404 the iSAAC Aligner (version 03.16.02.19) and Starling Small Variant Caller (version 2.4.7). Samples  
405 were aligned to the Homo Sapiens NCBI GRCh38 assembly with decoys.

406 A subset of the genomes from the Cancer program of the 100,000 Genomes Project were reprocessed  
407 (alignment and variants calling) using the same pipeline used for the COVID-19 cohorts (DRAGEN  
408 v3.2.22) for equity of alignment and variant calling.

### 409 **Aggregation**

410 Aggregation was conducted separately for the samples analysed with Genomics England pipeline 2.0  
411 (severe-cohort, mild-cohort, cancer-realigned-100K), and those analysed with the Illumina North  
412 Star Version 4 pipeline (100K-Genomes).

413 For the first three, the WGS data were aggregated from single sample gVCF files to multi-sample  
414 VCF files using GVCFFGenotyper (GG) v3.8.1, which accepts gVCF files generated via the DRAGEN  
415 pipeline as input. GG outputs multi-allelic variants (several ALT variants per position on the same  
416 row), and for downstream analyses the output was decomposed to bi-allelic variants per row using  
417 software vt v0.57721. We refer to the aggregate as aggCOVID\_vX, where X is the specific freeze.  
418 The analysis in this manuscript uses data from freeze v4.2 and the respective aggregate is referred  
419 to as aggCOVID\_v4.2.

420 Aggregation for the 100K-Genomes cohort was performed using Illumina's gvcfgenotyper v2019.02.26,  
421 merged with bcftools v1.10.2 and normalised with vt v0.57721.

### 422 **Sample Quality Control (QC)**

423 Samples that failed any of the following four BAM-level QC filters: freemix contamination (>3%),  
424 mean autosomal coverage (<25X), percent mapped reads (<90%), and percent chimeric reads (>5%)  
425 were excluded from the analysis.

426 Additionally, a set of VCF-level QC filters were applied post-aggregation on all autosomal bi-allelic  
427 SNVs (akin to gnomAD v3.1<sup>17</sup>). Samples were filtered out based on the residuals of eleven QC metrics  
428 (calculated using bcftools) after regressing out the effects of sequencing platform and the first three  
429 ancestry assignment principal components (including all linear, quadratic, and interaction terms)  
430 taken from the sample projections onto the SNP loadings from the individuals of 1000 Genomes  
431 Project phase 3 (1KGP3). Samples were removed that were four median absolute deviations (MADs)

432 above or below the median for the following metrics: ratio heterozygous-homozygous, ratio insertions-  
433 deletions, ratio transitions-transversions, total deletions, total insertions, total heterozygous snps,  
434 total homozygous snps, total transitions, total transversions. For the number of total singletons  
435 (snps), samples were removed that were more than 8 MADs above the median. For the ratio of  
436 heterozygous to homozygous alternate snps, samples were removed that were more than 4 MADs  
437 above the median.

438 After quality control, 79,803 individuals were included in the analysis with the breakdown according  
439 to cohort shown in Supplementary Table 2.

## 440 Selection of high-quality (HQ) independent SNPs

441 We selected high-quality independent variants for inferring kinship coefficients, performing PCA,  
442 assigning ancestry and for the conditioning on the Genetic Relatedness matrix by the logistic mixed  
443 model of SAIGE and SAIGE-GENE. To avoid capturing platform and/or analysis pipeline effects  
444 for these analyses, we performed very stringent variant QC as described below.

### 445 HQ common SNPs

446 We started with autosomal, bi-allelic SNPs which had frequency  $> 5\%$  in aggV2 (100K participant  
447 aggregate) and in the 1KGP3. We then restricted to variants that had missingness  $< 1\%$ , median  
448 genotype quality  $QC > 30$ , median depth (DP)  $\geq 30$  and  $\geq 90\%$  of heterozygote genotypes passing  
449 an ABRatio binomial test with  $P$ -value  $> 10^{-2}$  for aggV2 participants. We also excluded variants in  
450 complex regions from the list available in , and variants where the ref/alt combination was CG or AT  
451 (C/G, G/C, A/T, T/A). We also removed all SNPs which were out of Hardy Weinberg Equilibrium  
452 (HWE) in any of the AFR, EAS, EUR or SAS super-populations of aggV2, with a  $P$ -value cutoff of  
453  $pHWE < 10^{-5}$ . We then LD-pruned using plink v1.9 with an  $r^2 = 0.1$  and in 500kb windows. This  
454 resulted in a total of 63,523 high-quality sites from aggV2.

455 We then extracted these high-quality sites from the aggCOVID\_v4.2 aggregate and further applied  
456 variant quality filters (missingness  $< 1\%$ , median  $QC > 30$ , median depth  $\geq 30$  and  $\geq 90\%$  of  
457 heterozygote genotypes passing an ABRatio binomial test with  $P$ -value  $> 10^{-2}$ ), per batch of  
458 sequencing platform (i.e, HiseqX, NovaSeq6000).

459 After applying variant filters in aggV2 and aggCOVID\_v4.2, we merged the genomic data from the  
460 two aggregates for the intersection of the variants which resulted in a final total of 58,925 sites.

### 461 HQ rare SNPs

462 We selected high-quality rare ( $MAF < 0.005$ ) bi-allelic SNPs to be used with SAIGE for aggregate  
463 variant testing analysis. To create this set, we applied the same variant QC procedure as with  
464 the common variants: We selected variants that had missingness  $< 1\%$ , median  $QC > 30$ , median  
465 depth  $\geq 30$  and  $\geq 90\%$  of heterozygote genotypes passing an ABRatio binomial test with  $P$ -value  
466  $> 10^{-2}$  per batch of sequencing and genotyping platform (i.e, HiSeq+NSV4, HiSeq+Pipeline 2.0,  
467 NovaSeq+Pipeline 2.0). We then subsetted those to the following groups of MAC/MAF categories:  
468 MAC 1, 2, 3, 4, 5, 6-10, 11-20, MAC 20 - MAF 0.001, MAF 0.001 - 0.005.

## 469 **Relatedness, ancestry and principal components**

### 470 **Kinship**

471 We calculated kinship coefficients among all pairs of samples using software plink2 and its imple-  
472 mentation of the KING robust algorithm. We used a kinship cutoff  $< 0.0442$  to select unrelated  
473 individuals with argument “-king-cutoff”.

### 474 **Genetic Ancestry Prediction**

475 To infer the ancestry of each individual we performed principal components analysis (PCA) on  
476 unrelated 1KGP3 individuals with GCTA v1.93.1\_beta software using HQ common SNPs and  
477 inferred the first 20 PCs. We calculated loadings for each SNP which we used to project aggV2 and  
478 aggCOVID\_v4.2 individuals onto the 1KGP3 PCs. We then trained a random forest algorithm  
479 from R-package randomForest with the first 10 1KGP3 PCs as features and the super-population  
480 ancestry of each individual as labels. These were ‘AFR’ for individuals of African ancestry, ‘AMR’  
481 for individuals of American ancestry, ‘EAS’ for individuals of East Asian ancestry, ‘EUR’ for  
482 individuals of European ancestry, and ‘SAS’ for individuals of South Asian ancestry. We used  
483 500 trees for the training. We then used the trained model to assign probability of belonging to  
484 a certain super-population class for each individual in our cohorts. We assigned individuals to a  
485 super-population when class probability  $\geq 0.8$ . Individuals for which no class had probability  
486  $\geq 0.8$  were labelled as “unassigned” and were not included in the analyses.

### 487 **Principal component analysis**

488 After labelling each individual with predicted genetic ancestry, we calculated ancestry-specific PCs  
489 using GCTA v1.93.1\_beta, *i.e.*. We computed 20 PCs for each of the ancestries that were used in  
490 the association analyses (AFR, EAS, EUR, and SAS).

### 491 **Variant Quality Control**

492 Variant QC was performed to ensure high quality of variants and to minimise batch effects due to  
493 using samples from different sequencing platforms (NovaSeq6000 and HiseqX) and different variant  
494 callers (Strelka2 and DRAGEN). We first masked low-quality genotypes setting them to missing,  
495 merged aggregate files and then performed additional variant quality control separately for the two  
496 major types of association analyses, GWAS and AVT, which concerned common and rare variants,  
497 respectively.

### 498 **Masking**

499 Prior to any analysis we masked low quality genotypes using bcftools setGT module. Genotypes  
500 with  $DP < 10$ ,  $GQ < 20$ , and heterozygote genotypes failing an AB-ratio binomial test with  $P$ -value  $<$   
501  $10^{-3}$  were set to missing.

502 We then converted the masked VCF files to plink and bgen format using plink v.2.0.

### 503 **Merging of aggregate samples**

504 Merging of aggV2 and aggCOVID\_v4.2 samples was done using plink files with masked genotypes  
505 and the merge function of plink v.1.9.<sup>46</sup> for variants that were found in both aggregates.

### 506 **GWAS analyses**

#### 507 **Variant QC**

508 We restricted all GWAS analyses to common variants applying the following filters using plink v1.9:  
509  $MAF > 0$  in both cases and controls,  $MAF > 0.5\%$  and  $MAC > 20$ , missingness  $< 2\%$ , Differential  
510 missingness between cases and controls, mid- $P$ -value  $< 10^{-5}$ , HWE deviations on unrelated controls,  
511 mid- $P$ -value  $< 10^{-6}$ , Multi-allelic variants were additionally required to have  $MAF > 0.1\%$  in both  
512 aggV2 and aggCOVID\_v4.2.

#### 513 **Control-control QC filter**

514 100K aggV2 samples that were aligned and genotype called with the Illumina North Star Version 4  
515 pipeline represented the majority of control samples in our GWAS analyses, whereas all of the cases  
516 were aligned and called with Genomics England pipeline 2.0 (Supplementary Table 1). Therefore,  
517 the alignment and genotyping pipelines partially match the case/control status which necessitates  
518 additional filtering for adjusting for between-pipeline differences in alignment and variant calling. To  
519 control for potential batch effects, we used the overlap of 3,954 samples from the Genomics England  
520 100K participants that were aligned and called with both pipelines. For each variant, we computed  
521 and compared between platforms the inferred allele frequency for the population samples. We then  
522 filtered out all variants that had  $> 1\%$  relative difference in allele frequency between platforms. The  
523 relative difference was computed on a per-population basis for EUR ( $n=3,157$ ), SAS ( $n=373$ ), AFR  
524 ( $n=354$ ) and EAS ( $n=81$ ).

#### 525 **Model**

526 We used a 2-step logistic mixed model regression approach as implemented in SAIGE v0.44.5 for  
527 single variant association analyses. In step 1, SAIGE fits the null mixed model and covariates. In  
528 step 2, single variant association tests are performed with the saddlepoint approximation (SPA)  
529 correction to calibrate unbalanced case-control ratios. We used the HQ common variant sites for  
530 fitting the null model and *sex*, *age*, *age*<sup>2</sup>, *age* \* *sex* and 20 principal components as covariates in  
531 step 1. The principal components were computed separately by predicted genetic ancestry (i.e.,  
532 EUR-specific, AFR-specific, etc.), to capture subtle structure effects.

#### 533 **Analyses**

534 All analyses were done on unrelated individuals with pairwise kinship coefficient  $< 0.0442$ . We  
535 conducted GWAS analyses per genetic ancestry, for all populations for which we had  $>100$  cases  
536 and  $>100$  controls (AFR, EAS, EUR, and SAS).

#### 537 **Multiple testing correction**

538 As our study is testing variants that were directly sequenced by WGS and not imputed, we calculated  
539 the  $P$ -value significance threshold by estimating the effective number of tests. After selecting the

540 final filtered set of tested variants for each population, we LD-pruned in a window of 250Kb and  
541  $r^2 = 0.8$  with plink 1.9. We then computed the Bonferroni-corrected  $P$ -value threshold as 0.05  
542 divided by the number of LD-pruned variants. The  $P$ -value thresholds that were used for declaring  
543 statistical significance are given in Supplementary Table 3.

#### 544 **LD-clumping**

545 We used plink1.9 to do clumping of variants that were genome-wide significant for each analysis with  
546  $P1$  set to per-population  $P$ -value from table X,  $P2 = 0.01$ , clump distance 1500Mb and  $r^2 = 0.1$ .

#### 547 **Conditional analysis**

548 To find the set of independent variants in the per-population analyses, we performed a step-wise  
549 conditional analysis with the GWAS summary statistics for each population using GTCA 1.9.3  
550 `-cojo-slc` function. The parameters for the function were  $pval = 2.2 \times 10^{-8}$ , a distance of 10,000 kb  
551 and a colinear threshold of 0.9<sup>47</sup>.

#### 552 **Fine-mapping**

553 We performed fine-mapping for genome-wide significant signals using Rpackage SusieR v0.11.42<sup>48</sup>.  
554 For each genome-wide significant variant locus, we selected the variants 1.5 Mbp on each side and  
555 computed the correlation matrix among them with plink v1.9. We then run the susieR summary-  
556 statistics based function `susie_rss` and provided the summary z-scores from SAIGE (i.e, effect size  
557 divided by its standard error) and the correlation matrix computed with the same samples that  
558 were used for the corresponding GWAS. We required coverage  $>0.95$  for each identified credible set  
559 and minimum and median correlation coefficients (purity) of  $r=0.1$  and 0.5, respectively.

#### 560 **Functional annotation of credible sets**

561 We annotated all variants included in each credible set identified by SusieR using VEP v99. We also  
562 selected the worst consequence across transcripts using `bcftools +split-vep -s worst`. We also ranked  
563 each variant within each credible set according to the predicted consequence and the ranking was  
564 based on the table provided by Ensembl: [https://www.ensembl.org/info/genome/variation/prediction/predicted\\_data.html](https://www.ensembl.org/info/genome/variation/prediction/predicted_data.html).  
565

#### 566 **Trans-ancestry meta-analysis**

567 We performed a meta-analysis across all ancestries using an inverse-variance weighted method and  
568 control for population stratification for each separate analysis in the METAL software<sup>13</sup>. The  
569 meta-analysed variants were filtered for variants with heterogeneity  $P$ -value  $p < 2.22 \times 10^{-8}$  and  
570 variants that are not present in at least half of the individuals. We used the meta R package to plot  
571 forest plots of the clumped trans-ancestry meta-analysis variants<sup>49</sup>.

#### 572 **LD-based validation of lead GWAS signals**

In order to quantify the support for genome-wide significant signals from nearby variants in LD, we assessed the internal consistency of GWAS results of the lead variants and their surroundings. To this end, we compared observed z-scores at lead variants with the expected z-scores based on those

observed at neighbouring variants. Specifically, we computed the observed z-score for a variant  $i$  as  $s_i = \hat{\beta}/\hat{\sigma}_{\beta}$  and, following the approach of<sup>50</sup>, the imputed z-score at a target variant  $t$  as

$$\hat{s}_t = \mathbf{\Sigma}_{t,P}(\mathbf{\Sigma}_{P,P} + \lambda\mathbf{I})^{-1}\mathbf{s}_P$$

573 where  $\mathbf{s}_P$  are the observed z-scores at a set  $P$  of predictor variants,  $\mathbf{\Sigma}_{x,y}$  is the empirical correlation  
574 matrix of dosage coded genotypes computed on the GWAS sample between the variants in  $x$  and  $y$ ,  
575 and  $\lambda$  is a regularization parameter set to  $10^{-5}$ . The set  $P$  of predictor variants consisted of all  
576 variants within 100 kb of the target variant with a genotype correlation with the target variant  
577 greater than 0.25. This approach is similar to one proposed recently by Chen et al.<sup>51</sup>

## 578 Replication

579 We used the Host Genetic Initiative (HGI) GWAS meta-analysis round 6 hospitalised COVID vs  
580 population (B2 analysis), including all genetic ancestries. In order to remove overlapping signals  
581 we performed a mathematical subtraction of the GenOMICC GWAS of European genetic ancestry.  
582 The HGI data was downloaded from <https://www.covid19hg.org/results/r6/>. The subtraction was  
583 performed using MetaSubtract package (version 1.60) for R (version 4.0.2) after removing variants  
584 with the same genomic position and using the lambda.cohortswith genomic inflation calculated on  
585 the GenOMICC summary statistics. Then, we calculated a trans-ancestry meta-analysis for the three  
586 ancestries with summary statistics in 23andMe: African, Latino and European using variants that  
587 passed the 23andMe ancestry QC, with imputation score  $> 0.6$  and with maf  $> 0.005$ . And finally  
588 we performed a final meta-analysis of 23andMe and HGI B2 without GenOMICC to create the final  
589 replication set. Meta-analysis were performed using METAL<sup>13</sup>, with the inverse-variance weighting  
590 method (STDERR mode) and genomic control ON. We considered that a hit was replicating if the  
591 direction of effect in the GenOMICC-subtracted HGI summary statistics was the same as in our  
592 GWAS, and the  $P$ -value was significant after Bonferroni correction for the number of attempted  
593 replications ( $pval < 0.05/25$ ). If the main hit was not present in the HGI-23andMe meta-analysis or  
594 if the hit was not replicating we looked for replication in variants in high LD with the top variant  
595 ( $r^2 > 0.9$ ), which helped replicate two regions.

## 596 Stratified analysis

597 We also performed sex-specific analysis (male and females separately) as well as analysis stratified  
598 by age (*i.e.*, participants  $<60$  and  $\geq 60$  years old) for each super-population set. To compare effect  
599 of variants within groups for the age and sex stratified analysis we first adjusted the effect and error  
600 of each variant for the standard deviation of the trait in each stratified group and then used the  
601 following t-statistic, as in previous studies<sup>52;53</sup>

$$602 \quad t = \frac{b_1 - b_2}{\sqrt{se_1^2 + se_2^2 - 2 \cdot r \cdot se_1 \cdot se_2}}$$

603 where  $b_1$  is the adjusted effect for group 1,  $b_2$  is the adjusted effect for group 2,  $se_1$  and  $se_2$  are  
604 the adjusted standard errors for group 1 and 2 respectively and  $r$  is the Spearman rank correlation  
605 between groups across all genetic variants.

## 606 HLA Imputation and Association Analysis

607 HLA types were imputed at two field (4-digit) resolution for all samples within aggV2 and ag-  
608 gCOVID\_v4.2 for the following seven loci: HLA-A, HLA-C, HLA-B, HLA-DRB1, HLA-DQA1,

609 HLA-DQB1, and HLA-DPB1 using the HIBAG package in R<sup>19</sup>. At time of writing, HLA types  
610 were also imputed for 82% of samples using HLA\*LA<sup>54</sup>. Inferred HLA alleles between HIBAG and  
611 HLA\*LA were >96% identical at 4-digit resolution. HLA association analysis was run under an  
612 additive model using SAIGE; in an identical fashion to the SNV GWAS. The multi-sample VCF  
613 of aggregated HLA type calls from HIBAG were used as input where any allele call with posterior  
614 probability ( $T$ ) < 0.5 were set to missing.

## 615 **Aggregate variant testing (AVT)**

616 Aggregate variant testing on aggCOVID\_v4.2 was performed using SKAT-O as implemented in  
617 SAIGE-GENE v0.44.5<sup>16</sup> on all protein-coding genes. Variant and sample QC for the preparation  
618 and masking of the aggregate files has been described elsewhere. We further excluded SNPs with  
619 differential missingness between cases and controls (mid-P value <  $10^{-5}$ ) or a site-wide missingness  
620 above 5%. Only bi-allelic SNPs with a MAF < 0.5% were included.

621 We filtered the variants to include in the aggregate variant testing by applying two functional  
622 annotation filters: A putative loss of function (*pLoF*) filter, where only variants that are annotated  
623 by LOFTEE<sup>17</sup> as high confidence loss of function were included, and a more lenient (*missense*)  
624 filter where variants that have a consequence of missense or worse as annotated by VEP, with a  
625 CADD\_PHRED score of  $\geq 10$ , were also included. All variants were annotated using VEP v99.  
626 SAIGE-GENE was run with the same covariates used in the single variant analysis: *sex*, *age*, *age*<sup>2</sup>,  
627 *age \* sex* and 20 (population-specific) principal components generated from common variants (MAF  
628  $\geq 5\%$ ).

629 We ran the tests separately by genetically predicted ancestry, as well as across all four ancestries as  
630 a mega-analysis. We considered a gene-wide significant threshold on the basis of the genes tested  
631 per ancestry, correcting for the two masks (*pLoF* and *missense*, Supplementary Table 4).

## 632 **Post-GWAS analysis**

### 633 **Transcriptome-wide Association Studies (TWAS)**

634 We performed TWAS in the MetaXcan framework and the GTEExv8 eQTL and sQTL MASHR-M  
635 models available for download in (<http://predictdb.org/>). We first calculated, using the European  
636 summary statistics, individual TWAS for whole blood and lung with the S-PrediXcan function<sup>55;56</sup>.  
637 Then we performed a metaTWAS including data from all tissues to increase statistical power using  
638 s-MultiXcan<sup>57</sup>. We applied Bonferroni correction to the results in order to choose significant genes  
639 and introns for each analysis.

### 640 **Colocalisation analysis**

641 Significant genes from TWAS, splicing TWAS, metaTWAS and splicing metaTWAS, as well as genes  
642 where one of the top variants was a significant eQTL or sQTL were selected for a colocalisation  
643 analysis using the coloc R package<sup>58</sup>. We chose the lead SNPs from the European ancestry GWAS  
644 summary statistics and a region of  $\pm 200$  kb around each SNP to do the colocalisation with the  
645 identified genes in the region. GTEExv8 whole blood and lung tissue summary statistics and eqtlGen  
646 (which has blood eQTL summary statistics for > 30,000 individuals) were used for the analysis<sup>18;59</sup>.  
647 We first performed a sensitivity analysis of the posterior probability of colocalisation (PPH4) on the

648 prior probability of colocalisation ( $p_{12}$ ), going from  $p_{12} = 10^{-8}$  to  $p_{12} = 10^{-4}$  with the default  
649 threshold being  $p_{12} = 10^{-5}$ . eQTL signal and GWAS signals were deemed to colocalise if these  
650 two criteria were met: (1) At  $P_{12} = 5 \times 10^{-5}$  the probability of colocalisation  $PPH4 > 0.5$  and  
651 (2) At  $p_{12} = 10^{-5}$  the probability of independent signal (PPH3) was not the main hypothesis  
652 ( $PPH3 < 0.5$ ). These criteria were chosen to allow eQTLs with weaker  $P$ -values due to lack of  
653 power in GTE<sub>8</sub>, to be colocalised with the signal when the main hypothesis using small priors  
654 was that there was not any signal in the eQTL data.

655 As the chromosome 3 associated interval is larger than 200 kb, we performed additional colocalisation  
656 including a region up to 500 kb, but no further colocalisations were found.

## 657 References

- 658 [1] Pairo-Castineira, E. *et al.* Genetic mechanisms of critical illness in Covid-19. *Nature* 1–1 (2020).
- 659 [2] COVID-19 Host Genetics Initiative. Mapping the human genetic architecture of COVID-19.  
660 *Nature* (2021). URL <https://doi.org/10.1038/s41586-021-03767-x>.
- 661 [3] Zhang, Q. *et al.* Inborn errors of type I IFN immunity in patients with life-threatening COVID-  
662 19. *Science (New York, N.y.)* **370**, eabd4570 (2020). URL [https://www.ncbi.nlm.nih.gov/pmc](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7857407/)  
663 [articles/PMC7857407/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7857407/).
- 664 [4] Ellinghaus, D. *et al.* Genomewide association study of severe covid-19 with respiratory failure.  
665 *The New England journal of medicine* **383**, 1522–1534 (2020).
- 666 [5] Docherty, A. B. *et al.* Features of 20 133 UK patients in hospital with covid-19 using the  
667 ISARIC WHO Clinical Characterisation Protocol: Prospective observational cohort study. *BMJ*  
668 **369** (2020).
- 669 [6] Dorward, D. A. *et al.* Tissue-Specific Immunopathology in Fatal COVID-19. *American Journal*  
670 *of Respiratory and Critical Care Medicine* **203**, 192–201 (2021).
- 671 [7] Millar, J. E. *et al.* Robust, reproducible clinical patterns in hospitalised patients with COVID-19.  
672 *medRxiv* 2020.08.14.20168088 (2020).
- 673 [8] Horby, P. *et al.* Dexamethasone in Hospitalized Patients with Covid-19 — Preliminary Report.  
674 *New England Journal of Medicine* (2020).
- 675 [9] Degenhardt, F. *et al.* New susceptibility loci for severe COVID-19 by detailed GWAS analysis  
676 in European populations (2021).
- 677 [10] Kosmicki, J. A. *et al.* Pan-ancestry exome-wide association analyses of COVID-19 outcomes  
678 in 586,157 individuals. *American Journal of Human Genetics* **108**, 1350–1355 (2021). URL  
679 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8173480/>.
- 680 [11] Povysil, G. *et al.* Rare loss-of-function variants in type i ifn immunity genes are not associated  
681 with severe covid-19. *The Journal of clinical investigation* **131** (2021).
- 682 [12] Zhou, W. *et al.* Efficiently controlling for case-control imbalance and sample relatedness in  
683 large-scale genetic association studies. *Nature Genetics* **50**, 1335–1341 (2018). URL [http:](http://www.nature.com/articles/s41588-018-0184-y)  
684 [//www.nature.com/articles/s41588-018-0184-y](http://www.nature.com/articles/s41588-018-0184-y).

- 685 [13] Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of genomewide  
686 association scans. *Bioinformatics (Oxford, England)* **26**, 2190–2191 (2010).
- 687 [14] Rentzsch, P., Witten, D., Cooper, G. M., Shendure, J. & Kircher, M. CADD: predicting  
688 the deleteriousness of variants throughout the human genome. *Nucleic Acids Research* **47**,  
689 D886–D894 (2018). URL <https://doi.org/10.1093/nar/gky1016>.
- 690 [15] Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable  
691 selection in regression, with application to genetic fine mapping. *Journal of the Royal Statistical*  
692 *Society: Series B (Statistical Methodology)* **82**, 1273–1300 (2020). URL [https://rss.onlinelibrar](https://rss.onlinelibrary.wiley.com/doi/full/10.1111/rssb.12388)  
693 [y.wiley.com/doi/full/10.1111/rssb.12388](https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/rssb.12388)[https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/](https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/rssb.12388)  
694 [rssb.12388https://rss.onlinelibrary.wiley.com/doi/10.1111/rssb.12388](https://rss.onlinelibrary.wiley.com/doi/10.1111/rssb.12388).
- 695 [16] Zhou, W. *et al.* Scalable generalized linear mixed model for region-based association tests in  
696 large biobanks and cohorts. *Nature Genetics* **52**, 634–639 (2020). URL [https://www.nature.c](https://www.nature.com/articles/s41588-020-0621-6)  
697 [om/articles/s41588-020-0621-6](https://www.nature.com/articles/s41588-020-0621-6).
- 698 [17] Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in  
699 141,456 humans. *Nature* **581**, 434–443 (2020). URL [https://www.nature.com/articles/s41586-](https://www.nature.com/articles/s41586-020-2308-7)  
700 [020-2308-7](https://www.nature.com/articles/s41586-020-2308-7).
- 701 [18] Consortium, T. G. The GTEx Consortium atlas of genetic regulatory effects across human  
702 tissues. *Science* **369**, 1318–1330 (2020). URL [https://science.sciencemag.org/content/3](https://science.sciencemag.org/content/369/6509/1318)  
703 [69/6509/1318](https://science.sciencemag.org/content/369/6509/1318). Publisher: American Association for the Advancement of Science \_eprint:  
704 <https://science.sciencemag.org/content/369/6509/1318.full.pdf>.
- 705 [19] Zheng, X. *et al.* HIBAG - HLA genotype imputation with attribute bagging. *Pharmacogenomics*  
706 *Journal* **14**, 192–200 (2014).
- 707 [20] Dunning, J. W. *et al.* Open source clinical science for emerging infections. *The Lancet Infectious*  
708 *Diseases* **14**, 8–9 (2014).
- 709 [21] Repurposed Antiviral Drugs for Covid-19 — Interim WHO Solidarity Trial Results. *New*  
710 *England Journal of Medicine* **0**, null (2020).
- 711 [22] Dong, B. *et al.* Phospholipid scramblase 1 potentiates the antiviral activity of interferon.  
712 *Journal of virology* **78**, 8983–93 (2004).
- 713 [23] Luo, W. *et al.* Phospholipid scramblase 1 interacts with influenza a virus np, impairing its  
714 nuclear import and thereby suppressing virus replication. *PLoS pathogens* **14**, e1006851 (2018).
- 715 [24] Chen, M.-H. *et al.* Phospholipid Scramblase 1 Contains a Nonclassical Nuclear Localization  
716 Signal with Unique Binding Site in Importin A\*. *Journal of Biological Chemistry* **280**, 10599–  
717 10606 (2005).
- 718 [25] Chen, C.-W., Sowden, M., Zhao, Q., Wiedmer, T. & Sims, P. J. Nuclear phospholipid scramblase  
719 1 prolongs the mitotic expansion of granulocyte precursors during G-CSF-induced granulopoiesis.  
720 *Journal of Leukocyte Biology* **90**, 221–233 (2011).
- 721 [26] Bevers, E. M. & Williamson, P. L. Phospholipid scramblase: An update. *FEBS Letters* **584**,  
722 2724–2730 (2010).

- 723 [27] Yu, Y. *et al.* Bcl11a is essential for lymphoid development and negatively regulates p53. *The*  
724 *Journal of experimental medicine* **209**, 2467–83 (2012).
- 725 [28] Reizis, B. Plasmacytoid Dendritic Cells: Development, Regulation, and Function. *Immunity*  
726 **50**, 37–50 (2019).
- 727 [29] Zhang, Y., Lu, L., Furlonger, C., Wu, G. E. & Paige, C. J. Hemokinin is a hematopoietic-specific  
728 tachykinin that regulates b lymphopoiesis. *Nature immunology* **1**, 392–7 (2000).
- 729 [30] Wang, W. *et al.* Hemokinin-1 activates the mapk pathway and enhances b cell proliferation  
730 and antibody production. *Journal of immunology (Baltimore, Md. : 1950)* **184**, 3590–7 (2010).
- 731 [31] Janelins, B. M. *et al.* Proinflammatory tachykinins that signal through the neurokinin 1  
732 receptor promote survival of dendritic cells and potent cellular immunity. *Blood* **113**, 3017–26  
733 (2009).
- 734 [32] Thwaites, R. S. *et al.* Inflammatory profiles across the spectrum of disease reveal a distinct role  
735 for GM-CSF in severe COVID-19. *Science Immunology* **6** (2021).
- 736 [33] Lang, F. M., Lee, K. M.-C., Teijaro, J. R., Becher, B. & Hamilton, J. A. Gm-csf-based treatments  
737 in covid-19: reconciling opposing therapeutic approaches. *Nature reviews. Immunology* **20**,  
738 507–514 (2020).
- 739 [34] Moore, C. *et al.* Resequencing Study Confirms That Host Defense and Cell Senescence Gene  
740 Variants Contribute to the Risk of Idiopathic Pulmonary Fibrosis. *American Journal of*  
741 *Respiratory and Critical Care Medicine* **200**, 199–208 (2019). URL [https://www.atsjournals.or](https://www.atsjournals.org/doi/10.1164/rccm.201810-1891OC)  
742 [g/doi/10.1164/rccm.201810-1891OC](https://www.atsjournals.org/doi/10.1164/rccm.201810-1891OC). Publisher: American Thoracic Society - AJRCCM.
- 743 [35] Takatsu, H. *et al.* Phospholipid flippase activities and substrate specificities of human type iv  
744 p-type atpases localized to the plasma membrane. *The Journal of biological chemistry* **289**,  
745 33543–56 (2014).
- 746 [36] Bevers, E. M., Comfurius, P. & Zwaal, R. F. Changes in membrane phospholipid distribution  
747 during platelet activation. *Biochimica et biophysica acta* **736**, 57–66 (1983).
- 748 [37] Zwaal, R. F., Comfurius, P. & van Deenen, L. L. Membrane asymmetry and blood coagulation.  
749 *Nature* **268**, 358–60 (1977).
- 750 [38] Shrine, N. *et al.* New genetic signals for lung function highlight pathways and chronic obstructive  
751 pulmonary disease associations across multiple ancestries. *Nature genetics* **51**, 481–493 (2019).
- 752 [39] Mankelov, T. J. *et al.* Blood group type A secretors are associated with a higher risk of  
753 COVID-19 cardiovascular disease complications. *eJHaem* **2**, 175–187 (2021).
- 754 [40] Kelly, R. J., Rouquier, S., Giorgi, D., Lennon, G. G. & Lowe, J. B. Sequence and expression  
755 of a candidate for the human secretor blood group alpha(1,2)fucosyltransferase gene (fut2).  
756 homozygosity for an enzyme-inactivating nonsense mutation commonly correlates with the  
757 non-secretor phenotype. *The Journal of biological chemistry* **270**, 4640–9 (1995).
- 758 [41] Ferrer-Admetlla, A. *et al.* A natural history of fut2 polymorphism in humans. *Molecular biology*  
759 *and evolution* **26**, 1993–2003 (2009).

- 760 [42] Imbert-Marcille, B.-M. *et al.* A fut2 gene common polymorphism determines resistance to  
761 rotavirus a of the p[8] genotype. *The Journal of infectious diseases* **209**, 1227–30 (2014).
- 762 [43] Tian, C. *et al.* Genome-wide association and hla region fine-mapping studies identify suscepti-  
763 bility loci for multiple common infections. *Nature communications* **8**, 599 (2017).
- 764 [44] Kachuri, L. *et al.* The landscape of host genetic factors involved in immune response to common  
765 viral infections. *medRxiv : the preprint server for health sciences* (2020).
- 766 [45] Blackwell, C. C. *et al.* Non-secretion of abo antigens predisposing to infection by neisseria  
767 meningitidis and streptococcus pneumoniae. *Lancet (London, England)* **2**, 284–5 (1986).
- 768 [46] Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and Population-Based  
769 Linkage Analyses. *The American Journal of Human Genetics* **81**, 559–575 (2007). URL  
770 <https://www.sciencedirect.com/science/article/pii/S0002929707613524>.
- 771 [47] Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics  
772 identifies additional variants influencing complex traits. *Nature Genetics* **44**, 369–375 (2012).  
773 URL <https://doi.org/10.1038/ng.2213>.
- 774 [48] Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable  
775 selection in regression, with application to genetic fine mapping. *Journal of the Royal Statistical*  
776 *Society Series B (Statistical Methodology)* **82**, 1273–1300 (2020). URL [https://rss.onlinelibrary.](https://rss.onlinelibrary.wiley.com/doi/10.1111/rssb.12388)  
777 [wiley.com/doi/10.1111/rssb.12388](https://rss.onlinelibrary.wiley.com/doi/10.1111/rssb.12388).
- 778 [49] Balduzzi, S., Rücker, G. & Schwarzer, G. How to perform a meta-analysis with R: a practical  
779 tutorial. *Evidence-Based Mental Health* **22**, 153–160 (2019). URL [https://ebmh.bmj.com/con](https://ebmh.bmj.com/content/22/4/153)  
780 [tent/22/4/153](https://ebmh.bmj.com/content/22/4/153). Publisher: Royal College of Psychiatrists Section: Statistics in practice.
- 781 [50] Pasaniuc, B. *et al.* Fast and accurate imputation of summary statistics enhances evidence  
782 of functional enrichment. *Bioinformatics* **30**, 2906–2914 (2014). URL [https://doi.org/](https://doi.org/10.1093/bioinformatics/btu416)  
783 [10.1093/bioinformatics/btu416](https://doi.org/10.1093/bioinformatics/btu416). [https://academic.oup.com/bioinformatics/article-](https://academic.oup.com/bioinformatics/article-pdf/30/20/2906/17147061/btu416.pdf)  
784 [pdf/30/20/2906/17147061/btu416.pdf](https://academic.oup.com/bioinformatics/article-pdf/30/20/2906/17147061/btu416.pdf).
- 785 [51] Chen, W. *et al.* Improved analyses of GWAS summary statistics by reducing data heterogeneity  
786 and errors (2020).
- 787 [52] Winkler, T. W. *et al.* The influence of age and sex on genetic associations with adult body size  
788 and shape: A large-scale genome-wide interaction study. *PLOS Genetics* **11**, 1–42 (2015). URL  
789 <https://doi.org/10.1371/journal.pgen.1005378>.
- 790 [53] Bernabeu, E. *et al.* Sexual differences in genetic architecture in uk biobank. *bioRxiv* (2020).  
791 URL <https://www.biorxiv.org/content/early/2020/07/21/2020.07.20.211813>. [https:](https://www.biorxiv.org/content/early/2020/07/21/2020.07.20.211813.full.pdf)  
792 [//www.biorxiv.org/content/early/2020/07/21/2020.07.20.211813.full.pdf](https://www.biorxiv.org/content/early/2020/07/21/2020.07.20.211813.full.pdf).
- 793 [54] Dilthey, A. T. *et al.* HLA\*LA—HLA typing from linearly projected graph alignments. *Bioin-*  
794 *formatics* **35**, 4394–4396 (2019). URL <https://doi.org/10.1093/bioinformatics/btz235>.  
795 <https://academic.oup.com/bioinformatics/article-pdf/35/21/4394/30330845/btz235.pdf>.
- 796 [55] Gamazon, E. R. *et al.* A gene-based association method for mapping traits using reference  
797 transcriptome data. *Nature Genetics* **47**, 1091–1098 (2015). URL <https://doi.org/10.1038/ng.3>  
798 [367](https://doi.org/10.1038/ng.3).

- 799 [56] Barbeira, A. N. *et al.* Exploring the phenotypic consequences of tissue specific gene expression  
800 variation inferred from GWAS summary statistics. *Nature Communications* **9**, 1825 (2018).  
801 URL <https://doi.org/10.1038/s41467-018-03621-1>.
- 802 [57] Barbeira, A. N. *et al.* Integrating predicted transcriptome from multiple tissues improves  
803 association detection. *PLOS Genetics* **15**, 1–20 (2019). URL <https://doi.org/10.1371/journal.pgen.1007889>. Publisher: Public Library of Science.
- 805 [58] Giambartolomei, C. *et al.* Bayesian Test for Colocalisation between Pairs of Genetic Association  
806 Studies Using Summary Statistics. *PLOS Genetics* **10**, e1004383 (2014). URL [https://journals](https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1004383)  
807 [.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1004383](https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1004383). Publisher: Public Library of  
808 Science.
- 809 [59] Vösa, U. *et al.* Unraveling the polygenic architecture of complex traits using blood eQTL  
810 metaanalysis. *bioRxiv* 447367 (2018). URL [http://biorxiv.org/content/early/2018/10/19/447](http://biorxiv.org/content/early/2018/10/19/447367.abstract)  
811 [367.abstract](http://biorxiv.org/content/early/2018/10/19/447367.abstract).