

1 Common and rare variant analyses combined with single-cell 2 multiomics reveal cell-type-specific molecular mechanisms 3 of COVID-19 severity

4 Sai Zhang^{1,2,18}, Johnathan Cooper-Knock^{3,18}, Annika K. Weimer^{1,2}, Calum Harvey³,
5 Thomas H. Julian³, Cheng Wang^{4,5,6,7}, Jingjing Li^{4,5,6,7}, Simone Furini⁸, Elisa Frullanti^{8,9},
6 Francesca Fava^{8,9,10}, Alessandra Renieri^{8,9,10}, Cuiping Pan¹¹, Jina Song^{2,11}, Paul
7 Billing-Ross^{2,11}, Peng Gao^{1,2}, Xiaotao Shen^{1,2}, Ilia Sarah Timpanaro¹², Kevin P. Kenna¹²,
8 VA Million Veteran Program, GEN-COVID Network¹³, Mark M. Davis^{14,15,16}, Philip S.
9 Tsao^{11,17,*}, and Michael P. Snyder^{1,2,*}

10 ¹Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA

11 ²Center for Genomics and Personalized Medicine, Stanford University School of
12 Medicine, Stanford, CA, USA

13 ³Sheffield Institute for Translational Neuroscience, University of Sheffield, Sheffield, UK

14 ⁴Department of Neurology, School of Medicine, University of California, San Francisco,
15 CA, USA

16 ⁵Eli and Edythe Broad Center of Regeneration Medicine and Stem Cell Research,
17 University of California, San Francisco, CA, USA

18 ⁶Bakar Computational Health Sciences Institute, University of California, San Francisco,
19 CA, USA

20 ⁷Parker Institute for Cancer Immunotherapy, University of California, San Francisco, CA,
21 USA

22 ⁸Med Biotech Hub and Competence Center, Department of Medical Biotechnologies,
23 University of Siena, Siena, Italy

24 ⁹Medical Genetics, Department of Medical Biotechnologies, University of Siena, Siena,
25 Italy

26 ¹⁰Genetica Medica, Azienda Ospedaliero-Universitaria Senese, Siena, Italy

27 ¹¹VA Palo Alto Epidemiology Research and Information Center for Genomics, VA
28 Palo Alto Health Care System, Palo Alto, CA, USA

29 ¹²Department of Neurology, Brain Center Rudolf Magnus, University Medical Center
30 Utrecht, Utrecht, The Netherlands

31 ¹³A list of members and affiliations appears in the Supplementary file

32 ¹⁴Institute for Immunity, Transplantation and Infection, Stanford University School of
33 Medicine, Stanford, CA, USA

34 ¹⁵Department of Microbiology and Immunology, Stanford University School of Medicine,
35 Stanford, CA, USA

36 ¹⁶Howard Hughes Medical Institute, Stanford University School of Medicine, Stanford,
37 CA, USA

38 ¹⁷Department of Medicine, Stanford University School of Medicine, Stanford, CA,
39 USA

40 ¹⁸These authors contributed equally to this work.

41 *Correspondence: ptsao@stanford.edu (P.S.T.), mpsnyder@stanford.edu (M.P.S.)

42 **ABSTRACT**

43 The determinants of severe COVID-19 in non-elderly adults are poorly understood,
44 which limits opportunities for early intervention and treatment. Here we present novel
45 machine learning frameworks for identifying common and rare disease-associated
46 genetic variation, which outperform conventional approaches. By integrating single-cell
47 multiomics profiling of human lungs to link genetic signals to cell-type-specific functions,
48 we have discovered and validated over 1,000 risk genes underlying severe COVID-19
49 across 19 cell types. Identified risk genes are overexpressed in healthy lungs but
50 relatively downregulated in severely diseased lungs. Genetic risk for severe COVID-19,
51 within both common and rare variants, is particularly enriched in natural killer (NK) cells,
52 which places these immune cells upstream in the pathogenesis of severe disease.
53 Mendelian randomization indicates that failed NKG2D-mediated activation of NK cells
54 leads to critical illness. Network analysis further links multiple pathways associated with
55 NK cell activation, including type-I-interferon-mediated signalling, to severe COVID-19.
56 Our rare variant model, PULSE, enables sensitive prediction of severe disease in
57 non-elderly patients based on whole-exome sequencing; individualized predictions are
58 accurate independent of age and sex, and are consistent across multiple populations
59 and cohorts. Risk stratification based on exome sequencing has the potential to
60 facilitate post-exposure prophylaxis in at-risk individuals, potentially based around
61 augmentation of NK cell function. Overall, our study characterizes a comprehensive
62 genetic landscape of COVID-19 severity and provides novel insights into the molecular
63 mechanisms of severe disease, leading to new therapeutic targets and sensitive
64 detection of at-risk individuals.

65

66

67

68

69

70

71 INTRODUCTION

72 Infection with severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) giving
73 rise to coronavirus disease 2019 (COVID-19) has caused a global pandemic with
74 almost unprecedented morbidity and mortality¹. The severity of COVID-19 is markedly
75 variable ranging from an asymptomatic infection to fatal multiorgan failure. Severity
76 correlates with age and comorbidities² but not exclusively³. Indeed, host genetics has
77 been thought to be an essential determinant of severity⁴, but this is poorly understood.
78 Improved tools to identify individuals at risk of severe COVID-19 could facilitate
79 life-saving precision medicine.

80 There have been several efforts to address the genetic basis of COVID-19 severity^{5,6},
81 including large-scale genome-wide association studies (GWASs)^{7,8} and rare variant
82 approaches⁹⁻¹². However, the biological interpretation of those identified loci has been
83 difficult, partially because of the confounding effects of patient age and comorbidities¹³.
84 The development of novel therapies is likely to result from understanding and modifying
85 the host immune response to the SARS-CoV-2 virus, independent of immutable factors
86 such as age, sex, and general health.

87 A primary cause of morbidity and mortality in COVID-19 is respiratory disease and
88 specifically, a hyperinflammatory response within the lung that occurs in an
89 age-independent manner¹⁴. This is the basis of a number of interventions based on
90 immunosuppression¹⁵, which have repurposed treatments used for other diseases,
91 particularly autoimmune diseases. Efficacy and the side effect profile is likely to be
92 improved by a COVID-19-specific immunomodulatory approach.

93 Profiles of the immune response associated with severe COVID-19 have produced a
94 number of conflicting observations. These studies have variously linked COVID-19
95 severity to CD8 T cells¹⁶, CD19 B cells¹⁷, eosinophils¹⁸, and myeloid cells¹⁹. Single-cell
96 omic profiling has demonstrated the differential function of various immune cell types in
97 severe disease as opposed to mild disease or non-infected condition²⁰⁻²⁵. However,
98 these studies have focused on transcriptomics rather than the underlying genomics, and
99 have been observational rather than predictive.

100 Failure of the type I interferon response is linked to the incidence of severe COVID-19.
101 SARS-CoV-2 can initially inhibit the normal type I interferon response²⁶ in order to
102 facilitate viral replication. This delay is thought to be an essential determinant of a later
103 hyperinflammatory response and consequently of COVID-19 severity²⁷. Natural killer
104 (NK) cells form a crucial component of the innate immune response to viral infections.
105 Interestingly, NK cells are activated via the type I interferon response. Genetic evidence
106 suggests that NK cell function is a key determinant of severe COVID-19, including
107 loss-of-function (LoF) variants within an essential NK cell activating receptor, NKG2C, in
108 patients suffering severe COVID-19²⁸. A recent study of autoantibodies supports this
109 conclusion by showing that the impaired activation of NK cells, via the type I interferon
110 response in particular, is associated with severe COVID-19²⁹. All of this evidence is
111 suggestive of a role for NK cells in severe COVID-19 but not conclusive.

112 To understand the genetic basis of COVID-19 severity as well as gain insights into its
113 molecular mechanisms, we sought out to integrate the genetic architecture of severe
114 COVID-19, profiled in an age-independent manner, with single-cell-resolution functional
115 profiling of lung tissue. We developed two machine learning frameworks, RefMap and
116 PULSE, for common and rare variant analysis respectively, with increased discovery
117 power compared to traditional methodology. Using our approaches, we identified over
118 1,000 genes associated with critical illness across 19 cell types, and the
119 cell-type-specific molecular mechanisms underlying severe disease were uncovered.
120 Notably, both common and rare variant analyses underscored the importance of NK
121 cells in determining COVID-19 severity, which extends previous literature³⁰. We have
122 developed a prediction model for severe COVID-19 using rare variants profiled by
123 exome sequencing, which achieves sensitive and age- and population-independent risk
124 prediction across multiple cohorts. This prediction method could be particularly useful
125 for targeting medical interventions for individuals where SARS-CoV-2 vaccination is not
126 possible or is not effective³¹. Altogether, our study unveils a holistic genetic landscape of
127 COVID-19 severity and provides a better understanding of the disease pathogenesis,
128 implicating new prevention strategies and therapeutic targets.

129 **RESULTS**

130 **RefMap analysis of common variants uncovers cell-type-specific genetic basis of**
131 **COVID-19 severity**

132 We used the RefMap machine learning model (**Methods**) to identify the genomic
133 regions and genes associated with severe COVID-19. Briefly, RefMap is a Bayesian
134 network that combines genetic signals (e.g., allele Z-scores) with functional genomic
135 profiling (e.g., ATAC-seq and ChIP-seq) to fine-map risk regions for complex diseases.
136 With RefMap, we can scan the genome for functional regions in which
137 disease-associated genetic variation is significantly shifted from the null distribution. The
138 power of the RefMap model for gene discovery and recovery of missing heritability has
139 been demonstrated in our recent work³². Here, to achieve cell-type-specific resolution
140 within multicellular tissue, we modified RefMap to integrate single-cell multiomic profiling
141 of human lungs with COVID-19 GWAS data (**Fig. 1a**). In particular, we obtained
142 summary statistics (COVID-19 Host Genetics Initiative, Release 5, phenotype definition
143 A2; 5,101 cases versus 1,383,241 population controls) from the largest GWAS study of
144 COVID-19⁷, where age, sex, and 20 first principal components were included in the
145 analysis as covariates. Severe COVID-19 was defined by the requirement for
146 respiratory support or death attributed to COVID-19. Human lung single-cell multiomic
147 profiling, including snRNA-seq and snATAC-seq, was retrieved from a recent study of
148 healthy individuals³³. There are 19 cell types identified in both snATAC-seq and
149 snRNA-seq profiles, including epithelial (alveolar type 1 (AT1), alveolar type 2 (AT2),
150 club, ciliated, basal, and pulmonary neuroendocrine (PNEC)), mesenchymal
151 (myofibroblast, pericyte, matrix fibroblast 1 (matrix fib. 1), and matrix fibroblast 2 (matrix
152 fib. 2)), endothelial (arterial, lymphatic, capillary 1 (cap1), and capillary 2 (cap2)), and
153 hematopoietic (macrophage, B-cell, T-cell, NK cell, and enucleated erythrocyte) cell
154 types. We adopted these 19 cell types as the reference set within lung tissue throughout
155 our study. Based on snATAC-seq peaks called in one or more of the 19 cell types to
156 annotate functional regions, we used RefMap to identify disease-associated genomic
157 regions from the COVID-19 GWAS data, which resulted in 6,662 1kb regions passing
158 the 5% significance threshold ($Q^{+/-}$ -score>0.95, **Methods**; referred to as RefMap
159 COVID-19 regions). These identified regions were further intersected with open
160 chromatin in individual cell types based on corresponding snATAC-seq peaks, resulting

161 in cell-type-specific RefMap regions (mean per cell type =1732.8, standard deviation
162 (SD)=623.5; **Fig. 1b, Supplementary Table 1**). After removing RefMap regions present
163 in more than one cell type (mean per cell type =121.2, SD=142.7), we observed only a
164 weak correlation between the number of unique RefMap regions and the number of
165 snATAC-seq peaks detected per cell type (Spearman $\rho=0.40$, $P>0.05$; **Fig. 1b**),
166 indicating enrichment of genetic signals within certain cell types.

167 Next, we sought to map the target genes of RefMap COVID-19 regions in a
168 cell-type-specific manner. In particular, we identified the closest genes that are
169 expressed in the corresponding cell type for individual RefMap regions (**Methods**). In
170 total, we discovered 1,370 genes (referred to RefMap COVID-19 genes; mean per cell
171 type =279.9 and SD=80.3; **Fig. 1c, Supplementary Table 1**) associated with the severe
172 disease. Interestingly, hematopoietic cells have the largest number of unique RefMap
173 regions and genes among all major cell types (**Fig. 1d**); for example, there is a
174 significant enrichment of unique RefMap regions observed for hematopoietic cells
175 versus epithelial cells ($P=5.2e-03$, odds ratio (OR)=1.15, Fisher's exact test; **Fig. 1d**).
176 This indicates a critical role of immune cells, which are primarily hematopoietic, in the
177 development of severe COVID-19¹⁶⁻¹⁹. To profile the cell-cell interactions underlying
178 severe COVID-19 from a genetic perspective, we constructed a cell correlation matrix
179 based on the overlap of RefMap genes between cell types (**Fig. 1e**). We discovered that
180 the correlation is strongest between functionally related cells, demonstrating that the
181 RefMap signal is consistent with known biology³³.

182 To replicate our findings, we obtained SNPs associated with severe COVID-19 from a
183 GWAS for an entirely independent sample set (the 23andMe cohort, 15,434
184 COVID-19-positive cases and 1,035,598 population controls)⁵. The total union of
185 RefMap regions is significantly enriched with SNPs associated with multiple COVID-19
186 phenotypes defined in this new sample set (mean $P<5e-03$, Fisher's exact test; **Fig. 1f**,
187 **Supplementary Table 2; Methods**). Specifically, the most significant enrichment is with
188 SNPs associated with COVID-19 requiring respiratory support (mean $P=4.68e-04$,
189 Fisher's exact test; **Fig. 1f**). We further performed the enrichment analysis per cell type;
190 only RefMap regions associated specifically with T cells and NK cells are significantly

191 enriched with disease-associated SNPs across all measured COVID-19 phenotypes
192 (mean $P < 0.05$, Fisher's exact test; **Supplementary Table 2**).

193 **Heritability analysis and Mendelian randomization link NK cell function to** 194 **COVID-19 severity**

195 The LD score regression (LDSC)³⁴ has been used to measure the total SNP-based
196 heritability (h^2) from the GWAS study of severe COVID-19 (COVID-19 Host Genetics
197 Initiative, Release 5, phenotype definition A2)⁷. Here, we examined the partitioning of
198 SNP-based heritability for severe COVID-19 within RefMap genes (**Methods**). We
199 discovered that the heritability of severe COVID-19 is significantly enriched for RefMap
200 genes (OR=4.6, standard error (SE)=0.78, $P=1.55e-07$; **Fig. 2a, Supplementary Table**
201 **3**). We compared the proportion of SNP-based heritability captured by RefMap to other
202 methods (**Methods**), including naïve GWAS⁷ and MAGMA³⁵. The proportion of
203 heritability within naïve GWAS genes is 0.15 compared to 0.37 within MAGMA genes,
204 but 0.77 within RefMap genes (**Fig. 2b, Supplementary Table 3**), representing a
205 five-fold improvement in the recovered heritability based on RefMap over traditional
206 methods. The proportion of SNP-based heritability for hospitalized COVID-19
207 (COVID-19 Host Genetics Initiative, Release 5, phenotype definition B2) within RefMap
208 genes is 0.62 and within COVID-19 independent of severity (COVID-19 Host Genetics
209 Initiative, Release 5, phenotype definition C2) it is 0.52 (**Fig. 2b, Supplementary Table**
210 **3**). In both cases the improvement in captured heritability based on RefMap compared
211 to traditional methods is three-fold. Consistent with the design of our model, the
212 recovered heritability is highest in severe COVID-19.

213 Next, we used cell-type-specific RefMap genes to determine which cell types are
214 involved in the development of severe COVID-19. Specifically, we calculated the
215 partitioned heritability per cell type within the severe COVID-19 GWAS (A2) and also
216 within GWAS for hospitalized versus non-hospitalized COVID-19 (B2) and COVID-19
217 versus population (C2) (**Methods**). For severe COVID-19, of all 19 cell types tested, NK
218 cells are the most enriched with SNP-based heritability (OR=8.87, SE=3.68, $P=0.016$;
219 **Fig. 2a, Supplementary Table 3**). The same is also true for hospitalized COVID-19
220 (OR=10.57, SE=4.95, $P=0.039$), but not for COVID-19 irrespective of severity

221 (OR=5.74, SE=3.09, $P=0.077$; **Supplementary Table 3**). Thus, we conclude that NK
222 cell function is enriched with severe disease-associated genetic variation.

223 Two-sample Mendelian randomization (MR) facilitates identification of a causal
224 relationship between an exposure and an outcome³⁶. We examined whether NK cell
225 populations measured in the blood are causally related to severe COVID-19. In total, 46
226 GWAS measures of NK cell subtypes were identified³⁷ (**Methods**). After harmonizing
227 exposure and outcome genetic instruments, we excluded tests with less than five SNPs
228 (**Methods**). With MR, three exposures were shown to be causally related to severe
229 COVID-19 after correcting for multiple testing ($P<1e-03$, multiplicative random effects
230 (MRE), inverse-variance weighted (IVW)). All three exposures relate to NKG2D/CD314
231 expression on the cell surface, where a higher number of NKG2D/CD314- cells was
232 linked to severe COVID-19 (A2) (**Figs. 2c-d**) and a higher number of NKG2D/CD314+
233 cells is protective (**Fig. 2e**). Evidence of genetic pleiotropy (MR PRESSO intercept not
234 significantly different from zero, $P>0.05$; **Fig. 2f**) or instrument heterogeneity ($P>0.05$,
235 Cochran's Q test, and $I^2_{GX}>0.95$; **Fig. 2f**) are not evident. Moreover, robust measures
236 are significant for all three exposures (**Fig. 2f**). We also tested the identical phenotypes
237 with alternative COVID-19 phenotype GWAS; we discovered that CD335+ CD314- cell
238 counts are also causally associated with hospitalized COVID-19 (B2), and with
239 COVID-19 independent of severity (C2) (**Supplementary Fig. 1**), but in each case the
240 effect size is reduced compared to severe COVID-19 (A2). NKG2D/CD314 is a primary
241 receptor responsible for NK cell activation³⁸ and in light of this, we conclude that severe
242 COVID-19 is associated with a loss of NK cell cytotoxicity rather than a gain of function
243 linked to NKG2D/CD314- cells.

244 Inspired by our MR analysis, we further tested if the expression of RefMap genes
245 reflects a functional difference between NKG2D/CD314+ and NKG2D/CD314- cells. We
246 examined the expression levels of NK-cell RefMap genes based on scRNA-seq data
247 from healthy lungs³⁹, and discovered that RefMap gene expression is higher within
248 NKG2D/CD314+ cells than NKG2D/CD314- cells ($P=0.036$, one-tailed Wilcoxon
249 rank-sum test; **Fig. 2g**). This further supports the functional significance of RefMap

250 genes in COVID-19 severity and associates the genetic risk of severe COVID-19
251 directly with NK cell activity.

252 **Transcriptome analysis supports the functional significance of RefMap genes in** 253 **health and severe COVID-19**

254 To link RefMap COVID-19 genes to the underlying biology, we first performed functional
255 enrichment analyses based on gene ontology (GO) and KEGG pathways⁴⁰ (**Figs. 3a**
256 **and 3b, Supplementary Tables 4 and 5**). We observed that RefMap NK-cell genes are
257 enriched with pathways and ontologies related to intra- and intercellular signalling
258 important for NK cell activation, including "Phospholipase D signalling pathway"⁴¹,
259 "Antigen processing and presentation", "regulation of small GTPase mediated signal
260 transduction (GO:0051056)"⁴², and "regulation of intracellular signal transduction
261 (GO:1902531)" (adjusted $P < 0.1$; **Figs. 3a and 3b, Supplementary Tables 4 and 5**).
262 This is consistent with the hypothesis that COVID-19 severity is determined by failed
263 activation of NK cells. Furthermore, the pathway with the highest enrichment is "human
264 immunodeficiency virus (HIV) 1 infection" (adjusted $P = 3e-04$; **Fig. 3b**). Since HIV-1
265 works to suppress NK cell activation⁴³, and NK cell function has been associated with
266 an effective immune response to HIV⁴⁴, this result is also consistent with a role of NK
267 cells in severe COVID-19. Other cell-type-specific RefMap gene lists are also enriched
268 with relevant biological pathways. For example, AT2-cell genes are linked to pathways
269 associated with viral infection such as 'human papillomavirus infection' and 'viral
270 carcinogenesis' (adjusted $P < 0.1$; **Supplementary Table 5**), which is consistent with the
271 established role of AT2 cells as the initial site of SARS-CoV-2 entry into host cells⁴⁵.
272 T-cell genes are enriched with 'IL17-signalling pathway' (adjusted $P = 0.021$;
273 **Supplementary Table 5**) which is interesting in light of previous literature highlighting
274 the production of IL-17 by T cells from COVID-19 patients as a potential therapeutic
275 target⁴⁶.

276 Next, we investigated the baseline expression pattern of RefMap genes in healthy
277 lungs. In particular, we calculated mean expression levels of genes in different cell types
278 based on the lung snRNA-seq data from Wang et al.³³, and then compared the
279 expression of RefMap genes with the total set of expressed genes in each cell type.

280 Interestingly, although the gene expression level was not an input to the RefMap model,
281 RefMap genes are expressed at a higher level compared to the background
282 transcriptome in all 19 cell types, including immune and epithelial cells (false discovery
283 rate (FDR)<0.1, one-tailed Wilcoxon rank-sum test; **Fig. 3c**) with the exception of
284 pericytes (FDR=0.11, Z-score=1.25); it is interesting to note that pericytes may be
285 downstream in the pathogenesis of COVID-19 because they are protected by an
286 endothelial barrier⁴⁷. This supports the functional significance of RefMap genes across
287 multiple cell types in healthy human lungs. As a negative control, we performed a
288 similar expression comparison between non-developmental genes and all expressed
289 genes in lungs, which yielded no significant difference (**Supplementary Fig. 2**;
290 **Methods**). In summary, our transcriptome analyses indicate that RefMap genes are
291 expressed above background in relevant cell types, supporting their important role in
292 lung function.

293 To obtain further insights into the function of RefMap COVID-19 regions, we tested
294 whether RefMap regions are enriched with cell-type-specific candidate cis-regulatory
295 elements (cCREs, or enhancers and promoters) defined by H3K27ac and H3K4me3.
296 We obtained cCREs for lung tissues and primary cells from the ENCODE project⁴⁸
297 (**Supplementary Table 6**), and examined the overlap between those cCREs and
298 RefMap regions by permutation test⁴⁹. We discovered that RefMap regions specific to
299 immune cells (e.g., T cells, B cells, and NK cells) are significantly enriched with cCREs
300 in corresponding cell types (FDR<0.1; **Fig. 3d**). For other cell types, RefMap regions are
301 generally enriched with cCREs from lung tissue (FDR<0.1; **Fig. 3d**). These observations
302 are consistent with an important role of RefMap regions in the regulation of gene
303 expression. Moreover, the enrichment with cCREs across a variety of individuals and
304 datasets supports the generalizability of the genetic architecture defined by RefMap. A
305 similar association between genome-wide snATAC-seq peaks and cCREs was also
306 observed (**Supplementary Fig. 3**).

307 We have shown that RefMap COVID-19 regions are enriched within promoters and
308 enhancers responsible for regulating gene expression in healthy lung tissue. On this
309 basis, we hypothesized that genetic variation within RefMap regions would alter the

310 expression of corresponding target genes in the context of severe disease. Specifically,
311 we proposed that RefMap genes would be expressed at a lower level in lung tissue from
312 severe COVID-19 patients than moderately affected patients. To validate this
313 hypothesis, we obtained scRNA-seq data from the respiratory system for a large
314 COVID-19 cohort²³, including 12 bronchoalveolar lavage fluid (BALF) samples, 22
315 sputum samples, and 1 sample of pleural fluid mononuclear cells (PFMCs) from 27
316 severely and 8 mildly affected patients. Severity was classified based on the World
317 Health Organization (WHO) guidelines
318 (<https://www.who.int/publications/i/item/WHO-2019-nCoV-clinical-2021-1>). For individual
319 cell types, we compared the expression level of RefMap genes in severe patients
320 versus moderately affected patients (**Methods**). Compared to the background
321 transcriptome, we observed that RefMap genes are relatively lower expressed in severe
322 patients in corresponding cell types than in moderate patients (FDR<0.01, one-tailed
323 Wilcoxon rank-sum test; **Fig. 3e**), supporting the functional significance of RefMap
324 genes in severe COVID-19. As a replication experiment, we carried out a similar
325 analysis based on an independent COVID-19 scRNA-seq dataset²², including 9 BALF
326 samples from 6 severe patients and 3 moderate patients (**Methods**). The lower
327 expression of RefMap genes in severe patients is consistent across multiple cell types
328 (FDR<0.01, one-tailed Wilcoxon rank-sum test; **Fig. 3f**). Altogether, these
329 transcriptome-based orthogonal analyses are consistent with the hypothesis that
330 identified cell-type-specific RefMap genes are functionally linked to COVID-19 severity.

331 **PULSE analysis of rare variants enables population-independent prediction of** 332 **COVID-19 severity**

333 RefMap utilizes common genetic variation profiled in GWAS. Biological dysfunction can
334 also be determined by rare variants and therefore we assessed whether there is a
335 significant burden of severe-COVID-19-associated rare variants within RefMap genes.
336 First, we employed a standard methodology using SKAT⁵⁰ rare-variant burden testing
337 applied to whole-exome sequencing (WES) data from the GEN-COVID cohort⁵¹,
338 including non-elderly patients who suffered severe COVID-19 requiring respiratory
339 support, and individuals who suffered non-severe COVID-19 not requiring

340 hospitalization (**Methods**). No individual gene is enriched with significant genetic burden
341 after adjusting for multiple testing (**Supplementary Fig. 4**). This is true whether we
342 tested genome-wide or only for RefMap COVID-19 genes. However, for one subset of
343 cell-type-specific RefMap genes, the median P -value was lower than expected: NK cells
344 ($P < 0.05$, permutation test; **Fig. 4a**; **Methods**). This result from the analysis of rare
345 genetic variation in an independent cohort is convergent with our common variant
346 analysis, highlighting NK cell biology as a critical determinant of COVID-19 severity.

347 Traditional rare-variant burden testing failed to identify any enrichment of
348 COVID-19-associated variants within a single gene although there is significant
349 enrichment in the group of 236 NK-cell RefMap genes. This suggests that traditional
350 burden testing is underpowered when applied to the GEN-COVID dataset. We decided
351 to develop a new method with increased sensitivity. Here we present PULSE
352 (probabilistic burden analysis based on functional estimation), a discriminative Bayesian
353 network that integrates functional annotations of rare variants to model the relationship
354 between genotype and phenotype (**Fig. 4b**; **Methods**, **Supplementary Notes**). In
355 particular, PULSE combines multiple predictions of functional effects for different types
356 of variants, including missense, nonsense, splicing-site, and small insertion-deletion
357 (indel) mutations (**Supplementary Table 7**). After aggregating those functional scores
358 for individual genes, PULSE learns the importance of different annotations and genes
359 from the training data and maps the phenotype from the genotype in a bilinear form
360 (**Fig. 4b**; **Methods**). With PULSE as a discovery-by-prediction strategy, we are able to
361 (i) predict individual phenotypes from personal genotypes and (ii) discover
362 phenotype-associated genes by model interpretation.

363 We applied PULSE to study rare genetic variants associated with COVID-19 severity
364 based on the GEN-COVID cohort of whole-exome sequencing from 1,339 COVID-19
365 patients with 5 severity gradings⁵² (**Fig. 4b**, **Supplementary Table 8**; **Methods**). After
366 quality controls (QCs), we constructed a discovery cohort (training dataset) of
367 non-elderly European (EUR) adults (age >30 and <60 years) who were critically ill
368 (cases, $n=109$) or not hospitalized (controls, $n=269$) (**Methods**). There is no significant
369 age difference between cases and controls after filtration ($P=0.29$, Wilcoxon rank-sum

370 test; **Supplementary Fig. 5**). We then performed genome annotation⁵³ and feature
371 engineering (**Methods**), where only rare variants (i.e., absent from the EUR cohort
372 within the 1000 Genomes Project Phase 3⁵⁴) in autosomes were utilized for downstream
373 analysis. To test the prediction performance of PULSE, we first performed 5-fold
374 cross-validation (CV) based on the GEN-COVID discovery cohort, where a mean
375 AUROC (area under the receiver operating characteristics) of 0.631 was achieved
376 (SE=0.062; **Fig. 4c**). This demonstrates the predictability of COVID-19 severity from
377 personal genomes. The AUROC scores (0.629±0.073) of a logistic regression model
378 built from patient age and sex information are comparable to the PULSE genetic model
379 (**Fig. 4c**). However, combining scores (by averaging) of PULSE and age+sex produced
380 a further improvement in prediction performance (AUROC=0.653±0.072; **Fig. 4c**),
381 demonstrating that host genetics is relatively independent of the effect of age and sex
382 on disease severity. We note that since we removed the age bias in the discovery
383 cohort for genetic concentration, the largest contribution in the age+sex model came
384 from the sex information (model coefficients: 1.33±0.052 for sex versus 0.001±0.004 for
385 age; **Supplementary Fig. 6**). We also note that in our PULSE model only autosomal
386 variants were considered to remove the effect of sex in genetic modelling.

387 To further validate the prediction power of PULSE, we analyzed whole-genome
388 sequencing (WGS) data of an independent cohort from the Veterans Health
389 Administration (VA), consisting of 590 COVID-19 patients with variable disease severity
390 (**Fig. 4d, Supplementary Table 9; Methods**). Extensive QCs (without filtering based on
391 ancestry) resulted in 571 genomes (**Methods**). Genome annotation and feature
392 engineering were conducted as for the GEN-COVID cohort. Similarly, to remove the
393 effect of age, we focused on non-elderly adults (age >30 and <65 years) who were
394 critically ill or not hospitalized, yielding 243 samples (24 cases and 219 controls). In this
395 analysis, we relaxed the upper threshold of age from 60 to 65 years to include more
396 samples in testing. The PULSE model trained on the whole GEN-COVID cohort was
397 applied to predict severity within the VA EUR samples (14 cases versus 125 controls).
398 We found that PULSE succeeded in predicting severe disease solely from personal
399 genomes for this independent cohort with an AUROC of 0.675 (**Fig. 4e**).

400 Next, we asked if the prediction accuracy is generalizable across different populations.
401 We constructed a test set of non-EUR non-elderly adults (age >30 and <60 years) that
402 passed all other QC criteria within the GEN-COVID cohort, resulting in 12 cases
403 (critically ill) and 6 controls (not hospitalized). The PULSE model trained on EUR
404 samples was then applied to this non-EUR dataset, yielding AUROC of 0.667, which is
405 comparable to the prediction solely based on age and sex (AUROC=0.722; **Fig. 4e**).
406 Combining two scores further increased the prediction performance (AUROC=0.799;
407 **Fig. 4e**). Furthermore, we applied the same trained PULSE model to predict severe
408 disease for African (AFR) individuals (10 cases versus 92 controls) within the VA cohort.
409 Similarly, we discovered that PULSE succeeded in the cross-population prediction with
410 an AUROC of 0.784 (**Fig. 4e**). A similar result was observed for the whole VA dataset
411 with mixed populations (AUROC=0.716; **Fig. 4e**). These results demonstrate the
412 prediction power of PULSE and suggest that the rare-variant genetic architecture of
413 COVID-19 severity is conserved across multiple populations. Importantly, we observed
414 that the prediction in the VA cohort based on just age and sex information trained on the
415 GEN-COVID cohort is inferior to PULSE (AUROC=0.655, 0.474, and 0.577 for EUR,
416 AFR, and all samples, respectively; **Fig. 4e**). This may be linked to the different sex
417 distribution with fewer females in the VA cohort (**Fig. 4d, Supplementary Figs. 7 and**
418 **8**), but is further evidence of the robustness of host genetic signals in determining
419 COVID-19 severity and demonstrates that the PULSE prediction is independent of age
420 and sex.

421 We investigated additional performance measures including sensitivity and specificity
422 based on different cutoffs. Importantly, although the specificity scores are comparable
423 between PULSE and age+sex models cross cutoffs (**Supplementary Fig. 9**), we
424 discovered that PULSE yielded a significantly higher sensitivity (median values: 0.857
425 versus 0.688, 0.900 versus 0.525, and 0.875 versus 0.560 for EUR, AFR, and all VA
426 samples, respectively; **Fig. 4f**). High sensitivity is important for the clinical application of
427 severity prediction to guide the identification of at-risk individuals. Predictions for the
428 GEN-COVID non-EUR samples yielded similar sensitivity and specificity
429 (**Supplementary Fig. 10**).

430 **Common and rare variant analyses of severe COVID-19 converge on NK cell** 431 **function**

432 The trained PULSE model assigns a weighting to individual genes as a measure of
433 association between gene function and severe COVID-19 (**Supplementary Fig. 11**;
434 **Methods**), where a larger weight indicates a higher gene mutation burden in severe
435 patients. To test for the convergence between our common and rare variant analyses,
436 we compared the absolute values of model weights for RefMap genes per cell type with
437 all genes considered by the PULSE model. After correcting for multiple testing, we
438 concluded that for all cell types, RefMap genes tend to have weights with larger
439 absolute values, indicating an association with severe COVID-19 (FDR<0.01, one-tailed
440 Wilcoxon rank-sum test; **Fig. 5a**). Common and rare genetic variations are largely
441 independent⁵⁵⁻⁵⁷, and therefore, this convergence of common and rare variant signals
442 indicates shared biology underlying severe disease. Among all cell types, club-cell
443 RefMap genes are the most enriched with PULSE genes (FDR<1e-05, Z-score=5.33;
444 **Fig. 5a**). Interestingly, of hematopoietic cells, NK cell genes are the most enriched with
445 PULSE genes (FDR=1.1e-03, Z-score=3.16; **Fig. 5a**), consistent with our previous
446 conclusion that NK cells are an essential component of the immune response against
447 SARS-CoV-2.

448 To further validate the importance of genes captured by PULSE for severe COVID-19,
449 we identified 657 genes (referred to as PULSE COVID-19 genes) based on model
450 weights in the top 5% from all genes (**Supplementary Table 10**). We re-examined our
451 SKAT burden analysis results for the GEN-COVID cohort and observed that PULSE
452 genes are significantly enriched with severe-disease-associated rare variants (median
453 $P < 1e-5$, permutations test). Similar enrichment was confirmed based on an independent
454 rare-variant burden analysis from Regeneron¹², where the PULSE genes are
455 significantly enriched with Regeneron genes implicated in the analysis of severe
456 COVID-19 versus non-hospitalized COVID-19 ($n=68$ genes; $P=0.02$, OR=2.9, Fisher's
457 exact test; **Methods**).

458 To gain further insights into the cell-type-specificity of PULSE genes, we investigated
459 their expression levels in healthy lungs per cell type. We confirmed the function of

460 PULSE genes across different cell types by observing their non-random overlapping
461 with lung snATAC-seq peaks³³ (FDR<0.1, permutation test). Furthermore, the
462 expression levels of PULSE genes measured by scRNA-seq were examined, where we
463 found that PULSE genes are higher expressed in B cells, club cells, lymphatics, matrix
464 fibroblast 1, and NK cells (FDR<0.1, one-tailed Wilcoxon rank-sum test; **Fig. 5b**). This
465 supports the functional importance of PULSE genes in lung function.

466 PULSE genes carry a higher mutation burden in severe COVID-19 and therefore we
467 hypothesized that loss of function of PULSE genes leads to severe symptoms. To
468 validate this, we analyzed the expression levels of PULSE genes based on the
469 scRNA-seq data of COVID-19 patients²³. Consistent with our hypothesis, we observed a
470 down-regulation of PULSE genes in severe disease compared to moderate disease
471 across B cells, ciliated cells, macrophages, NK cells, and T cells (FDR<0.1, one-tailed
472 Wilcoxon rank-sum test; **Fig. 5c**). A similar analysis in another cohort²² led to the same
473 conclusion for macrophages, NK cells, and T cells (FDR<0.1, one-tailed Wilcoxon
474 rank-sum test; **Fig. 5d**). Our transcriptome study demonstrates the functional role of
475 PULSE genes in severe disease across multiple cell types. Notably, among all the cell
476 types we investigated, only NK cells are consistently associated with severe COVID-19
477 across all observations. This supports the conclusion of our common variant analysis,
478 suggesting that NK cells are vital determinants of COVID-19 severity.

479 **Systems analysis implicates association of NK cell activation with COVID-19** 480 **severity**

481 All of our analyses have suggested that NK cell dysfunction is a determinant of
482 COVID-19 severity. To obtain a comprehensive landscape of NK cell biology underlying
483 severe COVID-19, we examined the function of NK-cell genes identified by either
484 RefMap or PULSE (377 genes; **Supplementary Table 11**). Indeed, genes do not
485 function in isolation^{58,59} and therefore, rather than examining individual genes, we
486 mapped NK-cell genes to the global protein-protein interaction (PPI) network and then
487 inspected functional enrichment of COVID-19-associated network modules.

488 In particular, we extracted high-confidence (combined score >700) PPIs from STRING
489 v11.0⁶⁰, which include 17,161 proteins and 839,522 protein interactions. To eliminate the
490 bias of hub genes⁶¹, we performed the random walk with restart algorithm over the raw
491 PPI network to construct a smoothed network based on edges with weights in the top
492 5% (**Supplementary Table 12; Methods**). Next, this smoothed PPI network was
493 decomposed into non-overlapping subnetworks using the Leiden algorithm⁶². This
494 process yielded 1,681 different modules (**Supplementary Table 13**), in which genes
495 within modules are densely connected but sparsely connected with genes in other
496 modules.

497 NK-cell COVID-19 genes were mapped to individual modules, and four modules were
498 found to be significantly enriched with NK-cell genes: M237 ($n=471$ genes; $FDR<0.1$,
499 hypergeometric test; **Fig. 6a**), M1164 ($n=396$ genes; $FDR<0.1$, hypergeometric test;
500 **Fig. 6b**), M1311 ($n=14$ genes; $FDR<0.1$, hypergeometric test), and M1540 ($n=226$
501 genes; $FDR<0.1$, hypergeometric test; **Fig. 6c**) (**Supplementary Table 13**). We
502 excluded M1311 from our downstream analysis due to its limited size and lack of
503 functional enrichment.

504 Functionally, M237, M1164, and M1540 are all enriched with gene expression linked to
505 NK cells ($P<0.05$, Human Gene Atlas), demonstrating their specificity in the NK cell
506 function. Moreover, these three modules relate to different stages of NK cell activation.
507 M237 is enriched with GO/KEGG terms including ‘mRNA processing (GO:0006397)’
508 and ‘Spliceosome’, which are important for the transcriptional response involved in NK
509 cell activation (adjusted $P<0.1$; **Fig. 6d, Supplementary Tables 14 and 15**). M1164 is
510 enriched with GO/KEGG terms linked to intracellular signalling (e.g., ‘regulation of small
511 GTPase mediated signal transduction (GO:0051056)’), including pathways (e.g., ‘Rap1
512 signalling pathway’) key for NK cell activation (adjusted $P<0.1$; **Fig. 6e, Supplementary**
513 **Tables 14 and 15**). M1540 is highly enriched with GO/KEGG terms linked to type I
514 interferon signalling (e.g., ‘type I interferon signalling pathway (GO:0060337)’ and
515 ‘Antigen processing and presentation’) (adjusted $P<0.1$; **Fig. 6f, Supplementary Tables**
516 **14 and 15**). In summary, the functional enrichment of M237, M1164, and M1540 genes
517 includes extracellular, cytoplasmic, and nuclear processes necessary for NK cell

518 activation; thus the genetic architecture we have discovered places NK cell activation
519 upstream in determining severe COVID-19.

520 To further characterize the function of the identified NK-cell modules, we investigated
521 the expression of module genes based on scRNA-seq data from healthy and diseased
522 lung tissues. Genes in all three modules are relatively over-expressed in NK cells of
523 healthy lungs³³ than the background transcriptome (FDR<0.01, one-tailed Wilcoxon
524 rank-sum test; **Fig. 6g**). In contrast, in lung tissues infected with SARS-CoV-2, we
525 observed a down-regulation of M237 and M1540 genes in NK cells of severe disease²³
526 (FDR<0.01, one-tailed Wilcoxon rank-sum test; **Fig. 6h**). M1164 genes are also
527 down-regulated in NK cells from severe COVID-19 patients in another cohort²²
528 (FDR<0.01, one-tailed Wilcoxon rank-sum test; **Fig. 6i**) along with M237 and M1540
529 genes. These results are consistent with our previous findings and functionally link the
530 modules we have detected to NK cell biology in the context of severe COVID-19.

531 **DISCUSSION**

532 The COVID-19 pandemic is a global health crisis¹. Vaccination efforts have led to early
533 successes⁶³, but the prospect of evolving variants capable of immune-escape⁶⁴
534 highlights the importance of efforts to better understand the COVID-19 pathogenesis
535 and to develop effective treatments. Host genetic determinants of disease severity have
536 been investigated⁵⁻¹², but the findings and functional interpretations so far have been
537 limited¹³. In contrast, studies of the immune response accompanying severe
538 COVID-19¹⁶⁻¹⁹ have struggled to establish causality leading to a diverse array of
539 candidates and little consensus. Our contribution is an integrated analysis of common
540 and rare host genetic variation causally linked to severe COVID-19 in non-elderly
541 adults, together with biological interpretations via single-cell omics profiling of lung
542 tissue, and identification of >1,000 risk genes.

543 Our study of common and rare genetic variation associated with severe COVID-19
544 converges on common biology, despite non-overlapping datasets and orthogonal
545 analytical methods. We have achieved this because we have developed effective
546 machine learning methods which offer advantages over traditional methods: RefMap to

547 integrate common variants with epigenetic profiles³², and PULSE for rare variant
548 discovery by prediction. The evolution of clinical COVID-19 involves the interaction of
549 multiple viral and host factors in what is likely to be a nonlinear system; our work
550 supports this proposal and suggests that traditional methods may be inadequate given
551 current sample sizes. This study is the first time we have presented PULSE and we
552 have demonstrated a significant power advantage compared to standard methodology.
553 Both methods are ready for application in other disease areas.

554 Our network analysis highlights NK cell activation through type I interferon signalling
555 (**Fig. 6f**) as a key upstream determinant of COVID-19 severity. This links to previous
556 literature describing a delayed interferon response as a precursor of later
557 hyperinflammation associated with potentially fatal ARDS^{27,65}. NK cells can also be
558 activated via MHC signalling through NKG2 proteins. The CD94/NKG2C/HLA-E axis
559 has been shown to be key to the NK antiviral response⁶⁶ but so has the recognition of
560 induced-self antigens via the NKG2D receptor⁶⁷. Deletions of NKG2C have previously
561 been linked to severe COVID-19²⁸, whereas both our Mendelian randomization and
562 transcriptome analyses highlight a role for NKG2D+ NK cells. We suggest that all three
563 mechanisms for NK cell activation are critical to the host immune response to
564 SARS-CoV-2. Indeed, a recent study has revealed that autoantibodies which impair NK
565 cell activation are associated with severe COVID-19, and that manipulating the
566 activation of NK cells in a mouse model resulted in a significantly higher viral burden²⁹.
567 In the cancer field, NK cell stimulation has been postulated as a therapeutic strategy⁶⁸.
568 We propose that this strategy could protect at-risk individuals in future waves of
569 COVID-19.

570 It is important to note that our analyses also identified genetic risk of severe COVID-19
571 associated with non-NK cell types, including other immune cells and epithelial cells such
572 as AT2 cells, which is consistent with the previous literature⁶⁹. Indeed, the PULSE
573 prediction is based on a total genetic architecture and not limited to NK cell
574 genomics. Future work will determine how these other cell types are essential and how
575 they interact with NK cell activation.

576 We present a validated prediction of COVID-19 severity derived entirely from host
577 characteristics, including age, sex, and genetics. The average AUROC of ~0.72
578 outperforms all comparable strategies⁹; and we achieve a very high sensitivity of ~85%
579 with a specificity >50%. Our prediction could be applied in advance of infection or even
580 exposure, and thus has the potential to be very useful clinically. We anticipate future use
581 and refinement of our prediction model to guide administration of post-exposure
582 prophylaxis to at-risk individuals, in a similar manner to current standard practice for
583 HIV⁷⁰.

584 Our analyses are based on the largest available datasets to date but increasing sample
585 size could improve the precision of our discovery and prediction. In addition, the vast
586 majority of our data was taken from populations and at times when recently identified
587 SARS-CoV-2 variants were not prevalent in the population (before November 2020,
588 <https://covariants.org/per-country>). It is unlikely, but not impossible, that the NK cell
589 responses we have identified as essential determinants of severe COVID-19 are not
590 applicable to new variants.

591 In conclusion, we have uncovered a comprehensive genetic architecture of severe
592 COVID-19 integrated with single-cell-resolution biological functions. Both common and
593 rare variant analyses have highlighted NK cell activation as a potential key factor in
594 determining disease severity. Our novel rare variant method has also achieved age-,
595 sex-, and ancestry-independent prediction of COVID-19 severity from personal
596 genomes.

597 **FIGURES**

598 **Figure 1. Common variant analysis of COVID-19 severity integrated with lung**
599 **single-cell multiomics.**

600 **a**, Schematic of the study design for fine-mapping cell-type-specific genes from
601 COVID-19 GWAS (Panel 1). The diagram of the RefMap model is shown in Panel 2,
602 where grey nodes represent observations, green nodes are local hidden variables, and
603 pink nodes indicate global hidden variables (**Methods**). Cell-type-specific RefMap

604 genes are mapped using single-cell multiomic profiling (Panel 3). Heritability (Panel 4),
605 Mendelian randomization (Panel 5), and transcriptome analysis (Panel 6) validate the
606 functional importance of RefMap genes, particularly for NK cells, in severe COVID-19.
607 **b**, Total number and number of unique genomic regions containing genetic variation
608 associated with severe COVID-19 for different cell types. **c**, Total number and number of
609 unique genes implicated by genetic variation associated with severe COVID-19 for
610 different cell types. **d**, Fraction of unique genomic regions and genes associated with
611 severe COVID-19 for major cell types. **e**, Similarity between different cell types
612 quantified by the overlap of RefMap genes. Gene set overlapping was calculated by the
613 Jaccard index. **f**, RefMap regions overlap significantly with COVID-19-associated
614 genetic variation in an independent COVID-19 GWAS study. cCRE: candidate
615 cis-regulatory element. *: $P < 0.05$.

616 **Figure 2. Severe-COVID-19-associated common variants are linked to NK cell**
617 **function.**

618 **a**, Heritability enrichment estimated by LDSC for different cell types. Enrichment was
619 calculated as the proportion of total SNP-based heritability adjusted for SNP number. **b**,
620 Proportion of SNP-based heritability associated with risk genes identified using RefMap
621 or conventional methodology. **c**, **d**, **e**, Significant Mendelian randomization results for
622 three exposures linked to severe COVID-19, including blood counts of (**c**) CD335+
623 CD314-, (**d**) CCR7- CD314-, and (**e**) CD314+ NK cells. **f**, Sensitivity analyses and
624 robust tests for MR analyses (**Methods**). **g**, Comparative gene expression analysis of
625 NK-cell RefMap genes in NKG2D+ and NKG2D- NK cells. Fold change was calculated
626 as the ratio of gene expression levels in NKG2D+ NK cells to NKG2D- NK cells. The
627 transcriptome was defined by all the expressed genes (with at least one UMI (unique
628 molecular identifier)) in NK cells. Violin plots show the distributions of fold change
629 values within each group, and boxplots indicate the median, interquartile range (IQR),
630 $Q1 - 1.5 \times IQR$, and $Q3 + 1.5 \times IQR$. The red dashed line denotes the median value of fold
631 change distribution for the transcriptome.

632 **Figure 3. Functional enrichment and transcriptome analyses of RefMap COVID-19**
633 **genes.**

634 **a**, Gene Ontology (GO) terms that are significantly enriched in cell-type-specific RefMap
635 gene lists corresponding to hematopoietic cell types; only terms with adjusted $P < 0.05$,
636 $OR > 3$, and character number < 60 are visualized. **b**, KEGG Pathways that are
637 significantly enriched in cell-type-specific RefMap gene lists corresponding to
638 hematopoietic cell types; only terms with adjusted $P < 0.05$, $OR > 5$, and character
639 number < 50 are visualized. **c**, Gene expression analysis of RefMap genes across
640 different cell types in healthy lungs. The transcriptome was defined as the total set of
641 expressed genes for each cell type (**Methods**). Violin plots show the distributions of log
642 expression levels within each group, and point plots indicate the median and IQR. **d**,
643 Overlap between cell-type-specific RefMap regions and H3K27ac and H3K4me3
644 ChIP-seq peaks from ENCODE lung and immune cell samples. Z-scores calculated by
645 regionR⁴⁹ (1,000 permutations) were normalized into the 0-1 range for visualization. **e, f**,
646 Comparative gene expression analysis of cell-type-specific RefMap genes in severe
647 COVID-19 patients versus moderately affected patients based on scRNA-seq datasets
648 from **(e)** Ren et al. and **(f)** Liao et al., respectively. The Z-score of Wilcoxon rank-sum
649 test was used to indicate the gene expression change between severe and moderate
650 patient groups, where a positive value means higher gene expression in severe
651 patients. The Benjamini-Hochberg (BH) procedure was used to calculate FDRs
652 throughout the study. Violin plots show the distribution of gene expression changes
653 within each group, and boxplots indicate the median, IQR, $Q1 - 1.5 \times IQR$, and
654 $Q3 + 1.5 \times IQR$. *: $FDR < 0.1$. +: $FDR < 0.01$.

655 **Figure 4. Rare variant analysis informs individual risk of critical illness of**
656 **COVID-19.**

657 **a**, Enrichment analysis of cell-type-specific RefMap COVID-19 genes with rare variants
658 using SKAT burden testing. The red dashed line indicates $P = 0.05$. **b**, Schematic of the
659 study design for our rare variant analysis based on PULSE. We examine two
660 independent cohorts in which rare variants were profiled by different technologies:

661 whole-exome sequencing (WES) and whole-genome sequencing (WGS) (Panel 1).
662 Variants are annotated using ANNOVAR (Panel 2) and encoded as input for the PULSE
663 model (Panel 3, **Methods**), where grey nodes are observations and pink nodes
664 represent hidden variables. PULSE is trained to differentiate cases and controls (Panel
665 4), where the gene weights are useful for gene discovery (Panel 5). Functional
666 characterization of risk genes is performed based on scRNA-seq and PPIs (Panel 6). **c**,
667 Receiver operating characteristic (ROC) curves of different models, including PULSE,
668 age+sex, and integrative models, in the 5-fold cross-validation. Solid lines represent the
669 mean values, and the grey area indicates the standard errors. **d**, Summary statistics of
670 the VA COVID-19 cohort. **e**, AUROC (area under the receiver operating characteristics)
671 scores of predictions in multiple test datasets. Prediction performance is shown for
672 PULSE, age+sex, and integrative models. **f**, Comparison of prediction sensitivity
673 between PULSE and age+sex models. EHRs: electronic health records.

674 **Figure 5. Transcriptome analysis of PULSE COVID-19 genes.**

675 **a**, Analysis of convergence between PULSE and RefMap COVID-19 genes. The
676 Z-scores were calculated per cell-type by Wilcoxon rank-sum test of the difference in
677 PULSE weights between RefMap genes and the background transcriptome. Non-zero
678 Z-scores indicate biological overlap between common and rare variant architectures
679 detected by RefMap and PULSE, respectively. **b**, Gene expression analysis of PULSE
680 genes across different cell types in healthy lungs. The transcriptome was defined as the
681 total set of expressed genes for each cell type (**Methods**). Violin plots show the
682 distributions of log expression levels within each group, and point plots indicate the
683 median and IQR. **c**, **d**, Comparative gene expression analysis of cell-type-specific
684 PULSE genes in severe COVID-19 patients versus moderate patients based on
685 scRNA-seq datasets from (**c**) Ren et al. and (**d**) Liao et al., respectively. The Z-score of
686 Wilcoxon rank-sum test was used to indicate the gene expression change between
687 severe and moderate patient groups. Violin plots show the distribution of gene
688 expression changes within each group, and boxplots indicate the median, IQR,
689 $Q1-1.5 \times IQR$, and $Q3+1.5 \times IQR$. *: FDR<0.1. +: FDR<0.01.

690 **Figure 6. Network analysis of NK-cell genes identified in common and rare variant**
691 **analyses.**

692 **a, b, c**, Three PPI network modules, including (a) M237, (b) M1164, and (c) M1540, are
693 significantly enriched with NK-cell genes identified in either common or rare variant
694 analysis. Blue nodes represent NK-cell genes and yellow nodes indicate other genes
695 within each module. Edge thickness is proportional to STRING confidence score (>700).
696 **d, e, f**, Gene Ontology (GO) terms that are significantly enriched in modules (d) M237,
697 (e) M1164, and (f) M1540. Selected terms are shown for visualization and the complete
698 lists can be found in **Supplementary Tables 14 and 15**. **g**, Gene expression analysis of
699 module genes in NK cells. The transcriptome was defined as the total set of expressed
700 genes in NK cells (**Methods**). Violin plots show the distributions of log expression levels
701 within each group, and boxplots indicate the median, IQR, $Q1-1.5\times IQR$, and
702 $Q3+1.5\times IQR$. The red dashed line indicates the median expression level of the
703 transcriptome. **h, i**, Comparative gene expression analysis of module genes in severe
704 COVID-19 patients versus moderate patients based on scRNA-seq datasets from (h)
705 Ren et al. and (i) Liao et al., respectively. The *Z*-score of Wilcoxon rank-sum test was
706 used to indicate the gene expression change between severe and moderate patient
707 groups. Violin plots show the distribution of gene expression changes within each
708 group, and boxplots indicate the median, IQR, $Q1-1.5\times IQR$, and $Q3+1.5\times IQR$. The red
709 dashed line indicates the median expression change of the transcriptome. +: $FDR < 0.01$.
710 GOBP: gene ontology biological process.

711 METHODS

712 The RefMap model

713 Allele Z-scores were calculated as $Z=b/se$, where b and se are effect size and standard
 714 error, respectively, as reported by the COVID-19 GWAS⁷ (COVID-19 Host Genetics
 715 Initiative, Release 5, phenotype definition A2, EUR only) where the sample age, sex,
 716 and ancestry information were included as covariates. Given Z-scores and lung
 717 snATAC-seq peaks, we aim to identify functional genomic regions in which the Z-score
 718 distribution is significantly shifted from the null distribution. Suppose we have K 1Mb
 719 linkage disequilibrium (LD) blocks, where each LD block contains J_k ($k=1, \dots, K$) 1kb
 720 regions and each region harbors $I_{j,k}$ ($j=1, \dots, J_k, I_{j,k}>0$) SNPs, the Z-scores follow a
 721 multivariate normal distribution, i.e.,

$$722 \quad \mathbf{z}_k | \mathbf{u}_k \sim \mathcal{N}(\Sigma_k \mathbf{u}_k, \Sigma_k), \quad k=1, \dots, K, \quad (1)$$

723 in which the Z-score of the i -th SNP in the j -th region of the k -th block is denoted as $z_{i,j,k}$
 724 ($i=1, \dots, I_{j,k}$) and \mathbf{u}_k are the effect sizes that can be expressed as

$$725 \quad \mathbf{u}_k = \left[\mathbf{u}_{1:I_{1,k},1,k}^T, \dots, \mathbf{u}_{1:I_{j,k},j,k}^T, \dots, \mathbf{u}_{1:I_{J_k,k},J_k,k}^T \right]^T. \quad (2)$$

726 In addition, $\Sigma_k \in \mathbb{R}^{I_k \times I_k}$ in Eq. (1) represents the in-sample LD matrix comprising of the
 727 pairwise Pearson correlation coefficients between SNPs within the k -th block, where I_k is
 728 the total number of SNPs calculated by $I_k = \sum_{j=1}^{J_k} I_{j,k}$. Here, since we have no access to
 729 the individual-level data, we used EUR samples from the 1000 Genomes Project
 730 (Phase 3) to estimate Σ_k , yielding the out-sample LD matrix. A modified Cholesky
 731 algorithm⁷¹ was used to get a symmetric positive definite (SPD) approximation of the LD
 732 matrix.

733 Further, we assume $u_{i,j,k}$ ($i=1, \dots, I_{j,k}$) are independent and identically distributed (i.i.d.),
 734 following a normal distribution given by

$$735 \quad u_{i,j,k} | m_{j,k}, \lambda_{j,k} \sim \mathcal{N}(m_{j,k}, \lambda_{j,k}^{-1}), \quad i=1, \dots, I_{j,k}, \quad (3)$$

736 where the precision $\lambda_{j,k}$ follows a Gamma distribution, i.e.,

737

$$\lambda_{j,k} \sim \text{Gamma}(a_0, b_0) . \quad (4)$$

738 Moreover, to characterize the shift of the expectation in Eq. (3) from the null distribution,
739 we model $m_{j,k}$ by a three-component Gaussian mixture model given by

$$740 \quad m_{j,k} | t_{j,k}, v_{-1}, v_{+1}, \tau_0, \tau_{-1}, \tau_{+1} \sim \underbrace{\mathcal{N}(-v_{-1}, \tau_{-1}^{-1})^{t_{j,k}^{(-1)}}}_{\text{negative}} \underbrace{\mathcal{N}(0, \tau_0^{-1})^{t_{j,k}^{(0)}}}_{\text{zero}} \underbrace{\mathcal{N}(v_{+1}, \tau_{+1}^{-1})^{t_{j,k}^{(+1)}}}_{\text{positive}} , \quad (5)$$

741 where the precisions follow

$$742 \quad \tau_{-1}, \tau_0, \tau_{+1} \sim \text{Gamma}(a_0, b_0) , \quad (6)$$

743 and v_{-1} and v_{+1} are non-negative variables measuring the absolute values of effect size
744 shifts for the negative and positive components, respectively.

745 To impose non-negativity over v_{-1} and v_{+1} , we adopt the rectification nonlinearity
746 technique proposed previously⁷². In particular, we assume v_{-1} and v_{+1} follow

$$747 \quad v_{-1} | m_{-1}, \lambda_{-1} \sim \mathcal{R}^N(m_{-1}, \lambda_{-1}) , \quad (7)$$

$$748 \quad v_{+1} | m_{+1}, \lambda_{+1} \sim \mathcal{R}^N(m_{+1}, \lambda_{+1}) , \quad (8)$$

749 in which the rectified Gaussian distribution is defined via a dumb variable. In particular,
750 we first define v_{-1} and v_{+1} by

$$751 \quad v_{-1} = \max(r_{-1}, 0) , \quad (9)$$

$$752 \quad v_{+1} = \max(r_{+1}, 0) , \quad (10)$$

753 which guarantees that v_{-1} and v_{+1} are non-negative. The dumb variable r_{-1} and r_{+1} follow
754 the Gaussian distributions given by

$$755 \quad r_{-1} | m_{-1}, \lambda_{-1} \sim \mathcal{N}(m_{-1}, \lambda_{-1}^{-1}) , \quad (11)$$

$$756 \quad r_{+1} | m_{+1}, \lambda_{+1} \sim \mathcal{N}(m_{+1}, \lambda_{+1}^{-1}) , \quad (12)$$

757 where m_{\pm} and λ_{\pm} follow the Gaussian-Gamma distributions, i.e.,

$$758 \quad m_{-1}, \lambda_{-1} \sim \mathcal{N}(\mu_0, (\beta_0 \lambda_{-1})^{-1}) \text{Gamma}(a_0, b_0) , \quad (13)$$

$$759 \quad m_{+1}, \lambda_{+1} \sim \mathcal{N}(\mu_0, (\beta_0 \lambda_{+1})^{-1}) \text{Gamma}(a_0, b_0) . \quad (14)$$

760 The indicator variables in Eq. (5) denote whether that region is disease-associated or
 761 not. Indeed, we define the region to be disease-associated if $t_{j,k}^{(-1)} = 1$ or $t_{j,k}^{(+1)} = 1$, and
 762 to be non-associated otherwise. To simplify the analysis, we put a symmetry over $t_{j,k}^{(-1)}$
 763 and $t_{j,k}^{(+1)}$, and define the distribution by

$$764 \quad p(t_{j,k} | \pi_{j,k}) = (0.5\pi_{j,k})^{t_{j,k}^{(-1)}} (1 - \pi_{j,k})^{t_{j,k}^{(0)}} (0.5\pi_{j,k})^{t_{j,k}^{(+1)}}, \quad j=1, \dots, J_k, \quad k=1, \dots, K. \quad (15)$$

765 Furthermore, the probability parameter $\pi_{j,k}$ in Eq. (15) is given by

$$766 \quad \pi_{j,k} = \sigma(\mathbf{w}^T \mathbf{s}_{j,k}), \quad (16)$$

767 where $\sigma(\cdot)$ is the sigmoid function, $\mathbf{s}_{j,k}$ is the vector of epigenetic features for the j -th
 768 region in the k -th LD block, and the weight vector \mathbf{w} follows a multivariate normal
 769 distribution, i.e.,

$$770 \quad \mathbf{w} | \Lambda \sim \mathcal{N}(\mathbf{0}, \Lambda^{-1}), \quad (17)$$

771 and Λ follows

$$772 \quad \Lambda \sim \mathcal{W}(\mathbf{W}_0, \nu_0). \quad (18)$$

773 In this study, the epigenetic feature $\mathbf{s}_{j,k}$ was calculated as the overlapping ratios of that
 774 region with the snATAC-seq peaks detected in any of the cell types in healthy human
 775 lungs.

776 Based on the model defined in Eqs. (1) to (18), we are interested in calculating the
 777 posterior probability $p(\mathbf{T} | \mathbf{Z}, \mathbf{S})$, where the mean-field variational inference (MFVI)⁷³ was
 778 adopted to solve the intractability. More technical details, including a coordinate
 779 ascent-based inference algorithm, can be found in our previous work³².

780 In this study, we ran the MFVI algorithm per chromosome to accelerate the
 781 computation. The Q^+ - and Q^- -scores were defined as $q(t^{(+1)} = 1)$ and $q(t^{(-1)} = 1)$,
 782 respectively, and we also defined the Q -score as $Q = Q^+ + Q^-$. RefMap regions were
 783 identified by Q^+ - or Q^- -score > 0.95 .

784 Mapping cell-type-specific genes from RefMap regions

785 For each cell type within lung tissue, we defined cell-type-specific RefMap regions as
786 the overlap between RefMap regions and the total set of snATAC-seq peaks detected in
787 that cell type (**Supplementary Table 1**). Cell-type-specific RefMap genes were then
788 identified if the extended gene body (i.e., the region up to 10kb either side of the
789 annotated gene body) overlapped with any of the cell-type-specific regions. To get the
790 final gene lists, RefMap genes were further filtered based on their expression levels. In
791 particular, with the lung snRNA-seq data³³, we defined expressed genes in each cell
792 type as those with Seurat⁷⁴ log-normalized value>0.6931. In addition, we note that there
793 are non-adult samples (~30 weeks gestation and ~3 years) sequenced in the single cell
794 profiling data³³. To remove the bias towards lung development, we first calculated the
795 fold change of gene expression levels between the adult sample (~30 years) and
796 non-adult ones, and defined non-developmental genes (nDG) as those with FC>1.5.
797 Only RefMap genes that were identified as expressed and non-developmental in each
798 cell type were kept for downstream analysis (**Supplementary Table 1**).

799 **Validation of RefMap COVID-19 regions in the 23andMe dataset**

800 We calculated the overlap of total RefMap regions and of cell-type-specific RefMap
801 regions with genomic regions shown to contain COVID-19-associated SNPs ($P<1e-04$)
802 based on the GWAS of an independent cohort recruited by 23andMe⁵. To determine
803 whether the observed overlap is statistically significant, we examined the average
804 overlap with ten sets of control regions of equivalent length to RefMap regions. Control
805 regions were +/-1Mb-5Mb distant from the RefMap regions⁷⁵.

806 **Heritability analysis**

807 We used LD score regression (LDSC)³⁴ to calculate overall heritability for severe
808 COVID-19 (A1), hospitalized COVID-19 (B2), and COVID-19 overall (C2), respectively.
809 Heritability partitioning within genes identified by traditional methods and within
810 cell-type-specific RefMap genes was performed as previously described⁷⁶. Briefly, for all
811 gene lists, we examined the proportion of total SNP-based heritability carried by SNPs
812 +/-100kb from the transcription start site (TSS) of each gene in the list. Enrichment was

813 calculated by comparing the ratio of partitioned heritability to the quantity of genetic
814 materials.

815 **Mendelian randomization**

816 In total, 46 GWAS measures of NK cell subtypes were identified from the IEU Open
817 GWAS Project, including "prot-a-180", "met-b-124", "met-b-245", "met-b-242",
818 "prot-c-5244_12_3", "met-b-237", "met-b-258", "prot-a-1669", "prot-c-2917_3_2",
819 "met-b-246", "prot-a-1671", "met-b-249", "met-b-140", "met-b-240", "prot-a-3159",
820 "prot-c-5104_57_3", "prot-c-3056_11_1", "prot-a-13", "prot-a-3160", "met-b-123",
821 "met-b-250", "met-b-239", "met-b-120", "met-b-154", "prot-a-3162", "met-b-247",
822 "met-b-251", "met-b-238", "met-b-243", "prot-a-2487", "met-b-244", "prot-c-2734_49_4",
823 "met-b-153", "prot-a-3161", "prot-c-3003_29_2", "met-b-248", "prot-a-1674",
824 "prot-a-1675", "met-b-152", "met-b-122", "met-b-121", "prot-a-1670",
825 "prot-c-5424_55_3", "met-b-252", "prot-a-3233" and "met-b-241"^{37,77,78}. Exposure SNPs
826 or instrumental variables (IVs) are chosen based on an arbitrary P -value cutoff^{79,80}. A
827 cutoff that is too low will lose informative instruments, but a cutoff that is too high could
828 introduce non-informative instruments. We chose to set the cutoff at 5e-06 in line with
829 our previous work⁸¹. We employed a series of sensitivity analyses to ensure that our
830 analysis was not confounded by invalid IVs. Identified SNPs were clumped for
831 independence using PLINK clumping in the TwoSampleMR tool⁸². A stringent cutoff of
832 $R^2 \leq 0.001$ and a window of 10,000kb were used for clumping within a European
833 reference panel. Where SNPs were in LD, those with the lowest P -value were retained.
834 SNPs that were not present in the reference panel were excluded. Where an exposure
835 SNP was unavailable in the outcome dataset, a proxy with a high degree of LD ($R^2 \geq 0.9$)
836 was identified in LDlink within a European reference population⁸³. Where a proxy was
837 identified to be present in both datasets, the target SNP was replaced with the proxy in
838 both exposure and outcome datasets in order to avoid phasing issues⁸⁴. Where a SNP
839 was not present in both datasets and no SNP was available in sufficient LD, the SNP
840 was excluded from the analysis. The effects of SNPs on outcomes and exposures were
841 harmonized in order to ensure that the beta values were signed with respect to the
842 same alleles. For palindromic alleles, those with minor allele frequency (MAF) > 0.42

843 were omitted from the analysis in order to reduce the risk of errors due to strand
844 issues⁸⁴.

845 The MR measure with the greatest power is the inverse-variance weighted (IVW)
846 method, but this is contingent upon the exposure IV assumptions being satisfied⁸⁵. With
847 the inclusion of a large number of SNPs within the exposure IV, it is possible that not all
848 variants included are valid instruments and therefore, in the event of a significant result,
849 it is necessary to include a range of robust methods which provide valid results under
850 various violations of MR principles at the expense of power⁸⁶. Robust methods applied
851 in this study include MR-Egger, MR-PRESSO, weighted median, weighted mode, and
852 MR-Lasso.

853 With respect to the IVW analysis, a fixed-effects (FE) model is indicated in the case of
854 homogeneous data, whilst a multiplicative random effects (MRE) model is more suitable
855 for heterogeneous data. Burgess et al. recommended that an MRE model be
856 implemented when using GWAS summary data to account for heterogeneity in
857 variant-specific causal estimates⁸⁶. In the interest of transparency, we calculated both
858 results but present the MRE in the text.

859 MR analyses should include evaluation of exposure IV strength. In order to achieve this,
860 we provided the F -statistic, MR-Egger intercept, MR-PRESSO global test, Cochran's Q
861 test, and I^2 for our data. The F -statistic is a measure of instrument strength with >10
862 indicating a sufficiently strong instrument⁸⁷. We provided F -statistics for individual
863 exposure SNPs and the instrument as a whole. Cochran's Q test is an indicator of
864 heterogeneity in the exposure dataset and serves as a useful indicator that horizontal
865 pleiotropy is present as well as directing decisions to implement FE or MRE IVW
866 approaches⁸⁸. The MR-Egger intercept test determines whether there is directional
867 horizontal pleiotropy. The MR-PRESSO global test determines if there are statistically
868 significant outliers within the exposure-outcome analysis⁸⁹. I^2 was calculated as a
869 measure of heterogeneity between variant specific causal estimates, with a low I^2
870 indicating that Egger is more likely to be biased towards the null⁹⁰. Finally, we performed
871 a leave-one-out analysis using the method of best fit for each exposure SNP within the

872 IV in order to determine if any single variants were exerting a disproportionate effect
873 upon the results of our analysis⁸⁶.

874 **MAGMA analysis of COVID-19 GWAS data**

875 MAGMA (v1.08)⁹¹ was applied using default settings. Input consisted of summary
876 statistics for all SNPs genome-wide as measured in the COVID-19 GWAS⁷. We
877 estimated LD structure using EUR samples from the 1000 Genomes Project (Phase 3).
878 The top 50 MAGMA genes³⁵ were used for downstream analysis.

879 **DNA sequencing in rare variant analysis**

880 GEN-COVID cohort. The cohort was recruited by the GEN-COVID consortium
881 (<https://sites.google.com/dbm.unisi.it/gen-covid>) as described previously⁵². Briefly, adult
882 patients (>18 years) were recruited from 35 Italian hospitals starting on March 16, 2020.
883 Infection status was confirmed by SARS-CoV-2 viral RNA polymerase-chain-reaction
884 (PCR) test collected at least from nasopharyngeal swabs. Demographics and clinical
885 severity were assessed via an extensive questionnaire.

886 Sequencing and variant calling were performed as described previously⁵². Briefly,
887 sample preparation was performed following the Nextera Flex for Enrichment
888 manufacturer protocol. Whole-exome sequencing was performed with >97% coverage
889 at 20X using the Illumina NovaSeq 6000 System (Illumina, San Diego, CA, USA).
890 Reads were aligned to human reference genome build GRCh38 using BWA⁹². Variants
891 were called according to the GATK4 best practice guidelines⁹³. Duplicates were
892 removed by *MarkDuplicates*, and base qualities were recalibrated using
893 *BaseRecalibration* and *ApplyBQSR*. *HaplotypeCaller* was used to calculate Genomic
894 VCF files for each sample, which were then used for multi-sample calling by
895 *GenomicDBImport* and *GenotypeGVCF*. In order to improve the specificity-sensitivity
896 balance, variant quality scores were calculated by *VariantRecalibrator* and *ApplyVQSR*.
897 Variants with sequencing depth <20X were excluded.

898 VA cohort. Whole-genome sequence data on the VA COVID-19 cohort was derived from
899 the VA Million Veteran Program (MVP). The VA MVP is an ongoing national voluntary

900 research program that aims to better understand how genetic, lifestyle, and
901 environmental factors influence veteran health⁹⁴. Briefly, individuals aged 18 to over 100
902 years old have been recruited from over 60 VA Medical Centers nationwide since 2011
903 with current enrollment at >800,000. Informed consent is obtained from all participants
904 to provide blood for genomic analysis and access to their full electronic health record
905 (EHR) data within the VA prior to and after enrollment. The study received ethical and
906 study protocol approval from the VA Central Institutional Review Board in accordance
907 with the principles outlined in the Declaration of Helsinki. COVID-19 cases were
908 identified using an algorithm developed by the VA COVID National Surveillance Tool
909 based on reverse transcription polymerase chain reaction laboratory test results
910 conducted at VA clinics, supplemented with natural language processing on clinical
911 documents for SARS-CoV-2 tests conducted outside of the VA⁹⁵.

912 DNA isolated from peripheral blood samples was used for whole-genome sequencing.
913 Libraries were prepared using KAPA hyper prep kits, PCR-free according to
914 manufacturers' recommendations. Sequencing was performed using Illumina NovaSeq
915 6000 System (Illumina, San Diego, CA, USA) with paired-end 2x150bp read lengths,
916 and Illumina's proprietary reversible terminator-based method. The specimens were
917 sequenced to a minimum depth of 25X per specimen and an average coverage of 30X
918 per plate.

919 WGS data processing in the MVP was performed via the functional equivalence GATK
920 variant calling pipeline⁹⁶, which was developed by the Broad Institute and plugged into
921 our data and task management system Trellis. The human reference genome build was
922 GRCh38. We used BWA-MEM (v0.7.15) to align reads, Picard 2.15.0 to mark PCR
923 duplicates, and GATK 4.1.0.0 for BQSR and variant calling via the *haplotypeCaller*
924 function. We also used FASTQC (v0.11.4), SAMTools *flagstat* (v0.1.19), and RTG Tools
925 *vcfstats* (v3.7.1) to assess the qualities of the FASTQ, BAM, and gVCF files,
926 respectively. In addition, we used *verifybamID* in GATK 4.1.0.0 to estimate DNA
927 contamination rates for individual genomes and removed samples with 5% or more
928 contaminated reads.

929 **Data quality control**

930 GEN-COVID cohort. To guarantee high quality of the sequencing data, we performed
931 numerous quality control procedures. On the sample level, we (1) computed inbreeding
932 coefficients (Fhat1, Fhat2, and Fhat3 in GCTA⁹⁷) and removed genomes that resided
933 more than 3 standard deviation from the mean; (2) computed identity-by-descent (IBD)
934 and only kept one genome from pairs with proportion IBD>0.2; (3) computed missing
935 calls for each genome and removed those with missing rate larger than 10%; (4)
936 computed singleton calls, SNV count, indel count, Ti/Tv ratio, and heterozygous calls for
937 each genome and removed genomes that resided more than 3 standard deviation from
938 the mean.

939 On the variant level, we (1) removed multiallelic sites; (2) kept variants in autosomes;
940 (3) removed variants on blacklisted regions, compiled by the ENCODE Project
941 Consortium (Phase 4); (4) removed variants identified other than “PASS,” such as “low
942 quality,” “tranche99.0-99.5,” by VQSR in GATK; (5) removed variants with missing rate
943 larger than 10%. The samples which passed QCs were provided in **Supplementary**
944 **Table 8**.

945 VA cohort. For deriving high-quality variants for downstream analysis, we removed
946 samples with kinship >0.03, sample call rate <0.97, or mean sample coverage <=18X.
947 Genomic positions resided in low complexity regions or ENCODE blacklisted regions
948 were first removed. Next, we filtered out genotypes in individual samples that were
949 detected with too low or too high of read coverages (DP<5 or >1500). We required all
950 calls to have genotype quality (GQ) >=20, and for non-reference calls, sufficient portion
951 (>0.9) of reads was required to cover the alternate alleles. In addition, we removed
952 genomic positions with cohort-wise call rate <0.95 and computed Hardy-Weinberg
953 equilibrium (HWE), which was required to be <1e-05 for common variants and <1e-06
954 for rare variants. With all these filtering completed, we assessed the sample-level
955 genomic parameters, such as Ti/Tv ratios, het/hom ratios, and number of
956 singletons/SNVs/INDELs, and removed any sample that fell into the tail regions of the
957 distribution (>=3 standard deviation). The samples which passed QCs were provided in
958 **Supplementary Table 9**.

959 **Ancestry analysis**

960 We performed population admixture analysis using ADMIXTURE⁹⁸ (v1.3.0) referencing
961 five super populations, including AFR (African), AMR (Ad Mixed American), EAS (East
962 Asian), EUR (European), and SAS (South Asian), in the 1000 Genomes Project (Phase
963 3) and inferred the ancestry for each genome. For the GEN-COVID cohort, samples
964 with >90% EUR ancestry fraction were kept in the discovery cohort. For the VA
965 COVID-19 cohort, we relaxed the ancestry fraction cutoff to 70% for including more
966 samples in testing. Inferred sample ancestry can be found in **Supplementary Tables 8**
967 **and 9**.

968 **Variant- and gene-level annotations**

969 Genome annotation was performed by Annovar⁵³ integrating multiple databases. Variant
970 frequency was estimated using the 1000 Genomes Project (Phase 3). Nonsynonymous
971 (missense and nonsense) variants were annotated using dbNSFP⁹⁹ (v3.5). The
972 mutation effect of splicing-site variants was predicted by dbSNV¹⁰⁰ (v1.1) and
973 regSNP-intron¹⁰¹.

974 **Rare-variant burden testing**

975 Rare-variant burden testing was performed to determine whether any genes were
976 differentially enriched with rare variants between severe COVID-19 patients and
977 non-severe COVID-19-positive controls. We utilized whole-exome sequencing data from
978 the GEN-COVID cohort⁵¹, including 122 individuals aged ≤ 60 years who suffered severe
979 COVID-19 requiring respiratory support, and 465 individuals aged ≥ 20 years who
980 suffered non-severe COVID-19 not requiring hospitalisation. Variants were included if
981 they altered an amino acid, were rare (MAF<1%) and absent from the EUR cohort of the
982 1000 Genomes Project (Phase 3). Burden was calculated using SKAT⁵⁰ adjusted for
983 sample imbalance using a saddlepoint approximation¹⁰². Sex and the first ten principal
984 components were included as covariates. Genetic burden was compared with the
985 complete set of coding genes; genes carrying <10 variants were removed because of
986 insufficient data. After filtering a total set of 4,280 genes were tested for
987 severe-COVID-19-associated rare genetic variation of which 625 were also RefMap

988 COVID-19 genes. A QQ-plot confirmed that there was no significant genomic inflation
989 ($\lambda_{GC}=1.1$; **Supplementary Fig. 4**).

990 Regeneron's burden testing results were obtained from the Regeneron results browser
991 (<https://rgc-covid19.regeneron.com/results>), where only semi-significant genes
992 ($P < 1e-03$, REGENIE¹⁰³) were available. Data consists of exome-wide association
993 studies of various COVID-19 outcomes across 662,403 individuals (11,356 with
994 COVID-19) aggregated from four studies: UK Biobank (UKB; $n=455,838$), AncestryDNA
995 COVID-19 Research Study ($n=83,930$), Geisinger Health System (GHS; $n=113,731$),
996 and Penn Medicine BioBank (PMBB; $n=8,904$). For the Regeneron study of severe
997 COVID-19 versus non-hospitalized COVID-19, we obtained a list of 68 genes harboring
998 disease-associated missense mutations at a significance cutoff of $P < 1e-03$ in EUR
999 samples. Overlap between PULSE genes and Regeneron gene lists was tested by
1000 Fisher's exact test, assuming a background of 19,396 coding genes in the genome
1001 which is the total number profiled by Regeneron¹².

1002 **The PULSE model**

1003 Feature engineering. Given the variant annotations from ANNOVAR, we calculated
1004 gene-level mutation profiles for each individual. Here we only focused on rare
1005 nonsynonymous and splicing-site SNVs as well as frameshift and splicing-site indels.
1006 Rare variants were defined as those not present within 1000 Genomes Project (Phase
1007 3) samples. For nonsynonymous and splicing-site SNVs, we calculated the
1008 accumulative mutation burdens for each gene based on individual annotations (32 in
1009 total; **Supplementary Table 7**). For indels, the number of variants was counted for
1010 frameshift and splicing-site, respectively (**Supplementary Table 7**). Consequently, the
1011 mutation profile consists of 34 features per gene per individual.

1012 Mapping phenotype from genotype. Given the mutation profiles $\mathbf{X}_i \in \mathbb{R}^{K \times M}$ ($i = 1, \dots, N$)
1013 for the i -th sample and the corresponding disease status $y_i \in \{0, 1\}$ ($y_i = 1$ indicates a
1014 case, and $y_i = 0$ otherwise), PULSE models the conditional $P(y_i | \mathbf{X}_i)$, which is the
1015 probability of disease status for the i -th sample characterized by the genome. Here K ,
1016 M , and N are the numbers of annotation features, genes, and samples, respectively.

1017 Note that we have $K=34$ in this study. In particular, we aggregate the mutation profiles
1018 across the genome using a bilinear transformation and define the conditional as

1019
$$p(y_i | \mathbf{X}_i, \mathbf{w}_1, \mathbf{w}_2) = \text{Bern}(y_i; \sigma(\mathbf{w}_1^T \mathbf{X}_i \mathbf{w}_2)), \quad (19)$$

1020 where $\sigma(\cdot)$ denotes the sigmoid function, \mathbf{w}_1 are random variables weighing the
1021 importance of each annotation feature, and \mathbf{w}_2 effect sizes for individual genes. We
1022 model \mathbf{w}_1 by a multivariate Gaussian given by

1023
$$p(\mathbf{w}_1 | \Lambda) = \mathcal{N}(\mathbf{w}_1; \mathbf{0}, \Lambda^{-1}), \quad (20)$$

1024 in which the precision matrix Λ is characterized by a Wishart distribution, i.e.,

1025
$$p(\Lambda) = \mathcal{W}(\Lambda; \mathbf{W}_0, \nu_0), \quad (21)$$

1026 and the hyperparameters are set to $\mathbf{W}_0 = \mathbf{I}_K$ and $\nu_0 = K$ to introduce non-informative
1027 prior.

1028 To prevent overfitting, we introduce a spike-and-slab prior over \mathbf{w}_2 , i.e.,

1029
$$p(\mathbf{w}_{2j} | \pi, \lambda) = \pi \mathcal{N}(\mathbf{w}_{2j}; 0, \lambda^{-1}) + (1 - \pi) \delta_0(\mathbf{w}_{2j}), \quad (22)$$

1030 where π is the probability of being non-zero and $\delta_0(\cdot)$ is the Dirac function forcing \mathbf{w}_{2j}
1031 to be zero. Two additional conjugate priors are further used over distribution parameters
1032 in (22), i.e.,

1033
$$p(\pi) = \text{Beta}(\pi; \alpha_0, \beta_0), \quad (23)$$

1034 and

1035
$$p(\lambda) = \text{Gamma}(\lambda; a_0, b_0), \quad (24)$$

1036 in which we set $\alpha_0 = \beta_0 = 0.5$ (i.e., the Jeffrey prior) and $a_0 = b_0 = 10^{-6}$ to keep it
1037 non-informative. In this study, to prevent false positives, accelerate computation, and
1038 eliminate the sex bias in the genetic modelling, we only considered autosomal genes
1039 that are expressed in human lungs (TPM>1 in lung RNA-seq from GTEx¹⁰⁴), resulting in
1040 $M=13,129$. The diagram of the model structure is shown in Panel 3 of **Fig. 4b**.

1041 Model inference. The exact inference in PULSE is intractable. Here we adopt the
1042 mean-field variational inference (MFVI), an approximate but efficient way to perform
1043 inference in Bayesian models⁷³. Since the model posterior is difficult to calculate, MFVI
1044 aims to search for an optimal distribution closest to the model posterior from a family of

1045 regularized proposal distributions factorized with each other. Indeed, the solution of
1046 MFVI is given by minimizing the Kullback-Leibler (KL) divergence, i.e.,

$$1047 \quad q^*(\Phi) = \underset{q(\Phi) \in \mathcal{Q}}{\operatorname{argmin}} \operatorname{KL}(q(\Phi) \parallel p(\Phi | \mathbf{X}_{\text{train}}, \mathbf{Y}_{\text{train}})),$$

1048 where Φ represents the set of hidden variables in the model, and \mathcal{Q} is the family of
1049 factorized proposal distributions. It can be shown that minimizing the KL divergence is
1050 mathematically equivalent to maximizing the evidence lower bound (ELBO)⁷³, which is
1051 solvable in optimization. Further, in order to make the MFVI for PULSE tractable and
1052 efficient, several techniques were adopted.

1053 (i) *Local variational method*. The sigmoid function in Eq. (19) makes MFVI intractable.
1054 However, instead of dealing with the sigmoid directly, we can approximately calculate
1055 the posteriors of \mathbf{w}_1 and \mathbf{w}_2 with respect to its lower bound, which yields Gaussian or
1056 Gaussian-like distributions. Meanwhile, to make the approximation close to the true
1057 MFVI solution, we need to maximize the log-likelihood of observations that take the
1058 sigmoid lower bound into account with respect to local variational parameters
1059 introduced. This local variational method introduces a new objective function, which is
1060 consistent with the original MFVI. More technical details can be found in Supplementary
1061 Notes.

1062 (ii) *Reparameterization*. The spike-and-slab prior over \mathbf{w}_2 (Eq. (22)) also makes MFVI
1063 intractable. To solve this problem, we adopted the reparameterization trick introduced in
1064 ¹⁰⁵. In particular, \mathbf{w}_2 can be reparameterized by two other variables \mathbf{s} and $\bar{\mathbf{w}}_2$, whose
1065 joint distribution is given by

$$1066 \quad p(\mathbf{s}, \bar{\mathbf{w}}_2) = \mathcal{N}(\bar{\mathbf{w}}_2; \mathbf{0}, \lambda^{-1} \mathbf{I}) \prod_j \{ \pi^{s_j} (1 - \pi)^{1 - s_j} \}, \quad (25)$$

1067 and the new variable $s_j \bar{w}_{2j}$ follows the same distribution as in Eq. (22). Therefore, we
1068 can do MFVI over \mathbf{s} and $\bar{\mathbf{w}}_2$ instead of \mathbf{w}_2 . However, this still introduces another
1069 problem and makes the VI highly inefficient, where the approximate posteriors from
1070 reparameterization (unimodal) could badly deviate from the original posteriors
1071 (exponentially multimodal). To alleviate this issue, a partial factorization was taken by
1072 following¹⁰⁵, i.e., we assume

1073
$$q(\mathbf{s}, \bar{\mathbf{w}}_2) = \prod_j q(s_j, \bar{w}_{2j}), \quad (26)$$

1074 in proposal distributions, and performed MFVI over s_j and \bar{w}_{2j} jointly. More technical
1075 details can be found in Supplementary Notes.

1076 (iii) *Stochastic variational inference*. Conventional MFVI based on coordinate ascent
1077 (i.e., CAVI) updates variational parameters in batches. However, it is difficult to deploy
1078 such a batch algorithm in big data scenarios, where the sample size or feature
1079 dimension is large. Here, stochastic variational inference (SVI)¹⁰⁶ was used to scale up
1080 our model for the large amount of genome data. In fact, borrowing the idea from
1081 stochastic optimization, we can update parameters per epoch by using only one or a
1082 mini-batch of samples instead of the whole dataset. Specifically, with SVI we first
1083 calculated the natural gradient of ELBO with respect to the variational parameter whose
1084 update rule contains sample points. Thanks to the conditional conjugacy predefined in
1085 our model, the natural gradient enjoys a simple form (see Supplementary Notes for
1086 details). Then based on stochastic optimization, we sampled a minibatch and rescaled
1087 the term involving sample points, resulting in a noisy but cheaply computed and
1088 unbiased natural gradient. At last, the variational parameter was updated from this
1089 gradient according to the gradient-based optimization algorithm¹⁰⁷. This SVI update can
1090 be easily embedded into CAVI without many changes. In implementation, we followed
1091 ¹⁰⁶ and set the learning rate as

1092
$$\epsilon_t = (t + \tau)^{-\kappa}, \quad \tau = 1, \kappa = 0.9, \quad (27)$$

1093 where t is the iteration index, τ is the delay, and κ is the forgetting rate.

1094 We integrated all above techniques into our VI algorithm. Details on the update rules for
1095 both local and global variational parameters and the VI algorithm are provided in
1096 Supplementary Notes.

1097 MAP prediction. The exact Bayesian prediction for test samples needs to integrate out
1098 all hidden variables, which is computationally intense and usually not necessary. Here,
1099 we adopted maximum a posteriori (MAP) and predicted new coming sample by

1100
$$p(y_{\text{new}} | \mathbf{X}_{\text{new}}, \mathbf{X}_{\text{train}}, y_{\text{train}}) \approx p(y_{\text{new}} | \mathbf{X}_{\text{new}}, \boldsymbol{\Theta}^*), \quad (28)$$

1101 where the optimal hidden variables are given by

$$\begin{aligned} \Theta^* &= \operatorname{argmax} q(\Theta) \\ &\approx \operatorname{argmax} p(\Theta | \mathbf{X}_{\text{train}}, \mathbf{Y}_{\text{train}}). \end{aligned} \quad (29)$$

1103 Similarly, the importance weights for individual genes (referred to as PULSE gene
1104 weights) were also estimated by the MAP of $q(w_{2j})$.

1105 **Network analysis**

1106 We first downloaded the human PPIs from STRING v11, including 19,567 proteins and
1107 11,759,455 protein interactions. To eliminate the bias caused by hub proteins, we first
1108 carried out the random walk with restart algorithm¹⁰⁸ over the PPI network, wherein the
1109 restart probability was set to 0.5, resulting in a smoothed network after retaining the top
1110 5% predicted edges. To decompose the network into different subnetworks/modules, we
1111 performed the Leiden algorithm⁶², a community detection algorithm that searches for
1112 densely connected modules by optimizing the modularity. After the algorithm converged,
1113 we obtained 1,681 modules with an average size of 9.98 nodes (SD=53.35;
1114 **Supplementary Table 13**).

1115 **Transcriptome analysis**

1116 Four single-cell RNA-seq datasets were used in the transcriptome analyses, including
1117 human healthy lungs^{33,39} and COVID-19 patients^{22,23}. Data after QC was acquired for
1118 each study. Only samples from the respiratory system were considered in the analyses.
1119 For the healthy lung data, a cutoff of 0.6931 was used to define expressed genes in the
1120 transcriptome throughout the study if not specified. For the disease samples, we
1121 removed the overlap of severe patients between the two cohorts^{22,23}. In the comparative
1122 expression analysis of severe versus moderate patients, to stabilize the analysis we
1123 estimated the change of gene expression levels using the Z-score estimated from
1124 Wilcoxon rank-sum test, wherein a positive Z-score indicates a higher expression level
1125 in severe patients and a negative value suggests the lower expression. The
1126 Benjamini-Hochberg (BH) procedure was used for multiple testing correction throughout
1127 the study.

1128 **ACKNOWLEDGEMENTS**

1129 We acknowledge the Stanford Genetics Bioinformatics Service Center (GBSC) for
1130 providing computational infrastructure for this study. This study was also supported by
1131 the National Institutes of Health (1S10OD023452-01 to GBSC;
1132 CECS 5P50HG00773504, 1P50HL083800, 1R01HL101388, 1R01-HL122939,
1133 S10OD025212, P30DK116074, and UM1HG009442 to M.P.S.), the Wellcome Trust
1134 (216596/Z/19/Z to J.C.K.), and Million Veteran Program, Office of Research and
1135 Development, Veterans Health Administration (MVP001). This publication does not
1136 represent the views of the Department of Veteran Affairs or the United States
1137 Government. Figures 1a and 4b were created with BioRender.com.

1138 This study is part of the GEN-COVID Multicenter Study
1139 (<https://sites.google.com/dbm.unisi.it/gen-covid>), the Italian multicenter study aimed at
1140 identifying the COVID-19 host genetic bases. Specimens were provided by the
1141 COVID-19 Biobank of Siena, which is part of the Genetic Biobank of Siena, member of
1142 BBMRI-IT, of Telethon Network of Genetic Biobanks (project no. GTB18001), of
1143 EuroBioBank, and of RDConnect. We thank the CINECA consortium for providing
1144 computational resources and the Network for Italian Genomes (NIG)
1145 (<http://www.nig.cineca.it>) for its support. We thank private donors for the support
1146 provided to A.R. (Department of Medical Biotechnologies, University of Siena) for the
1147 COVID-19 host genetics research project (D.L n.18 of March 17, 2020). We also thank
1148 the COVID-19 Host Genetics Initiative ([https:// www.covid19hg.org/](https://www.covid19hg.org/)), MIUR project
1149 “Dipartimenti di Eccellenza 2018-2020” to the Department of Medical Biotechnologies
1150 University of Siena, Italy, and “Bando Ricerca COVID-19 Toscana” project to Azienda
1151 Ospedaliero Universitaria Senese. We also thank Intesa San Paolo for the 2020 charity
1152 fund dedicated to the project “N. B/2020/0119 Identificazione delle basi genetiche
1153 determinanti la variabilità clinica della risposta a COVID-19 nella popolazione italiana”
1154 and “Bando FISR 2020” in COVID-19 from Italian Ministry of University e Research.

1155 **AUTHOR CONTRIBUTIONS**

1156 S.Z., J.C.K. and M.P.S. conceived and designed the study. S.Z. contributed to the
1157 design, implementation, training and testing of RefMap and PULSE. S.Z., J.C.K.,
1158 A.K.W., C.H., T.H.J., S.F., E.F., F.F., A.R., C.P., J.S., P.B.R., P.S.T. and M.P.S. were
1159 responsible for data acquisition. S.Z., J.C.K., C.H., T.H.J., C.W., J.L. and C.P. were
1160 responsible for data analysis. S.Z., J.C.K., A.K.W., C.H., T.H.J., C.W., J.L., S.F., E.F.,
1161 F.F., A.R., C.P., P.G., X.S., I.S.T., K.P.K., M.M.D., P.S.T. and M.P.S. were responsible for
1162 the interpretation of the findings. S.Z., J.C.K., P.S.T. and M.P.S. drafted the manuscript
1163 with assistance from all authors. All authors meet the four ICMJE authorship criteria,
1164 and were responsible for revising the manuscript, approving the final version for
1165 publication, and for accuracy and integrity of the work.

1166 **COMPETING INTERESTS**

1167 M.P.S. is a cofounder of Personalis, Qbio, Sensomics, Filtricine, Mirvie, and January. He
1168 is on the scientific advisory of these companies and Genapsys. J.L. is a cofounder of
1169 Sensomics. No other authors have competing interests.

1170 **SUPPLEMENTARY INFORMATION**

1171 **SUPPLEMENTARY FIGURES**

1172 **Supplementary Figure 1**

1173 Mendelian randomization for COVID-19 GWAS with phenotypes B2 and C2

1174 **Supplementary Figure 2**

1175 Expression levels of non-developmental genes in healthy lungs. nDG is short for
1176 non-developmental gene

1177 **Supplementary Figure 3**

1178 Overlap between lung snATAC-seq peaks and ENCODE CHIP-seq peaks

1179 **Supplementary Figure 4**

1180 Q-Q plot of P -value distribution for SKAT analysis on the GEN-COVID cohort

1181 **Supplementary Figure 5**

1182 Age distribution of the GEN-COVID cohort after sample filtering

1183 **Supplementary Figure 6**

1184 Coefficients of the age+sex logistic regression models in 5-fold cross-validation

1185 **Supplementary Figure 7**

1186 Age distribution of the VA COVID-19 cohort after sample filtering

1187 **Supplementary Figure 8**

1188 Sex distributions of the GEN-COVID and VA cohorts after sample filtering

1189 **Supplementary Figure 9**

1190 Normalized rank versus specificity of PULSE prediction for the VA COVID-19 cohort

1191 **Supplementary Figure 10**

1192 Normalized rank versus specificity and sensitivity of PULSE prediction for GEN-COVID

1193 non-EUR samples

1194 **Supplementary Figure 11**

1195 Distribution of gene weights of the PULSE model trained on GEN-COVID EUR samples

1196 **SUPPLEMENTARY TABLES**

1197 **Supplementary Table 1**

1198 RefMap COVID-19 regions and genes

1199 **Supplementary Table 2**

1200 Enrichment of disease-associated SNPs in RefMap regions based on the 23andMe

1201 study

1202 **Supplementary Table 3**

1203 Partitioned heritability analysis by LDSC for COVID-19 GWAS phenotypes A2, B2, and
1204 C2

1205 **Supplementary Table 4**

1206 GO enrichment for RefMap genes per cell type

1207 **Supplementary Table 5**

1208 Pathway enrichment for RefMap genes per cell type

1209 **Supplementary Table 6**

1210 Accession identifiers for ENCODE samples

1211 **Supplementary Table 7**

1212 Variant annotations and their weights learned by PULSE

1213 **Supplementary Table 8**

1214 Clinical characteristics and QC results of 1,339 samples in the GEN-COVID cohort

1215 **Supplementary Table 9**

1216 Clinical characteristics and QC results of 590 samples in the VA COVID-19 cohort

1217 **Supplementary Table 10**

1218 Genes with top 5% weights in the PULSE model

1219 **Supplementary Table 11**

1220 Genes predicted by either RefMap or PULSE to be associated with NK cells

1221 **Supplementary Table 12**

1222 PPI network after network smoothing. Gene identifiers were given in Supplementary
1223 Table 13.

1224 **Supplementary Table 13**

1225 Modules detected by the Leiden algorithm and modules significantly enriched with
1226 NK-cell COVID-19 genes

1227 **Supplementary Table 14**

1228 GO enrichment for modules enriched with NK-cell COVID-19 genes

1229 **Supplementary Table 15**

1230 Pathway enrichment for modules enriched with NK-cell COVID-19 genes

1231 **Supplementary Notes**

1232 Technical details on the PULSE model

1233 **Data availability**

1234 The GEN-COVID WES and clinical data are available by consultation (A.R.). The VA
1235 WGS and clinical data are available upon request from the corresponding authors
1236 (P.S.T. and M.P.S.); these data are not publicly available due to US Government and
1237 Department of Veteran's Affairs restrictions relating to participant privacy and consent.
1238 All other data used in this study are available from the original studies.

1239 **Code availability**

1240 The computer codes generated in this study are available from the authors upon
1241 request (P.S.T. and M.P.S.).

1242

1243 REFERENCES

- 1244 1. Dong, E., Du, H. & Gardner, L. An interactive web-based dashboard to track
1245 COVID-19 in real time. *Lancet Infect. Dis.* **20**, 533–534 (2020).
- 1246 2. Shang, Y. *et al.* Scoring systems for predicting mortality for severe patients with
1247 COVID-19. *EClinicalMedicine* **24**, 100426 (2020).
- 1248 3. Li, X. *et al.* Predictive indicators of severe COVID-19 independent of comorbidities
1249 and advanced age: a nested case- control study. *Epidemiology & Infection* **148**,
1250 (2020).
- 1251 4. Initiative, T. C.-19 H. G. & The COVID-19 Host Genetics Initiative. The COVID-19
1252 Host Genetics Initiative, a global initiative to elucidate the role of host genetic
1253 factors in susceptibility and severity of the SARS-CoV-2 virus pandemic. *European*
1254 *Journal of Human Genetics* vol. 28 715–718 (2020).
- 1255 5. Shelton, J. F. *et al.* Trans-ancestry analysis reveals genetic and nongenetic
1256 associations with COVID-19 susceptibility and severity. *Nature Genetics* (2021)
1257 doi:10.1038/s41588-021-00854-7.
- 1258 6. Genomewide Association Study of Severe Covid-19 with Respiratory Failure. *N.*
1259 *Engl. J. Med.* **383**, 1522–1534 (2020).
- 1260 7. Initiative, C.-19 H. G. & Others. Mapping the human genetic architecture of
1261 COVID-19 by worldwide meta-analysis. *MedRxiv* (2021).
- 1262 8. Pairo-Castineira, E. *et al.* Genetic mechanisms of critical illness in COVID-19.
1263 *Nature* **591**, 92–98 (2021).
- 1264 9. Wang, F. *et al.* Initial whole-genome sequencing and analysis of the host genetic
1265 contribution to COVID-19 severity and susceptibility. *Cell Discovery* vol. 6 (2020).

- 1266 10. Benetti, E. *et al.* Clinical and molecular characterization of COVID-19 hospitalized
1267 patients. *PLoS One* **15**, e0242534 (2020).
- 1268 11. Novelli, A. *et al.* Analysis of ACE2 genetic variants by direct exome sequencing in
1269 99 SARS-CoV-2 positive patients. (2020).
- 1270 12. Kosmicki, J. A. *et al.* A catalog of associations between rare coding variants and
1271 COVID-19 outcomes. *medRxiv* (2021) doi:10.1101/2020.10.28.20221804.
- 1272 13. Huang, C. *et al.* Clinical features of patients infected with 2019 novel coronavirus in
1273 Wuhan, China. *Lancet* **395**, 497–506 (2020).
- 1274 14. Brodin, P. Immune determinants of COVID-19 disease presentation and severity.
1275 *Nat. Med.* **27**, 28–33 (2021).
- 1276 15. Mehta, P. *et al.* COVID-19: consider cytokine storm syndromes and
1277 immunosuppression. *Lancet* **395**, 1033–1034 (2020).
- 1278 16. Mathew, D. *et al.* Deep immune profiling of COVID-19 patients reveals distinct
1279 immunotypes with therapeutic implications. *Science* **369**, (2020).
- 1280 17. Sosa-Hernández, V. A. *et al.* B Cell Subsets as Severity-Associated Signatures in
1281 COVID-19 Patients. *Front. Immunol.* **11**, 611004 (2020).
- 1282 18. Lucas, C. *et al.* Longitudinal analyses reveal immunological misfiring in severe
1283 COVID-19. *Nature* **584**, 463–469 (2020).
- 1284 19. Arunachalam, P. S. *et al.* Systems biological assessment of immunity to mild versus
1285 severe COVID-19 infection in humans. *Science* **369**, 1210–1220 (2020).
- 1286 20. Zhang, J.-Y. *et al.* Single-cell landscape of immunological responses in patients
1287 with COVID-19. *Nat. Immunol.* **21**, 1107–1118 (2020).
- 1288 21. Stephenson, E. *et al.* The cellular immune response to COVID-19 deciphered by

- 1289 single cell multi-omics across three UK centres. *medRxiv* (2021).
- 1290 22. Liao, M. *et al.* Single-cell landscape of bronchoalveolar immune cells in patients
1291 with COVID-19. *Nat. Med.* **26**, 842–844 (2020).
- 1292 23. Ren, X. *et al.* COVID-19 immune features revealed by a large-scale single-cell
1293 transcriptome atlas. *Cell* **184**, 1895–1913.e19 (2021).
- 1294 24. Melms, J. C. *et al.* A molecular single-cell lung atlas of lethal COVID-19. *Nature*
1295 (2021) doi:10.1038/s41586-021-03569-1.
- 1296 25. Delorey, T. M. *et al.* COVID-19 tissue atlases reveal SARS-CoV-2 pathology and
1297 cellular targets. *Nature* 1–8 (2021).
- 1298 26. Miorin, L. *et al.* SARS-CoV-2 Orf6 hijacks Nup98 to block STAT nuclear import and
1299 antagonize interferon signaling. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 28344–28354
1300 (2020).
- 1301 27. Blanco-Melo, D. *et al.* Imbalanced Host Response to SARS-CoV-2 Drives
1302 Development of COVID-19. *Cell* **181**, 1036–1045.e9 (2020).
- 1303 28. Vietzen, H. *et al.* Deletion of the NKG2C receptor encoding KLRC2 gene and
1304 HLA-E variants are risk factors for severe COVID-19. *Genet. Med.* (2021)
1305 doi:10.1038/s41436-020-01077-7.
- 1306 29. Wang, E. Y. *et al.* Diverse Functional Autoantibodies in Patients with COVID-19.
1307 *Nature* (2021) doi:10.1038/s41586-021-03631-y.
- 1308 30. Maucourant, C. *et al.* Natural killer cell immunotypes related to COVID-19 disease
1309 severity. *Sci Immunol* **5**, (2020).
- 1310 31. Azzi, Y., Bartash, R., Scalea, J., Loarte-Campos, P. & Akalin, E. COVID-19 and
1311 Solid Organ Transplantation: A Review Article. *Transplantation* **105**, 37–55 (2021).

- 1312 32. Zhang, S., Cooper-Knock, J., Weimer, A. K., Shi, M. & Moll, T. Genome-wide
1313 Identification of the Genetic Basis of Amyotrophic Lateral Sclerosis. (2020).
- 1314 33. Wang, A. *et al.* Single-cell multiomic profiling of human lungs reveals
1315 cell-type-specific and age-dynamic control of SARS-CoV2 host genes. *Elife* **9**,
1316 (2020).
- 1317 34. Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from
1318 polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).
- 1319 35. Delorey, T. M. *et al.* A single-cell and spatial atlas of autopsy tissues reveals
1320 pathology and cellular targets of SARS-CoV-2. *bioRxiv* (2021)
1321 doi:10.1101/2021.02.25.430130.
- 1322 36. Smith, G. D. Mendelian Randomization for Strengthening Causal Inference in
1323 Observational Studies. *Perspectives on Psychological Science* vol. 5 527–545
1324 (2010).
- 1325 37. Roederer, M. *et al.* The genetic architecture of the human immune system: a
1326 bioresource for autoimmunity and disease pathogenesis. *Cell* **161**, 387–403 (2015).
- 1327 38. Raulet, D. H. Roles of the NKG2D immunoreceptor and its ligands. *Nat. Rev.*
1328 *Immunol.* **3**, 781–790 (2003).
- 1329 39. Travaglini, K. J. *et al.* A molecular cell atlas of the human lung from single-cell RNA
1330 sequencing. *Nature* **587**, 619–625 (2020).
- 1331 40. Kuleshov, M. V. *et al.* Enrichr: a comprehensive gene set enrichment analysis web
1332 server 2016 update. *Nucleic Acids Res.* **44**, W90–7 (2016).
- 1333 41. Balboa, M. A., Balsinde, J., Aramburu, J., Mollinedo, F. & López-Botet, M.
1334 Phospholipase D activation in human natural killer cells through the Kp43 and

- 1335 CD16 surface antigens takes place by different mechanisms. Involvement of the
1336 phospholipase D pathway in tumor necrosis factor alpha synthesis. *J. Exp. Med.*
1337 **176**, 9–17 (1992).
- 1338 42. Watzl, C. & Long, E. O. Signal transduction during activation and inhibition of
1339 natural killer cells. *Curr. Protoc. Immunol.* **Chapter 11**, Unit 11.9B (2010).
- 1340 43. Mikulak, J., Oriolo, F., Zaghi, E., Di Vito, C. & Mavilio, D. Natural killer cells in HIV-1
1341 infection and therapy. *AIDS* **31**, 2317–2330 (2017).
- 1342 44. Nuvor, S. V., van der Sande, M., Rowland-Jones, S., Whittle, H. & Jaye, A. Natural
1343 Killer Cell Function Is Well Preserved in Asymptomatic Human Immunodeficiency
1344 Virus Type 2 (HIV-2) Infection but Similar to That of HIV-1 Infection When CD4
1345 T-Cell Counts Fall. *Journal of Virology* vol. 80 2529–2538 (2006).
- 1346 45. Hoffmann, M. *et al.* SARS-CoV-2 Cell Entry Depends on ACE2 and TMPRSS2 and
1347 Is Blocked by a Clinically Proven Protease Inhibitor. *Cell* **181**, 271–280.e8 (2020).
- 1348 46. De Biasi, S. *et al.* Marked T cell activation, senescence, exhaustion and skewing
1349 towards TH17 in patients with COVID-19 pneumonia. *Nat. Commun.* **11**, 3434
1350 (2020).
- 1351 47. He, L. *et al.* Pericyte-specific vascular expression of SARS-CoV-2 receptor ACE2 –
1352 implications for microvascular inflammation and hypercoagulopathy in COVID-19.
1353 doi:10.1101/2020.05.11.088500.
- 1354 48. ENCODE Project Consortium *et al.* Expanded encyclopaedias of DNA elements in
1355 the human and mouse genomes. *Nature* **583**, 699–710 (2020).
- 1356 49. Gel, B. *et al.* regioneR: an R/Bioconductor package for the association analysis of
1357 genomic regions based on permutation tests. *Bioinformatics* **btv562** (2015)

1358 doi:10.1093/bioinformatics/btv562.

1359 50. Lee, S. *et al.* Optimal unified approach for rare-variant association testing with
1360 application to small-sample case-control whole-exome sequencing studies. *Am. J.*
1361 *Hum. Genet.* **91**, 224–237 (2012).

1362 51. Benetti, E. *et al.* ACE2 gene variants may underlie interindividual variability and
1363 susceptibility to COVID-19 in the Italian population. *Eur. J. Hum. Genet.* **28**,
1364 1602–1614 (2020).

1365 52. Daga, S. *et al.* Employing a systematic approach to biobanking and analyzing
1366 clinical and genetic data for advancing COVID-19 research. *Eur. J. Hum. Genet.*
1367 (2021) doi:10.1038/s41431-020-00793-7.

1368 53. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic
1369 variants from high-throughput sequencing data. *Nucleic Acids Res.* **38**, e164
1370 (2010).

1371 54. 1000 Genomes Project Consortium *et al.* A global reference for human genetic
1372 variation. *Nature* **526**, 68–74 (2015).

1373 55. Aschard, H. *et al.* Combining effects from rare and common genetic variants in an
1374 exome-wide association study of sequence data. *BMC Proc.* **5 Suppl 9**, S44
1375 (2011).

1376 56. Pritchard, J. K. Are rare variants responsible for susceptibility to complex diseases?
1377 *Am. J. Hum. Genet.* **69**, 124–137 (2001).

1378 57. Pritchard, J. K. & Cox, N. J. The allelic architecture of human disease genes:
1379 common disease–common variant... or not? *Hum. Mol. Genet.* **11**, 2417–2423
1380 (2002).

- 1381 58. Li, J. *et al.* Decoding the Genomics of Abdominal Aortic Aneurysm. *Cell* **174**,
1382 1361–1372.e10 (2018).
- 1383 59. Li, J., Li, X., Zhang, S. & Snyder, M. Gene-Environment Interaction in the Era of
1384 Precision Medicine. *Cell* **177**, 38–44 (2019).
- 1385 60. Szklarczyk, D. *et al.* STRING v11: protein-protein association networks with
1386 increased coverage, supporting functional discovery in genome-wide experimental
1387 datasets. *Nucleic Acids Res.* **47**, D607–D613 (2019).
- 1388 61. Krishnan, A. *et al.* Genome-wide prediction and functional characterization of the
1389 genetic basis of autism spectrum disorder. *Nat. Neurosci.* **19**, 1454–1462 (2016).
- 1390 62. Traag, V. A., Waltman, L. & van Eck, N. J. From Louvain to Leiden: guaranteeing
1391 well-connected communities. *Sci. Rep.* **9**, 5233 (2019).
- 1392 63. Shilo, S., Rossman, H. & Segal, E. Signals of hope: gauging the impact of a rapid
1393 national vaccination campaign. *Nat. Rev. Immunol.* **21**, 198–199 (2021).
- 1394 64. Darby, A. C. & Hiscox, J. A. Covid-19: variants and vaccination. *BMJ* vol. 372 n771
1395 (2021).
- 1396 65. Petrilli, C. M. *et al.* Factors associated with hospital admission and critical illness
1397 among 5279 people with coronavirus disease 2019 in New York City: prospective
1398 cohort study. *BMJ* **369**, m1966 (2020).
- 1399 66. Rölle, A. *et al.* IL-12–producing monocytes and HLA-E control HCMV-driven
1400 NKG2C+ NK cell expansion. *J. Clin. Invest.* **124**, 5305–5316 (2014).
- 1401 67. Medzhitov, R. & Janeway, C. A., Jr. Decoding the patterns of self and nonself by the
1402 innate immune system. *Science* **296**, 298–300 (2002).
- 1403 68. Hu, W., Wang, G., Huang, D., Sui, M. & Xu, Y. Cancer Immunotherapy Based on

- 1404 Natural Killer Cells: Current Progress and New Opportunities. *Front. Immunol.* **10**,
1405 1205 (2019).
- 1406 69. Chua, R. L. *et al.* COVID-19 severity correlates with airway epithelium–immune cell
1407 interactions identified by single-cell analysis. *Nat. Biotechnol.* **38**, 970–979 (2020).
- 1408 70. Siedner, M. J., Tumarkin, E. & Bogoch, I. I. HIV post-exposure prophylaxis (PEP).
1409 *BMJ* k4928 (2018) doi:10.1136/bmj.k4928.
- 1410 71. Cheng, S. H. & Higham, N. J. A Modified Cholesky Algorithm Based on a
1411 Symmetric Indefinite Factorization. *SIAM Journal on Matrix Analysis and*
1412 *Applications* vol. 19 1097–1110 (1998).
- 1413 72. Harva, M. & Kabán, A. Variational learning for rectified factor analysis. *Signal*
1414 *Processing* vol. 87 509–527 (2007).
- 1415 73. Blei, D. M., Kucukelbir, A. & McAuliffe, J. D. Variational Inference: A Review for
1416 Statisticians. *Journal of the American Statistical Association* vol. 112 859–877
1417 (2017).
- 1418 74. Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* (2021)
1419 doi:10.1016/j.cell.2021.04.048.
- 1420 75. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing
1421 genomic features. *Bioinformatics* **26**, 841–842 (2010).
- 1422 76. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using
1423 genome-wide association summary statistics. *Nat. Genet.* **47**, 1228–1235 (2015).
- 1424 77. Sun, B. B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73–79
1425 (2018).
- 1426 78. Suhre, K. *et al.* Connecting genetic risk to disease end points through the human

- 1427 blood plasma proteome. *Nat. Commun.* **8**, 14357 (2017).
- 1428 79. Choi, K. W. *et al.* Assessment of Bidirectional Relationships Between Physical
1429 Activity and Depression Among Adults: A 2-Sample Mendelian Randomization
1430 Study. *JAMA Psychiatry* **76**, 399–408 (2019).
- 1431 80. Wootton, R. E. *et al.* Evaluation of the causal effects between subjective wellbeing
1432 and cardiometabolic health: mendelian randomisation study. *BMJ* **362**, k3788
1433 (2018).
- 1434 81. Julian, T. H. *et al.* Physical exercise is a risk factor for amyotrophic lateral sclerosis:
1435 Convergent evidence from mendelian randomisation, transcriptomics and risk
1436 genotypes. doi:10.1101/2020.11.24.20238063.
- 1437 82. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and
1438 population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 1439 83. Machiela, M. J. & Chanock, S. J. LDlink: a web-based application for exploring
1440 population-specific haplotype structure and linking correlated alleles of possible
1441 functional variants. *Bioinformatics* **31**, 3555–3557 (2015).
- 1442 84. Hartwig, F. P., Davies, N. M., Hemani, G. & Davey Smith, G. Two-sample
1443 Mendelian randomization: avoiding the downsides of a powerful, widely applicable
1444 but potentially fallible technique. *Int. J. Epidemiol.* **45**, 1717–1726 (2016).
- 1445 85. Burgess, S. & Thompson, S. G. Interpreting findings from Mendelian randomization
1446 using the MR-Egger method. *Eur. J. Epidemiol.* **32**, 377–389 (2017).
- 1447 86. Burgess, S. *et al.* Guidelines for performing Mendelian randomization
1448 investigations. *Wellcome Open Research* **4**, (2019).
- 1449 87. Burgess, S., Thompson, S. G. & CRP CHD Genetics Collaboration. Avoiding bias

- 1450 from weak instruments in Mendelian randomization studies. *Int. J. Epidemiol.* **40**,
1451 755–764 (2011).
- 1452 88. Bowden, J., Hemani, G. & Smith, G. D. Invited Commentary: Detecting Individual
1453 and Global Horizontal Pleiotropy in Mendelian Randomization—A Job for the
1454 Humble Heterogeneity Statistic? *American Journal of Epidemiology* (2018)
1455 doi:10.1093/aje/kwy185.
- 1456 89. Verbanck, M., Chen, C.-Y., Neale, B. & Do, R. Detection of widespread horizontal
1457 pleiotropy in causal relationships inferred from Mendelian randomization between
1458 complex traits and diseases. *Nat. Genet.* **50**, 693–698 (2018).
- 1459 90. Bowden, J. *et al.* Assessing the suitability of summary data for two-sample
1460 Mendelian randomization analyses using MR-Egger regression: the role of the I²
1461 statistic. *International Journal of Epidemiology* dyw220 (2016)
1462 doi:10.1093/ije/dyw220.
- 1463 91. de Leeuw, C. A., Mooij, J. M., Heskes, T. & Posthuma, D. MAGMA: generalized
1464 gene-set analysis of GWAS data. *PLoS Comput. Biol.* **11**, e1004219 (2015).
- 1465 92. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows–Wheeler
1466 transform. *Bioinformatics* **26**, 589–595 (2010).
- 1467 93. Poplin, R. *et al.* Scaling accurate genetic variant discovery to tens of thousands of
1468 samples. *bioRxiv* 201178 (2018) doi:10.1101/201178.
- 1469 94. Gaziano, J. M. *et al.* Million Veteran Program: A mega-biobank to study genetic
1470 influences on health and disease. *J. Clin. Epidemiol.* **70**, 214–223 (2016).
- 1471 95. Song, R. J. *et al.* Phenome-wide association of 1809 phenotypes and COVID-19
1472 disease progression in the Veterans Health Administration Million Veteran Program.

- 1473 *PLoS One* **16**, e0251651 (2021).
- 1474 96. Regier, A. A. *et al.* Functional equivalence of genome sequencing analysis
1475 pipelines enables harmonized variant calling across human genetics projects. *Nat.*
1476 *Commun.* **9**, 4038 (2018).
- 1477 97. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for
1478 genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
- 1479 98. Alexander, D. H. & Lange, K. Enhancements to the ADMIXTURE algorithm for
1480 individual ancestry estimation. *BMC Bioinformatics* **12**, 246 (2011).
- 1481 99. Liu, X., Wu, C., Li, C. & Boerwinkle, E. dbNSFP v3.0: A One-Stop Database of
1482 Functional Predictions and Annotations for Human Nonsynonymous and
1483 Splice-Site SNVs. *Human Mutation* vol. 37 235–241 (2016).
- 1484 100. Jian, X., Boerwinkle, E. & Liu, X. In silico prediction of splice-altering single
1485 nucleotide variants in the human genome. *Nucleic Acids Res.* **42**, 13534–13544
1486 (2014).
- 1487 101. Lin, H. *et al.* RegSNPs-intron: a computational framework for predicting pathogenic
1488 impact of intronic single nucleotide variants. *Genome Biol.* **20**, 254 (2019).
- 1489 102. Wu, B., Guan, W. & Pankow, J. S. On efficient and accurate calculation of
1490 significance P-values for sequence kernel association testing of variant set. *Ann.*
1491 *Hum. Genet.* **80**, 123–135 (2016).
- 1492 103. Mbatchou, J., Barnard, L., Backman, J. & Marcketta, A. Computationally efficient
1493 whole genome regression for quantitative and binary traits. *bioRxiv* (2020).
- 1494 104. Consortium, G. & GTEx Consortium. Genetic effects on gene expression across
1495 human tissues. *Nature* vol. 550 204–213 (2017).

- 1496 105. Titsias, M. K. & Lázaro-Gredilla, M. Spike and Slab Variational Inference for
1497 Multi-Task and Multiple Kernel Learning. in *Advances in Neural Information*
1498 *Processing Systems 24* (eds. Shawe-Taylor, J., Zemel, R. S., Bartlett, P. L., Pereira,
1499 F. & Weinberger, K. Q.) 2339–2347 (Curran Associates, Inc., 2011).
- 1500 106. Hoffman, M. D., Blei, D. M., Wang, C. & Paisley, J. Stochastic variational inference.
1501 (2013).
- 1502 107. Robbins, H. & Monro, S. A Stochastic Approximation Method. *Herbert Robbins*
1503 *Selected Papers* 102–109 (1985) doi:10.1007/978-1-4612-5110-1_9.
- 1504 108. Wang, S., Cho, H., Zhai, C., Berger, B. & Peng, J. Exploiting ontology graph for
1505 predicting sparsely annotated gene function. *Bioinformatics* **31**, i357–64 (2015).
- 1506

Figure 1

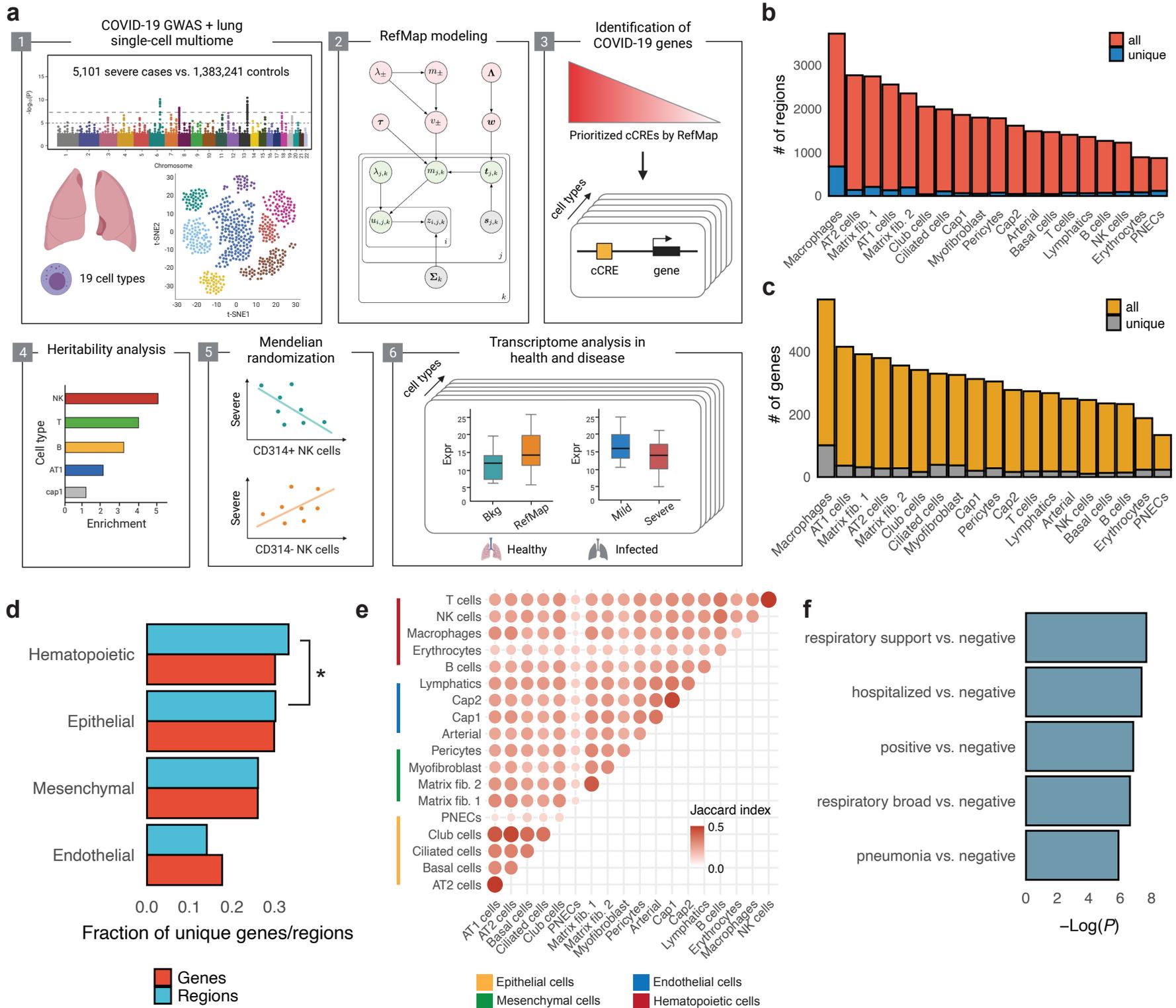
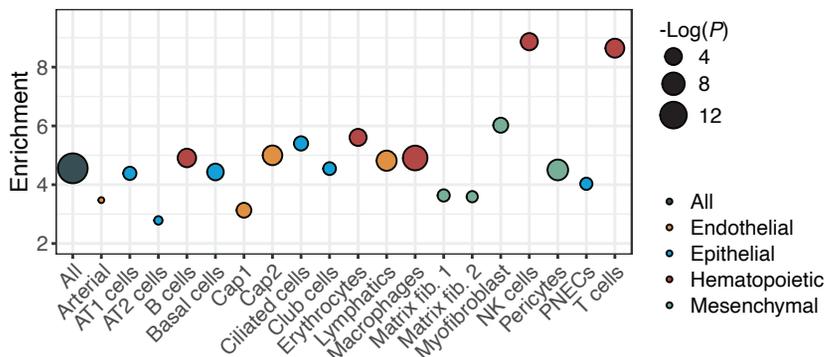
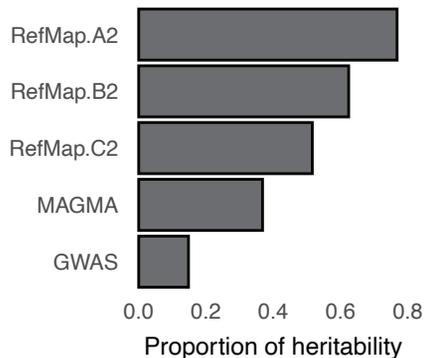


Figure 2

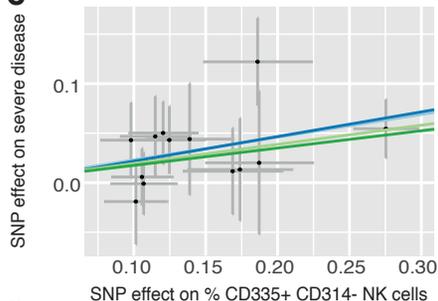
a



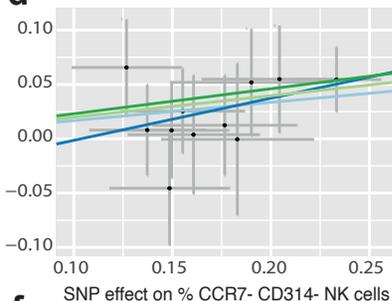
b



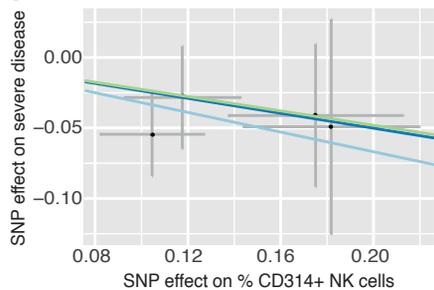
c



d



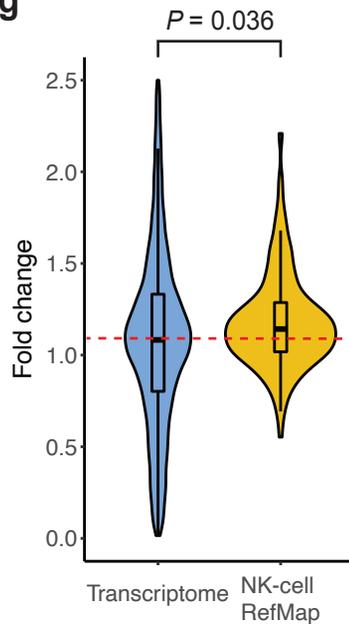
e



f

| Test | CD314+ | CCR7- CD314- | CD335+ CD314- |
|-----------------|----------|-----------------|------------------|
| IVW (mre) | 8.11E-06 | 5.71E-04 | 1.66E-05 |
| MR Egger | 9.21E-01 | 3.23E-01 | 2.10E-01 |
| Weighted median | 1.93E-01 | 4.87E-02 | 4.06E-02 |
| Weighted mode | 3.83E-01 | 6.59E-02 | 1.20E-01 |
| MR Lasso | 3.58E-02 | 2.51E-02 | 7.26E-04 |
| F<10 (n var) | 0 | 0 | 0 |
| Egger Q test | 0.86 | 0.92 | 0.77 |
| IVW Q test | 0.88 | 0.94 | 0.83 |
| MR PRESSO | 0.87 | 0.95 | 0.84 |
| Egger intercept | 0.61 | 0.56 | 0.91 |
| I2 | 0.96 | 0.97 | 0.97 |
| LOO (n var) | 0 | 0 | 0 |

g



▬ Inverse variance weighted (fixed effects) ▬ MR Egger
▬ Inverse variance weighted (multiplicative random effects) ▬ Weighted median

Figure 3

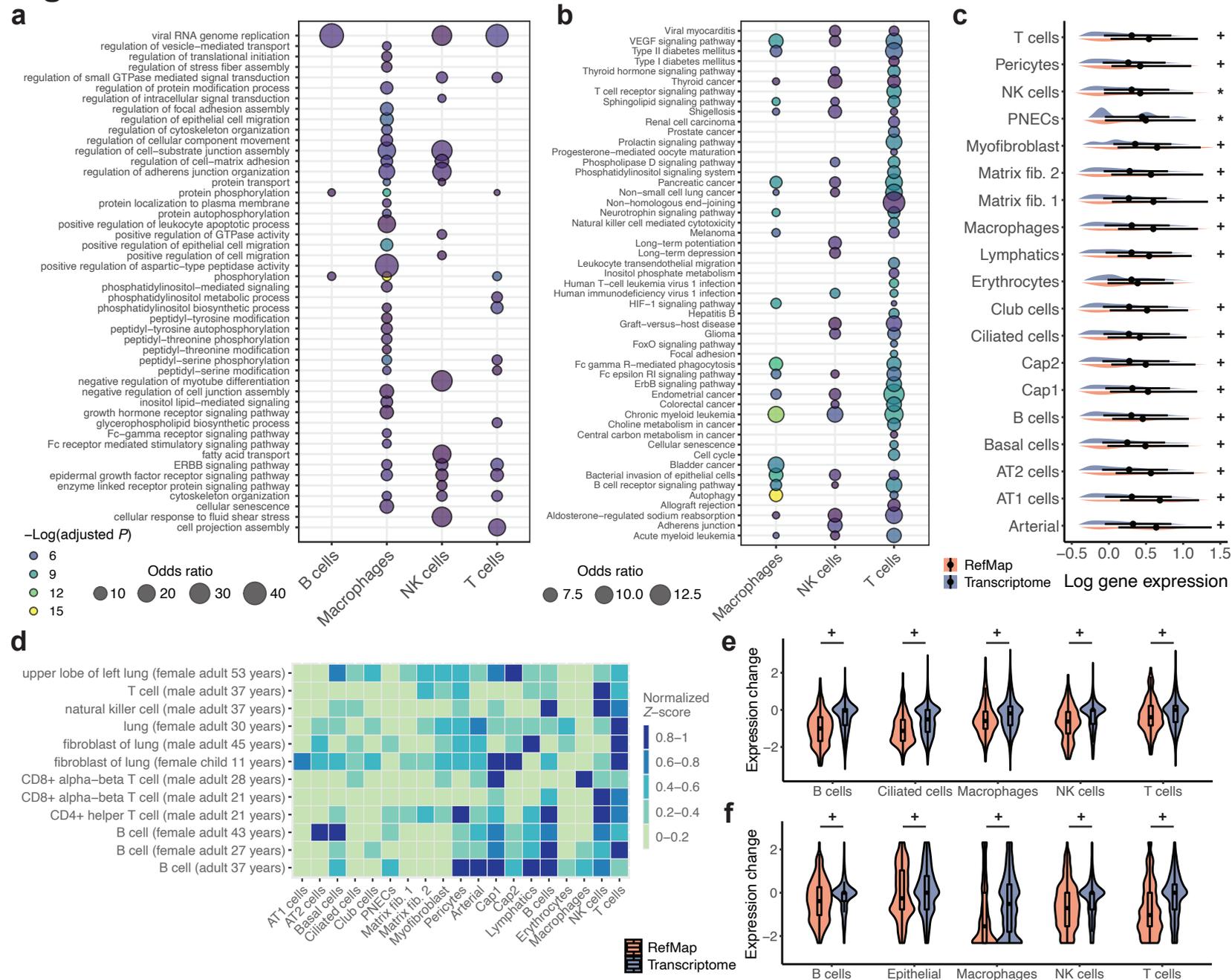


Figure 4

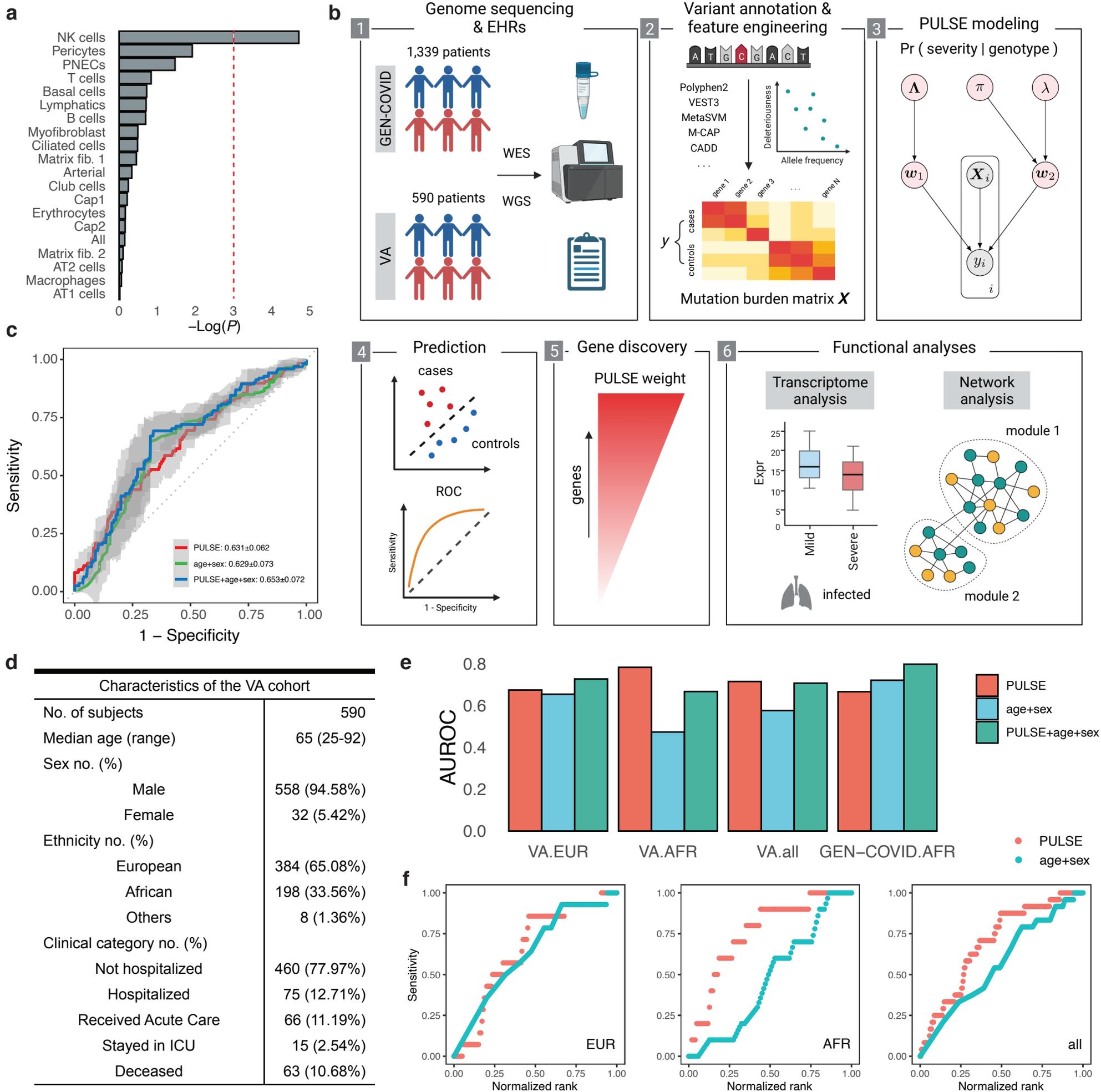
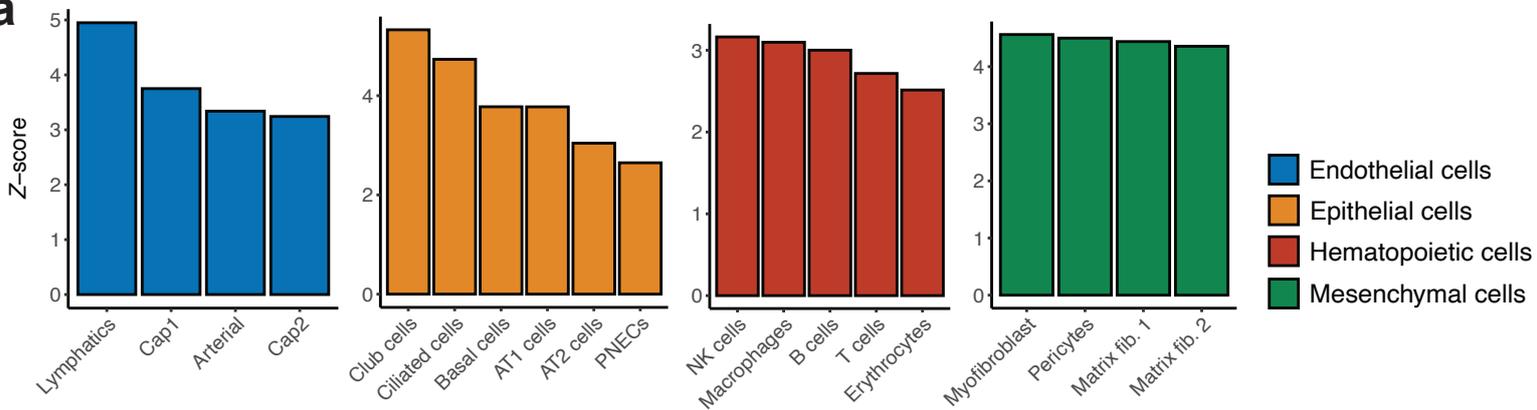
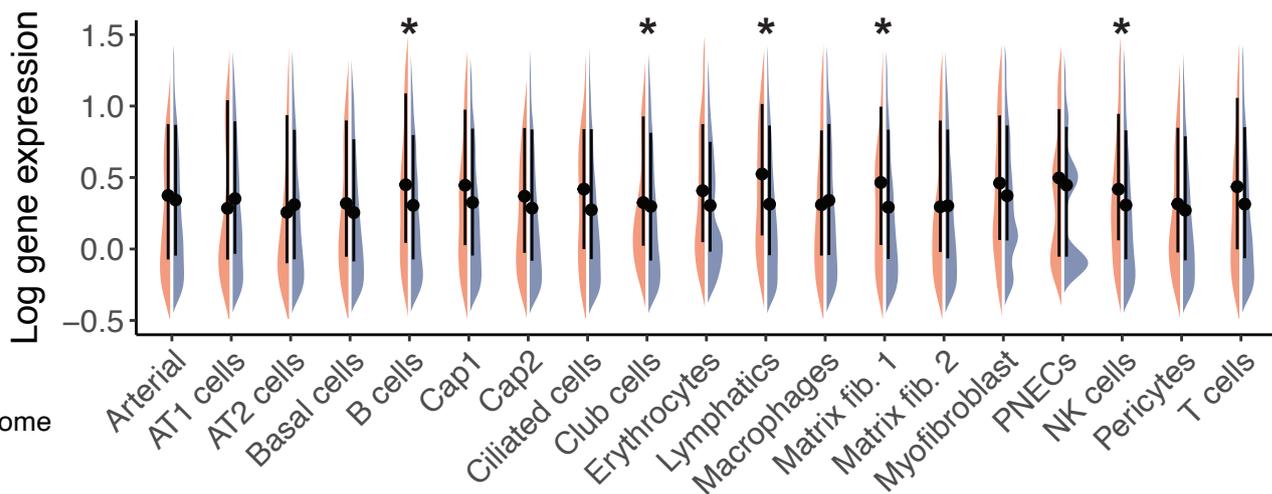


Figure 5

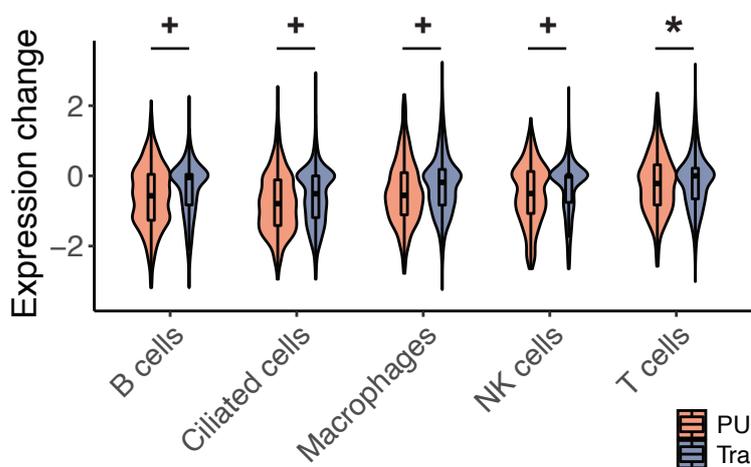
a



b



c



d

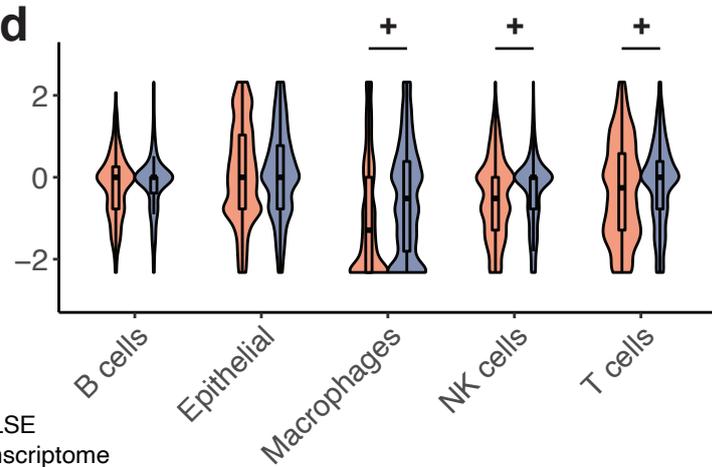
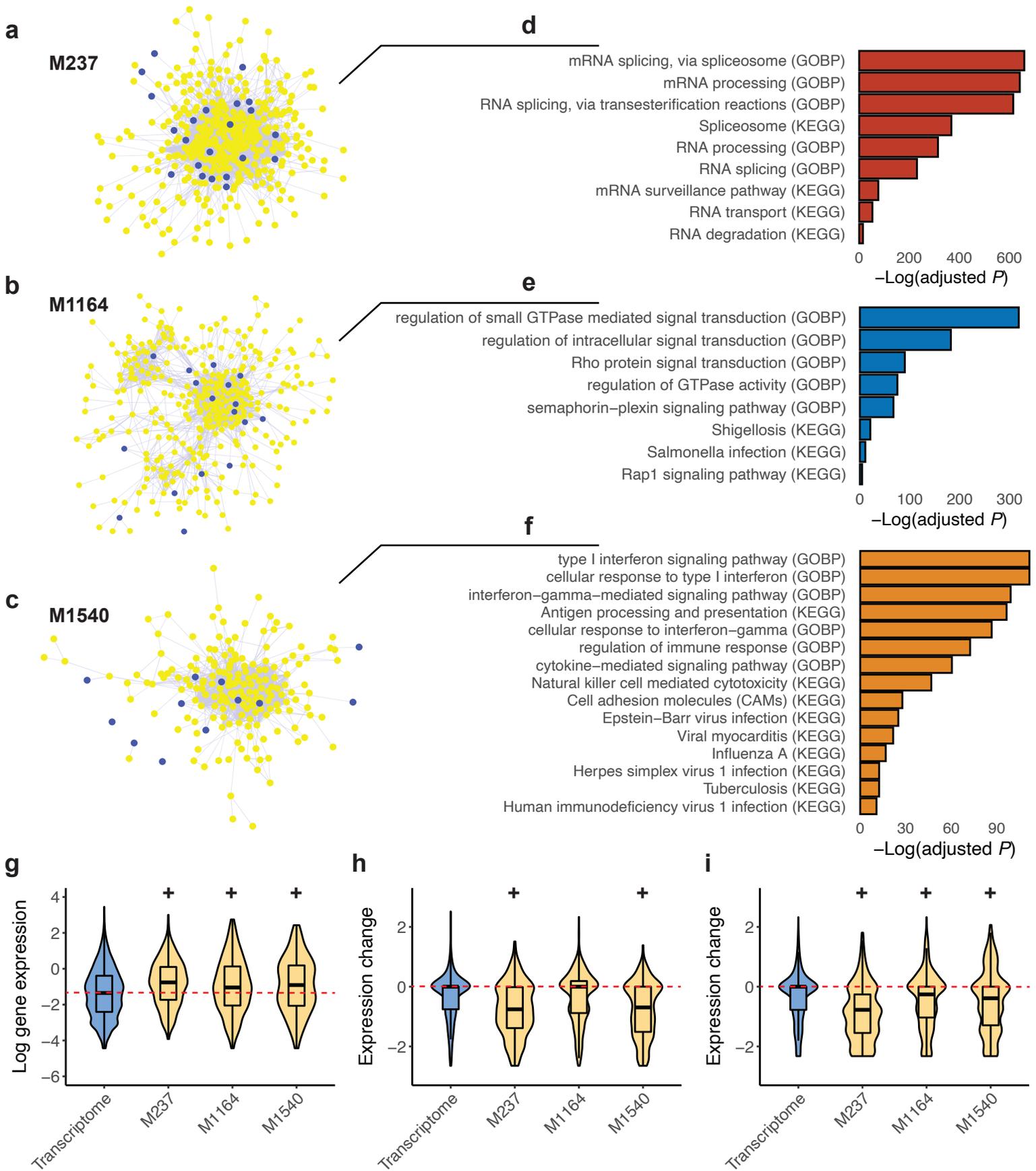
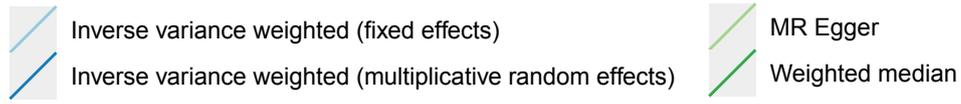


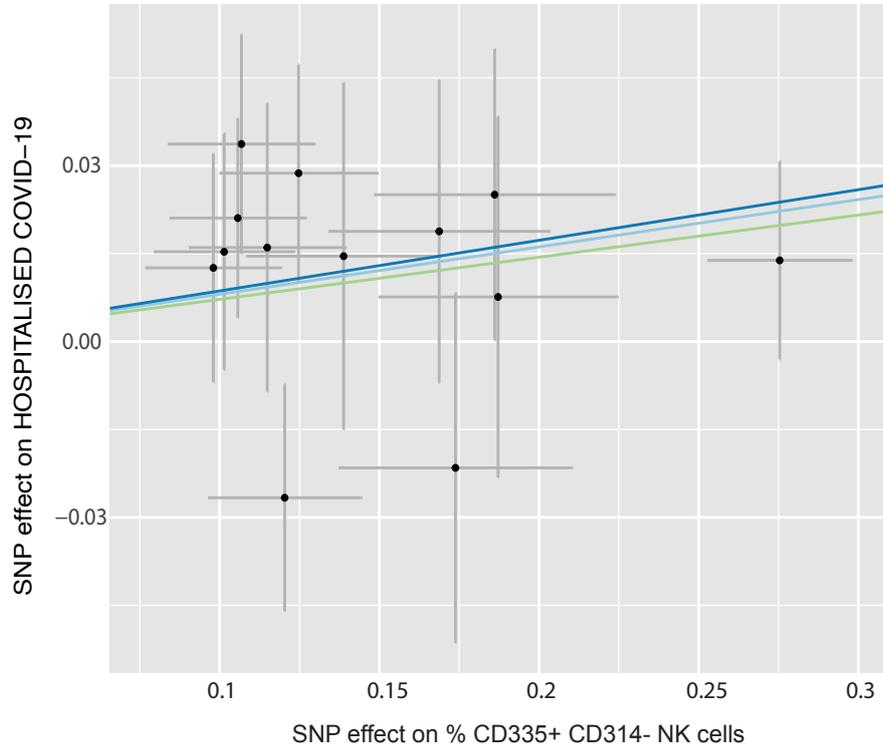
Figure 6



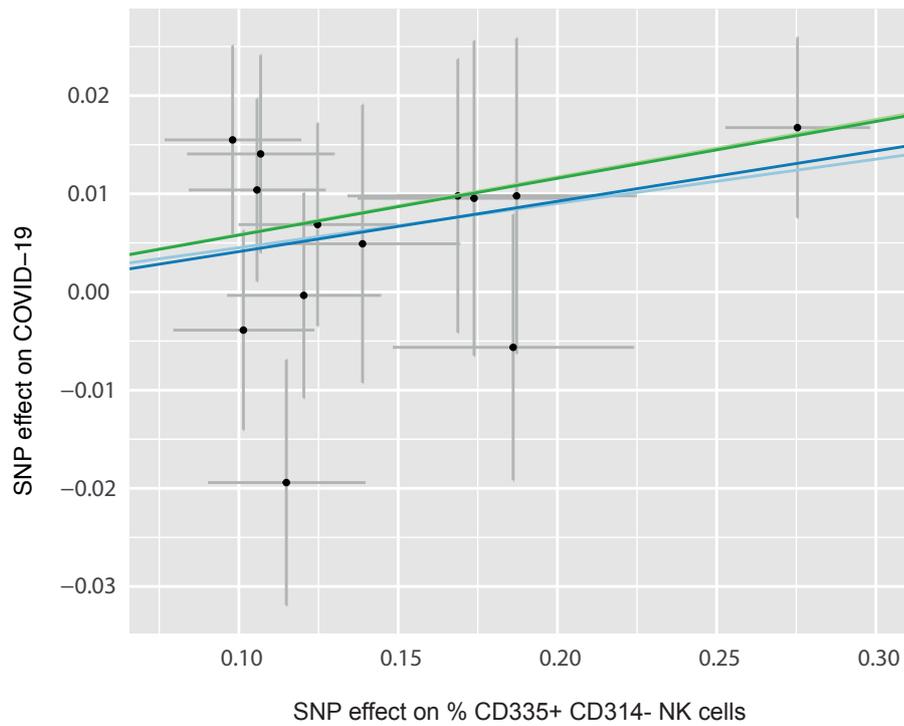
Supplementary Figure 1



a



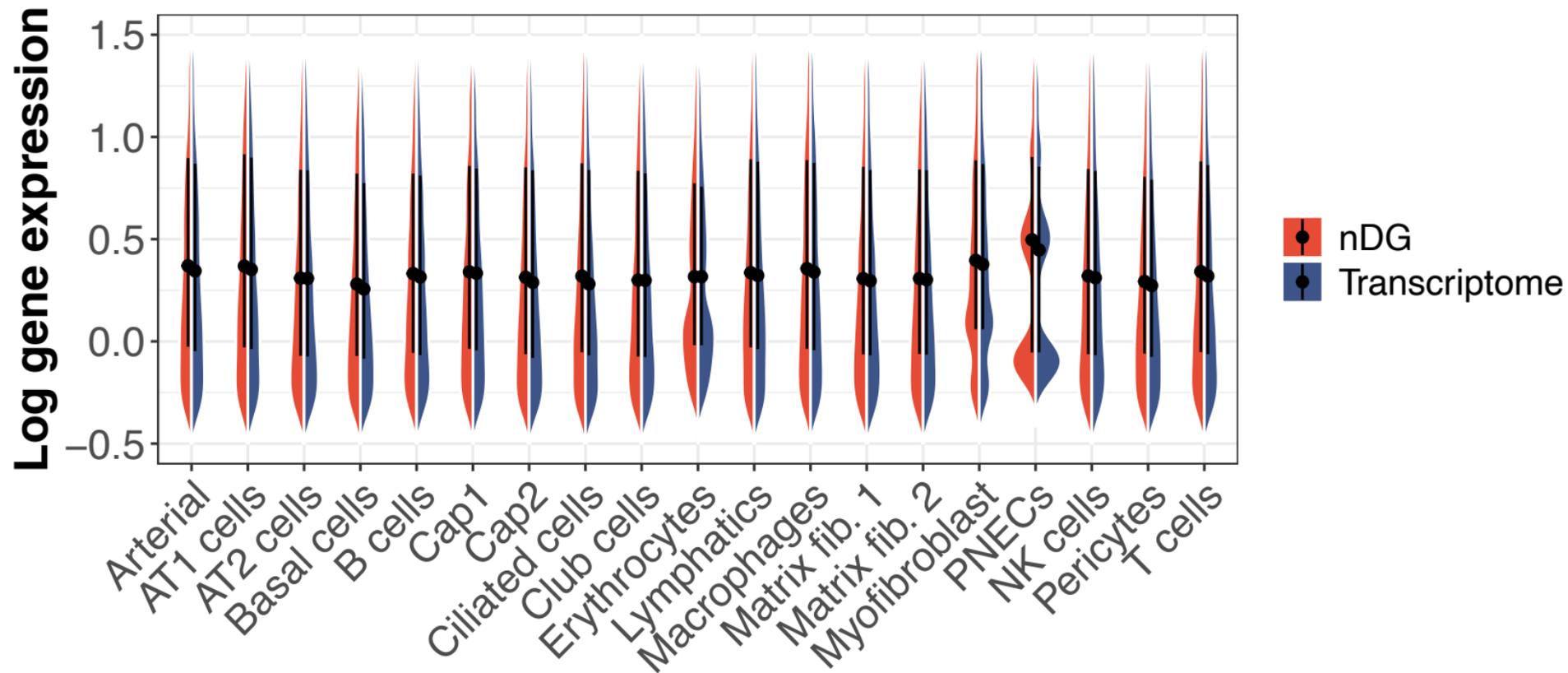
b



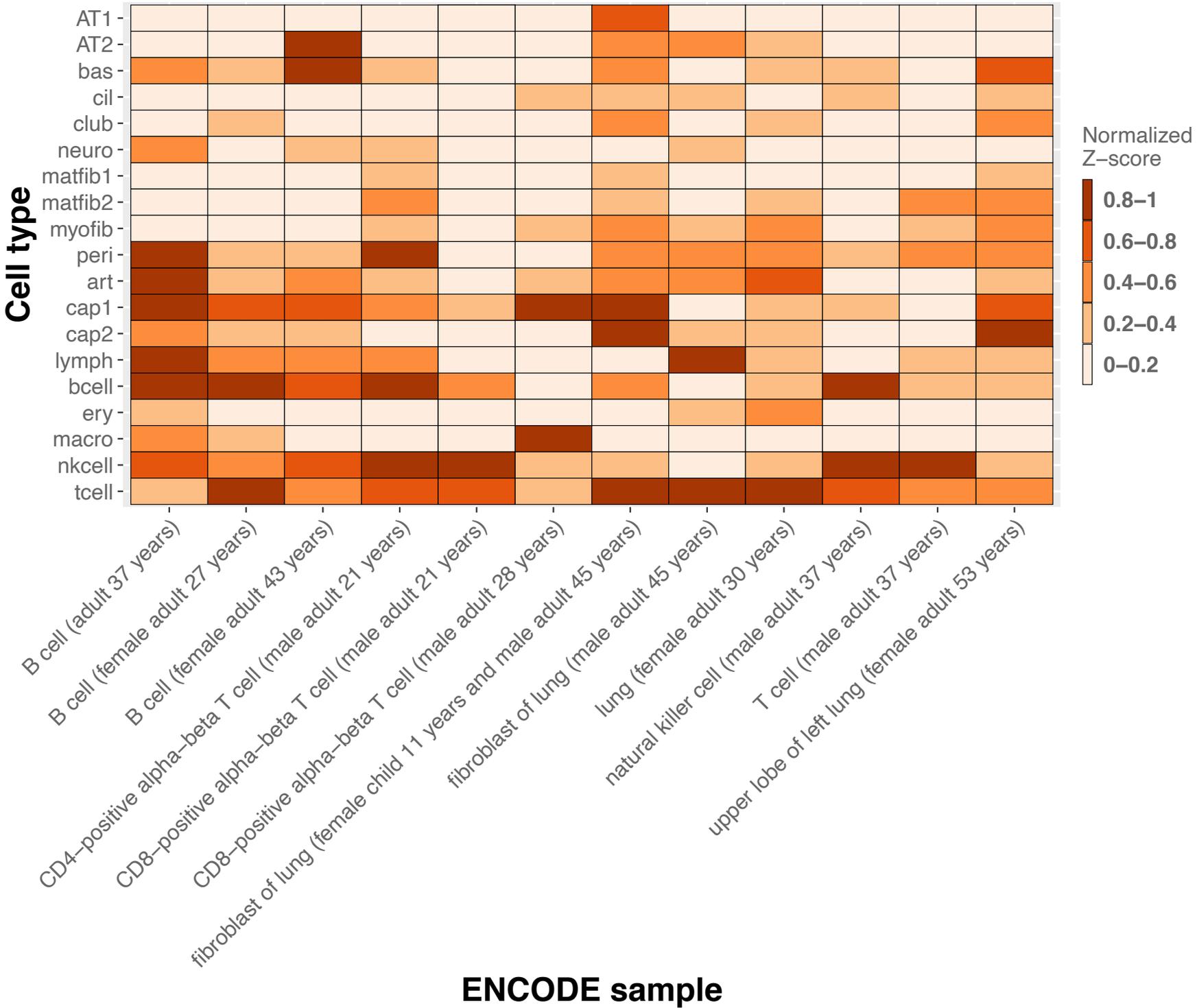
c

| Test | Hospitalised COVID-19 | | | COVID-19 | | |
|-----------------|-----------------------|----------------|-----------------|----------|----------------|-----------------|
| | % CD314+ | % CCR7- CD314- | % CD335+ CD314- | % CD314+ | % CCR7- CD314- | % CD335+ CD314- |
| IVW (mre) | 1.70E-01 | 8.01E-01 | 1.52E-02 | 1.74E-01 | 4.42E-01 | 9.90E-03 |
| MR Egger | 7.72E-01 | 4.38E-01 | 8.39E-01 | 9.25E-02 | 1.87E-01 | 3.83E-01 |
| Weighted Median | 5.22E-01 | 5.50E-01 | 1.01E-01 | 3.34E-01 | 1.59E-01 | 3.50E-02 |
| Weighted Mode | 5.82E-01 | 6.52E-01 | 2.36E-01 | 1.97E-01 | 2.63E-01 | 1.00E-01 |
| MR Lasso | 3.87E-01 | 8.71E-01 | 3.46E-02 | 2.67E-01 | 4.42E-01 | 2.73E-02 |

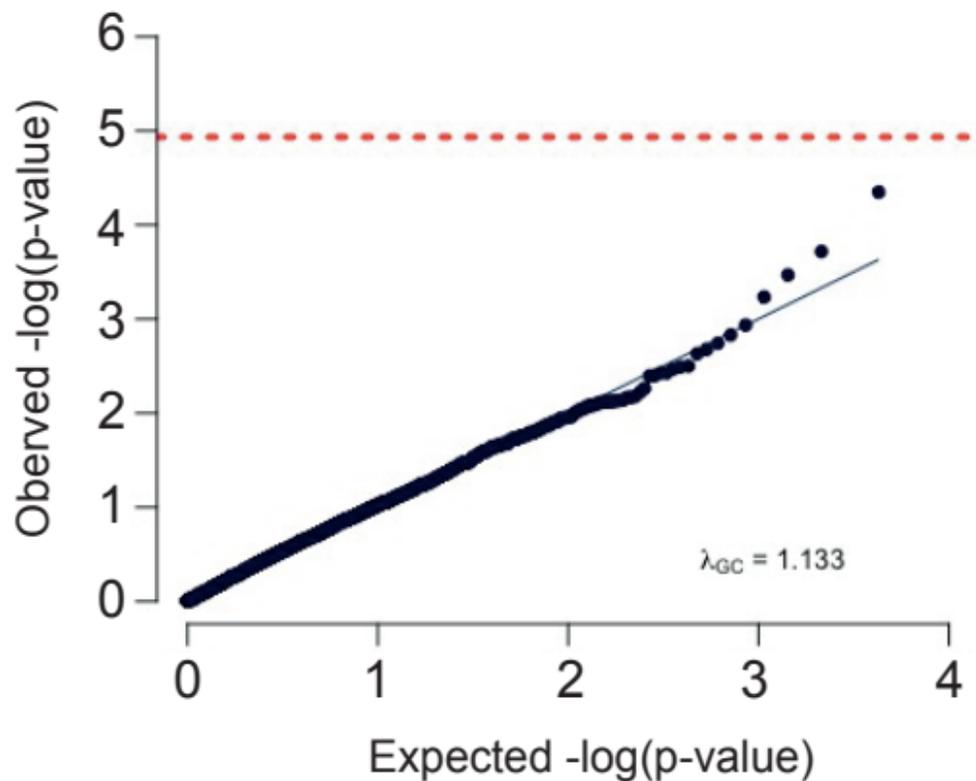
Supplementary Figure 2



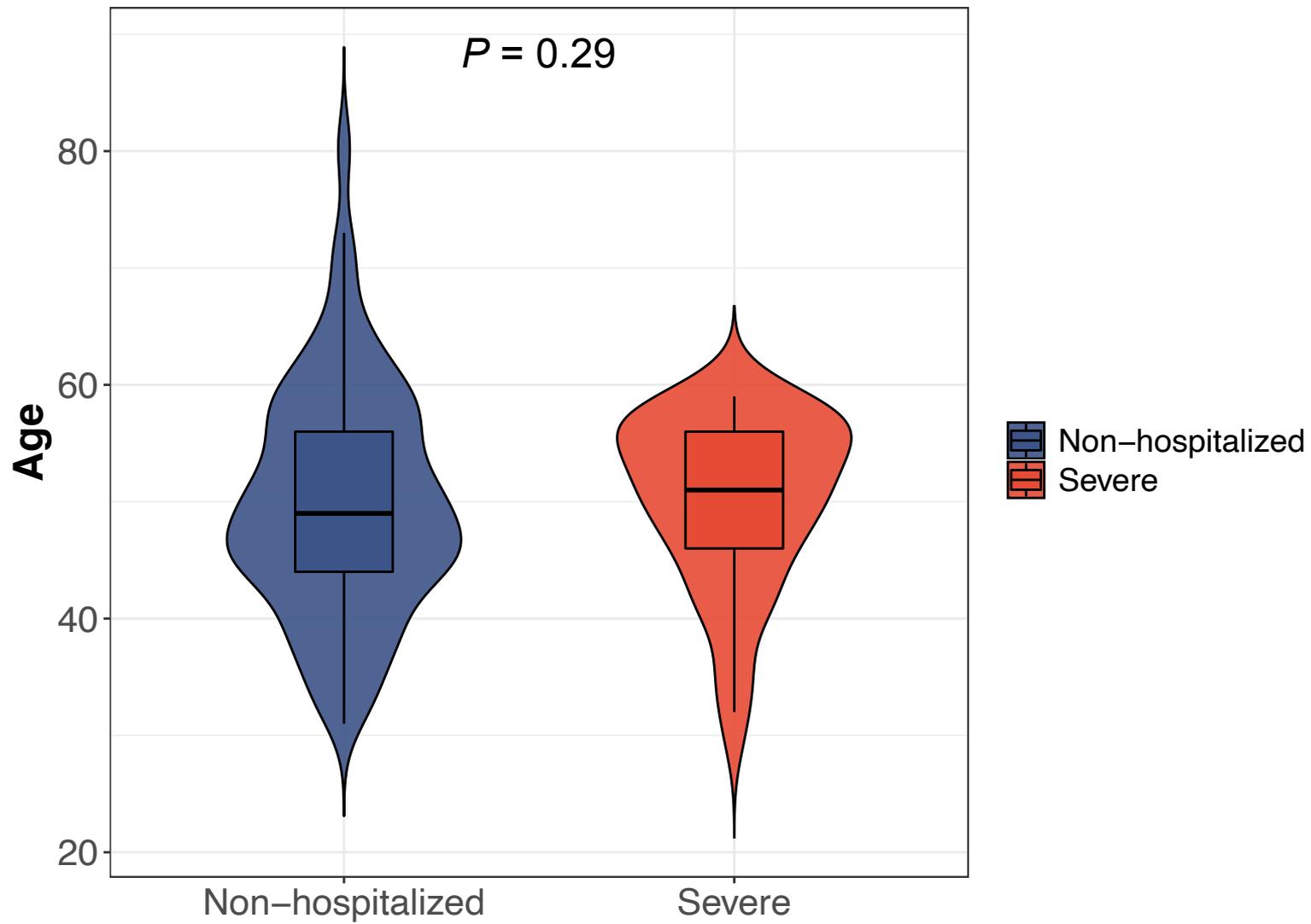
Supplementary Figure 3



Supplementary Figure 4

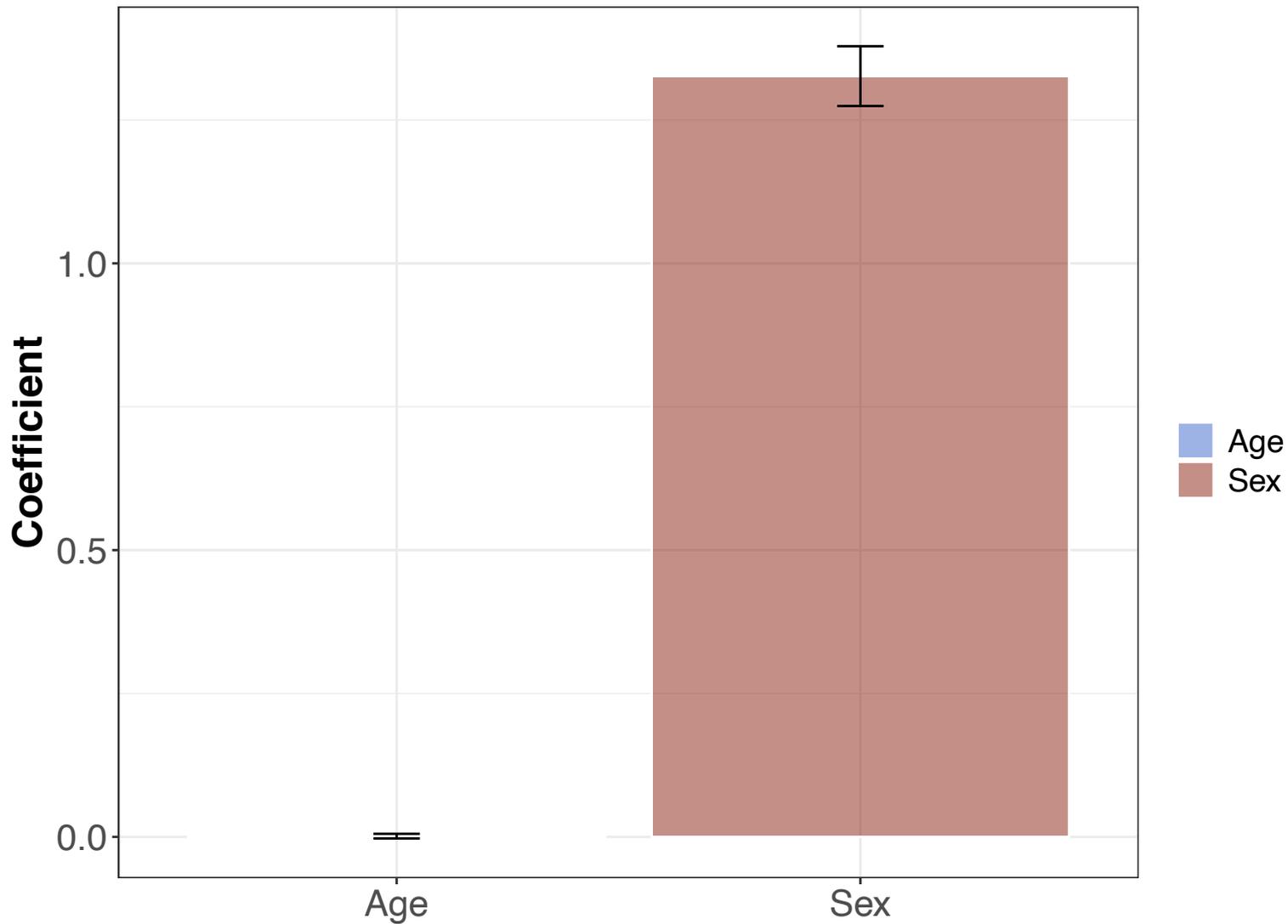


Supplementary Figure 5



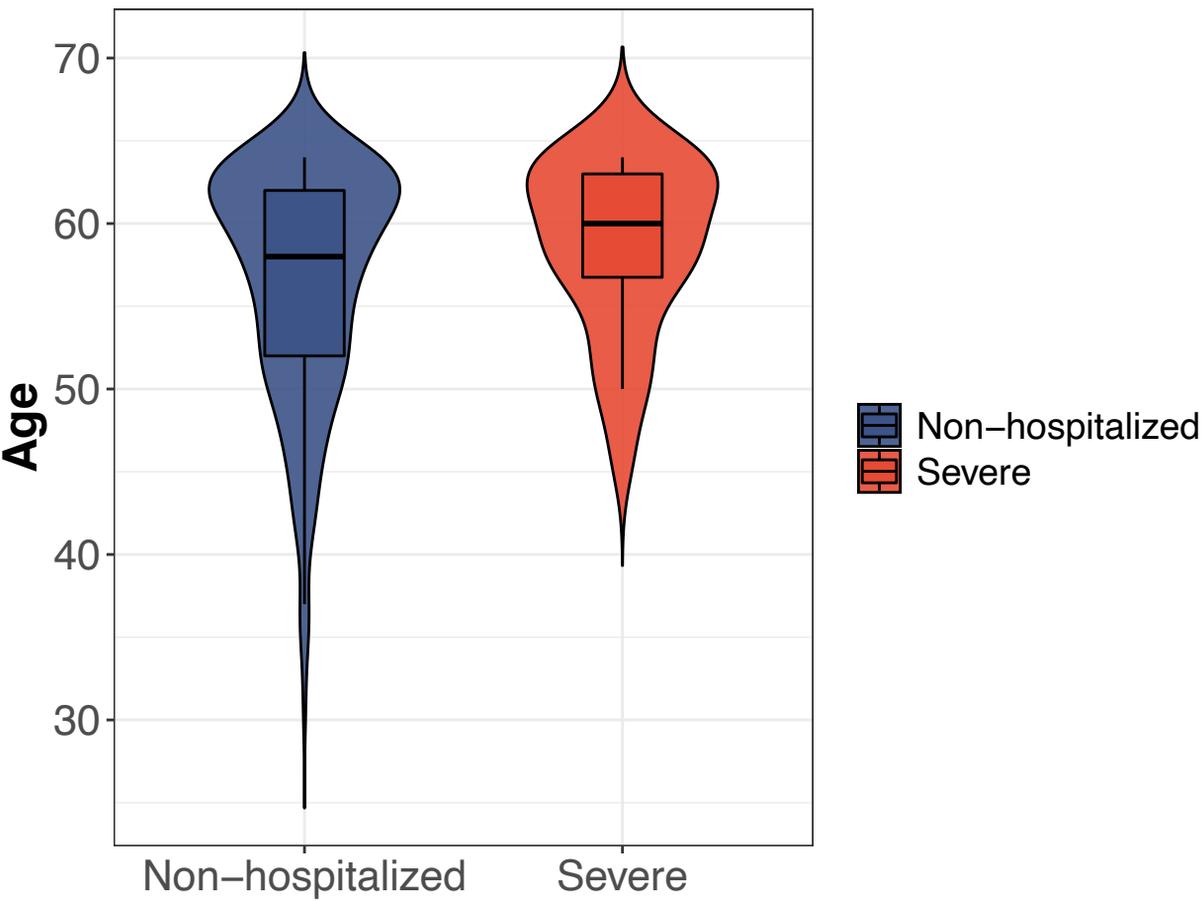
Supplementary Figure 6

Coefficients in logistic regression (MATLAB mnrfit)



Supplementary Figure 7

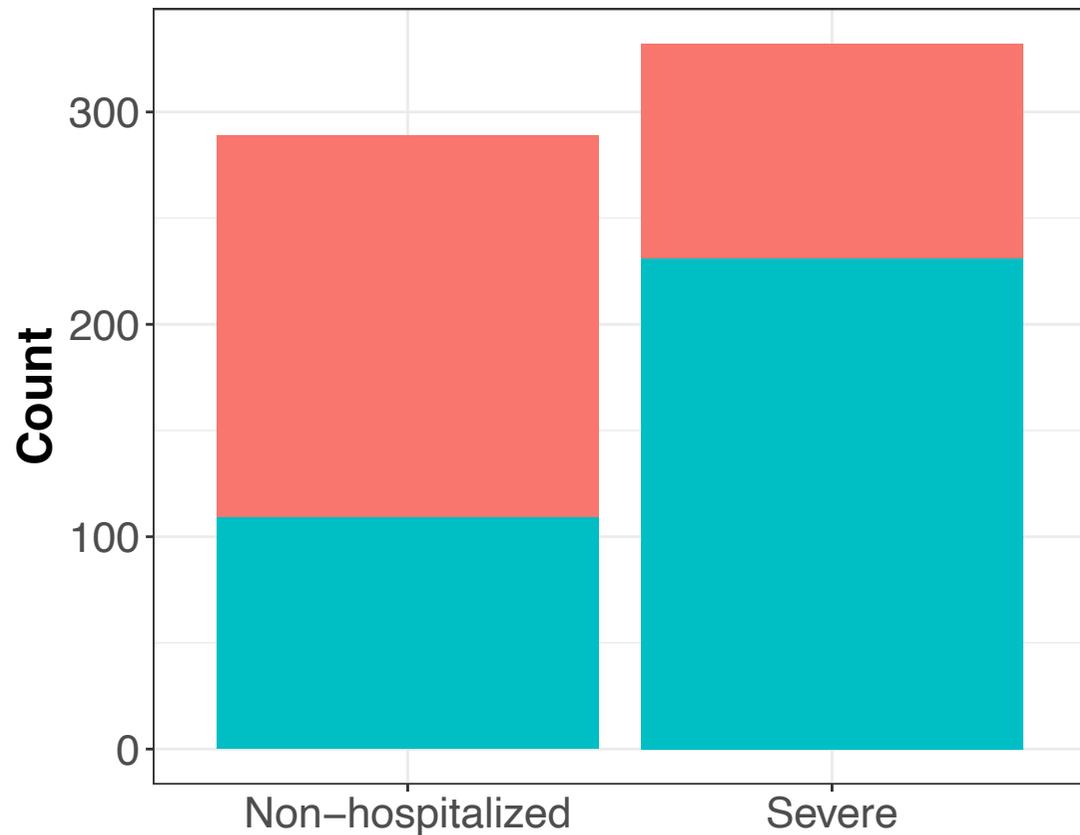
$P = 0.224$



Supplementary Figure 8

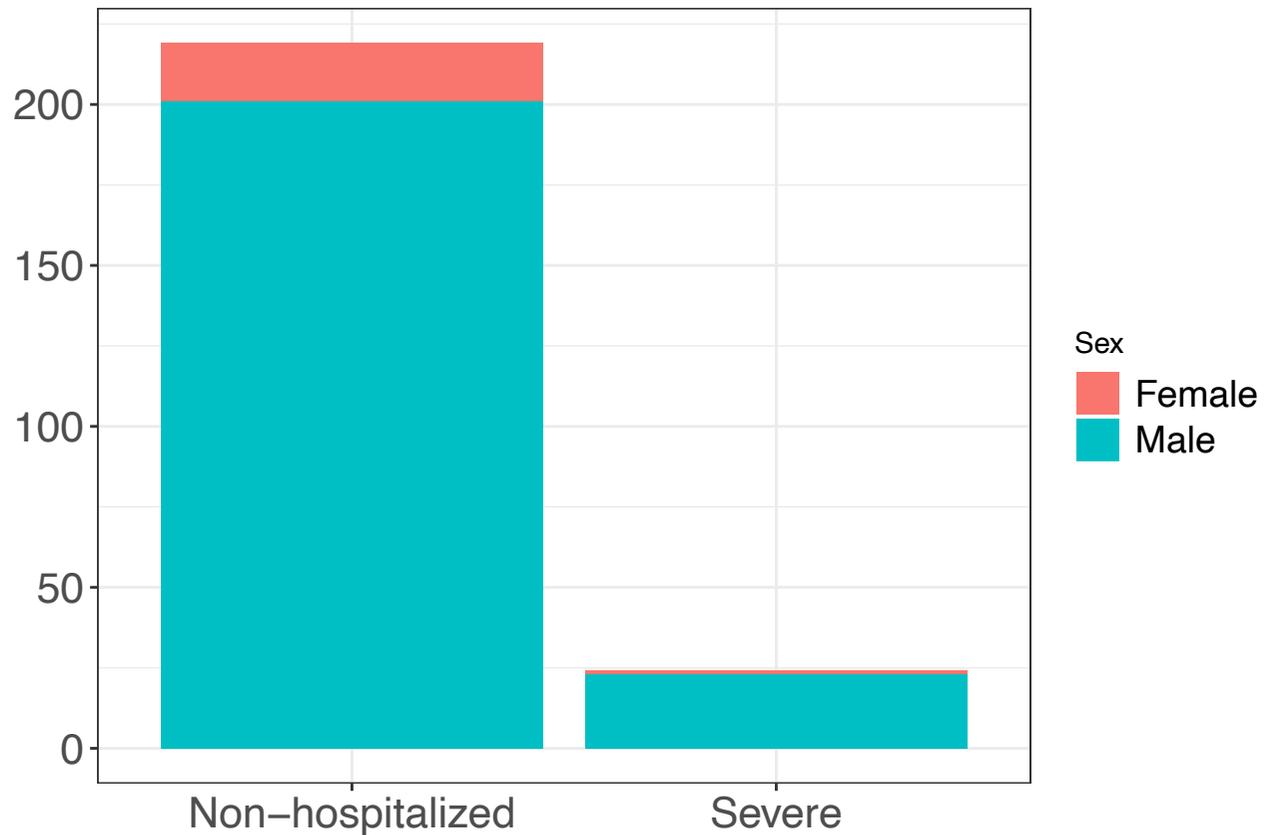
a

GEN-COVID

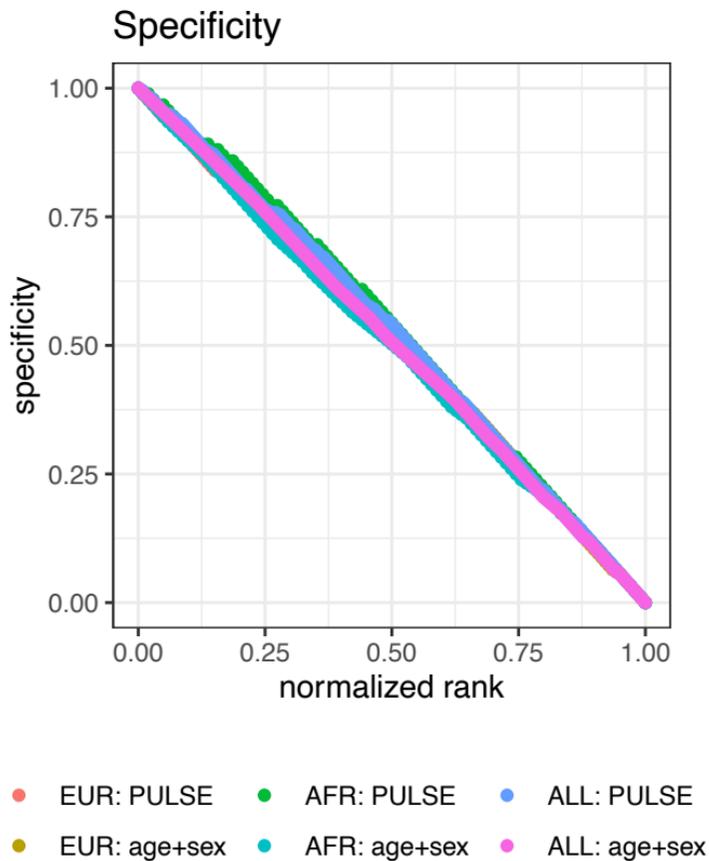


b

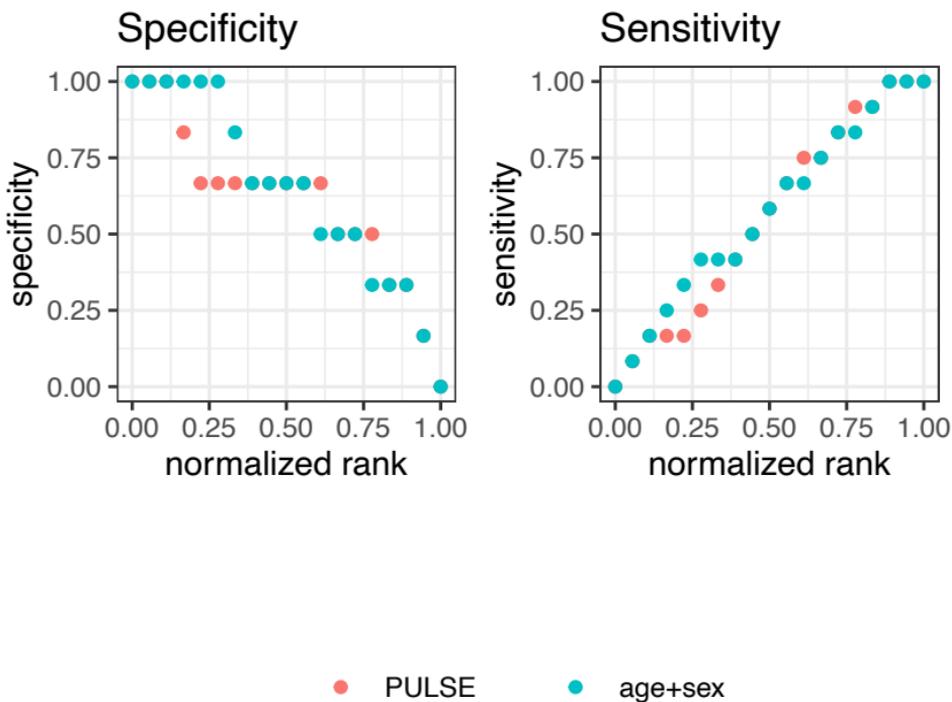
VA



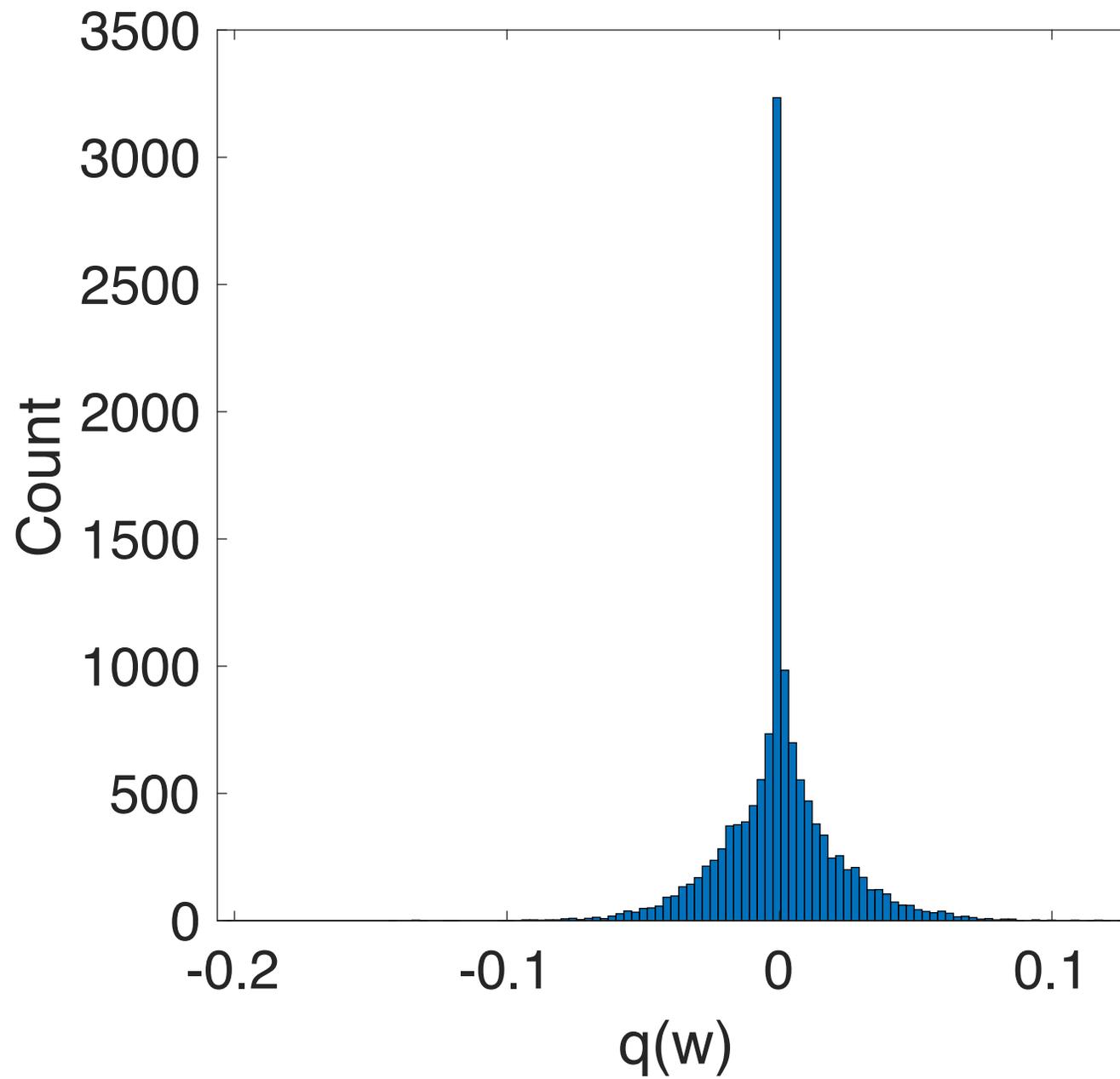
Supplementary Figure 9



Supplementary Figure 10



Supplementary Figure 11



Supplementary Notes for “Common and rare variant analyses combined with single-cell multiomics reveal cell-type-specific molecular mechanisms of COVID-19 severity”

1 Update rules of variational inference for PULSE

We provide update rules for the local and global variational parameters in PULSE.

1.1 Local variational method

As described in the Methods section in our main text, we used the local variational method [1] to handle the sigmoid function in variational inference (VI). Indeed, the sigmoid function involved in the Bernoulli distribution in Eq. 19 in the Methods section can be lower bounded by

$$\sigma(c_i) \geq h(c_i, \xi_i) = \sigma(\xi_i) \exp \left\{ (c_i - \xi_i)/2 - \chi(\xi_i)(c_i^2 - \xi_i^2) \right\}, \quad (1)$$

where

$$\chi(\xi) = \frac{1}{2\xi} \left(\sigma(\xi) - \frac{1}{2} \right), \quad (2)$$

$c_i = \mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2$, and ξ_i is a local variational parameter introduced to control the bound tightness. Therefore, the log-likelihood of observations is also lower bounded, i.e.,

$$\begin{aligned} \ln p(y_{1:N} | \mathbf{X}_{1:N}) &= \ln \int p(y_{1:N} | \mathbf{X}_{1:N}, \Theta) p(\Theta) d\Theta \\ &= \ln \int \left(\prod_{i=1}^N p(y_i | \mathbf{X}_i, \Theta) \right) p(\Theta) d\Theta \\ &= \ln \int \left(\exp \left\{ \sum_{i=1}^N c_i y_i \right\} \prod_{i=1}^N \sigma(-c_i) \right) p(\Theta) d\Theta \\ &\geq \ln \int \left(\exp \left\{ \sum_{i=1}^N c_i y_i \right\} \prod_{i=1}^N h(-c_i, \xi_i) \right) p(\Theta) d\Theta \\ &= \mathcal{L}(\xi_{1:N}). \end{aligned} \quad (3)$$

On one hand, we aim to perform variational inference based on the tractability of the lower bound $h(c_i, \xi_i)$, matching the proposal distribution with the true posterior. On the other hand, the variational parameters ξ_i 's need to be optimized by maximizing the lower bound $\mathcal{L}(\xi_{1:N})$ of the marginal likelihood, which achieves a better approximation after each

update. Therefore, we adopted the variational expectation-maximization (VEM) algorithm that solves both optimization problems simultaneously.

We first note that the “joint distribution”, denoted by \hat{p}^* , after lower bounding is not a proper density function, but by normalization, the inequality may not hold any more. Indeed, after normalizing, we get

$$p^* = \frac{1}{A(\xi_{1:N})} \hat{p}^*, \quad (4)$$

with

$$A(\xi_{1:N}) = \sum_{y_{1:N}} \int \left(\exp \left\{ \sum_{i=1}^N c_i y_i \right\} \prod_{i=1}^N h(-c_i, \xi_i) \right) p(\Theta) d\Theta. \quad (5)$$

Then, we can rewrite the lower bound $\mathcal{L}(\xi_{1:N})$ as

$$\begin{aligned} \mathcal{L}(\xi_{1:N}) &= \ln \int p^* d\Theta + \ln A(\xi_{1:N}) \\ &= \text{ELBO}_{p^*}(q, \xi_{1:N}) + \text{KL}(q \parallel p^*) + \ln A(\xi_{1:N}) \\ &= \text{ELBO}_{\hat{p}^*}(q, \xi_{1:N}) + \text{KL}(q \parallel p^*), \end{aligned} \quad (6)$$

resulting in a similar decomposition of the marginal log-likelihood to that in conventional VI.

As a consequence, we can perform the VEM as follows. (i) In the E-step where the variational parameters $\xi_{1:N}$ are fixed, the standard variational inference is performed to maximize the computationally feasible $\text{ELBO}_{\hat{p}^*}(q, \xi_{1:N})$ with respect to q . Here, everything in the mean-field variational inference (MFVI) keeps unchanged except replacing the sigmoid functions in the joint distribution by their lower bounds given by Eq. 1. This computes the approximate distribution best matching the true posterior, i.e., minimizing the KL divergence between q and p^* (see the last equation in Eq. 6). After the E-step, we approximately tighten the gap between $\mathcal{L}(\xi_{1:N})$ and the ELBO, and obtain $\mathcal{L}(\xi_{1:N}) \approx \text{ELBO}_{\hat{p}^*}(q, \xi_{1:N})$. (ii) In the M-step, we fix q and maximize the ELBO with respect to $\xi_{1:N}$, which increases $\mathcal{L}(\xi_{1:N})$ accordingly, as it is obvious that the inequality $\mathcal{L}(\xi_{1:N}) \geq \text{ELBO}_{\hat{p}^*}(q, \xi_{1:N})$ holds. Using VEM, we update q 's and ξ_i 's iteratively, gradually increasing the log-likelihood lower bound until reaching a local optimum and simultaneously yielding approximate posteriors with performance guarantee.

According to above discussions, we first get the lower bound of the log-likelihood of the conditional distribution over observations, i.e.,

$$\begin{aligned} \ln \prod_{i=1}^N p(y_i | \mathbf{X}_i, \mathbf{w}_1, \mathbf{w}_2) &= \sum_{i=1}^N y_i \ln \sigma(\mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2) + (1 - y_i) \ln(1 - \sigma(\mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2)) \\ &= \sum_{i=1}^N \mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2 y_i + \ln \sigma(-\mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2) \\ &\geq \sum_{i=1}^N \mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2 \left(y_i - \frac{1}{2} \right) - \chi(\xi_i) (\mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2)^2 + \ln \sigma(\xi_i) - \frac{1}{2} \xi_i + \chi(\xi_i) \xi_i^2, \end{aligned} \quad (7)$$

which serves as the basis for the inference of \mathbf{w}_1 and \mathbf{w}_2 . Then based on this lower bound

and the update principle of MFVI, the logarithm of $q(\mathbf{w}_1)$ can be calculated as

$$\begin{aligned} \ln q(\mathbf{w}_1) &\propto \mathbb{E}_{-\mathbf{w}_1} \left[-\frac{1}{2} \mathbf{w}_1^\top \boldsymbol{\Lambda} \mathbf{w}_1 + \sum_{i=1}^N \left(\mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2 \left(y_i - \frac{1}{2} \right) - \chi(\xi_i) (\mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2)^2 \right) \right] \\ &= -\frac{1}{2} \mathbf{w}_1^\top \left(\mathbb{E}[\boldsymbol{\Lambda}] + 2 \sum_{i=1}^N \chi(\xi_i) \mathbf{X}_i \mathbb{E}[\mathbf{w}_2 \mathbf{w}_2^\top] \mathbf{X}_i^\top \right) \mathbf{w}_1 + \mathbf{w}_1^\top \sum_{i=1}^N \left(y_i - \frac{1}{2} \right) \mathbf{X}_i \mathbb{E}[\mathbf{w}_2]. \end{aligned} \quad (8)$$

This indicates that $q(\mathbf{w}_1)$ follows a Gaussian defined as

$$q(\mathbf{w}_1; \tilde{\boldsymbol{\mu}}_{w_1}, \tilde{\boldsymbol{\Lambda}}_{w_1}) = \mathcal{N}(\mathbf{w}_1; \tilde{\boldsymbol{\mu}}_{w_1}, \tilde{\boldsymbol{\Lambda}}_{w_1}^{-1}), \quad (9)$$

where

$$\tilde{\boldsymbol{\mu}}_{w_1} = \tilde{\boldsymbol{\Lambda}}_{w_1}^{-1} \sum_{i=1}^N \left(y_i - \frac{1}{2} \right) \mathbf{X}_i \mathbb{E}[\mathbf{w}_2], \quad (10)$$

$$\tilde{\boldsymbol{\Lambda}}_{w_1} = \mathbb{E}[\boldsymbol{\Lambda}] + 2 \sum_{i=1}^N \chi(\xi_i) \mathbf{X}_i \mathbb{E}[\mathbf{w}_2 \mathbf{w}_2^\top] \mathbf{X}_i^\top. \quad (11)$$

The update rules expressed in Eqs. 10 and 11 are batch based, which is inefficient for large sample size or large feature dimension. We will transform this batch update into stochastic or mini-batch one based on the stochastic variational inference (SVI), scaling up the inference algorithm to big data. More details are shown in Section 1.3.

1.2 Reparameterization

To perform VI over the spike-and-slab prior defined in Eq. 22 in the Methods section of the main text, we adopted the reparameterization trick introduced in [4]. In particular, as discussed in the Methods section, \mathbf{w}_2 can be reparameterized by two additional variables \mathbf{s} and $\bar{\mathbf{w}}_2$ with

$$\mathbf{w}_2 = \bar{\mathbf{w}}_2 \circ \mathbf{s}, \quad (12)$$

where \circ means element-wise product. It can be easily shown that the new variable constructed by $\bar{w}_{2j} s_j$ follows the same distribution as w_{2j} . Then we can perform MFVI over $\bar{\mathbf{w}}_2$ and \mathbf{s} . However, the solution derived from a direct application of the fully factorized MFVI will deviate from the true posterior $q(\mathbf{w}_2)$ a lot, as the former is unimodal while the latter exponentially multimodal. To solve this problem, we followed [4], in which \bar{w}_{2j} and s_j are bundled together in the factorization. In particular, we assume the proposal distributions factorize as

$$q(\bar{\mathbf{w}}_2, \mathbf{s}) = \prod_{j=1}^M q(\bar{w}_{2j}, s_j), \quad (13)$$

Given the MFVI principle, after substituting w_{2j} with $\bar{w}_{2j}s_j$ in Eq. 7, we get

$$\begin{aligned}
\ln q(\bar{w}_{2j}, s_j) &\propto \mathbb{E}_{-\{\bar{w}_{2j}, s_j\}} \left[\sum_{i=1}^N \mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2 \left(y_i - \frac{1}{2} \right) - \chi(\xi_i) (\mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2)^2 \right. \\
&\quad \left. - \frac{1}{2} \lambda \bar{w}_{2j}^2 + s_j \ln \pi + (1 - s_j) \ln(1 - \pi) - \frac{1}{2} \sum_{l=1}^L r_l^{-1} \left(\hat{\mathbf{w}}_2^{(l)} - \mathbf{w}_2 \right)^\top \left(\hat{\mathbf{w}}_2^{(l)} - \mathbf{w}_2 \right) \right] \\
&\propto \mathbb{E}_{-\{\bar{w}_{2j}, s_j\}} \left[\sum_{i=1}^N \left(y_i - \frac{1}{2} \right) X_{i1j} \bar{w}_{2j} s_j - \chi(\xi_i) \left(X_{i1j}^2 \bar{w}_{2j}^2 s_j + 2X_{i1j} \left(\sum_{k \neq j} X_{i1k} w_{2k} \right) \bar{w}_{2j} s_j \right) \right. \\
&\quad \left. - \frac{1}{2} \lambda \bar{w}_{2j}^2 + s_j \ln \pi + (1 - s_j) \ln(1 - \pi) - \frac{1}{2} \sum_{l=1}^L r_l^{-1} \left(\bar{w}_{2j} s_j - 2\hat{w}_{2j}^{(l)} \bar{w}_{2j} s_j \right) \right] \\
&= -\frac{1}{2} \left(2 \sum_{i=1}^N \chi(\xi_i) \mathbb{E} [X_{i1j}^2] s_j + \mathbb{E}[\lambda] + \sum_{l=1}^L r_l^{-1} s_j \right) \bar{w}_{2j}^2 \\
&\quad + \left(\sum_{i=1}^N \left(y_i - \frac{1}{2} \right) \mathbb{E}[X_{i1j}] - 2\chi(\xi_i) \mathbb{E}[X_{i1j}] \mathbb{E} \left[\sum_{k \neq j} X_{i1k} w_{2k} \right] + \sum_{l=1}^L r_l^{-1} \mathbb{E} \left[\hat{w}_{2j}^{(l)} \right] \right) s_j \bar{w}_{2j} \quad (14)
\end{aligned}$$

where we define

$$X_{i1j} = \mathbf{w}_1^\top \mathbf{X}_i \mathbf{1}_j, \quad (15)$$

and $\mathbf{1}_j$ is a vector with all zeros but the j -th element one.

Since $q(\bar{w}_{2j}|s_j) \propto q(\bar{w}_{2j}, s_j)$, based on Eq. 14, we have

$$q(\bar{w}_{2j}|s_j = 0) = \mathcal{N} \left(\bar{w}_{2j}; \tilde{\mu}_{\bar{w}_{2j}|s_j=0}, \tilde{\lambda}_{\bar{w}_{2j}|s_j=0}^{-1} \right), \quad (16)$$

where

$$\tilde{\mu}_{\bar{w}_{2j}|s_j=0} = 0, \quad (17)$$

$$\tilde{\lambda}_{\bar{w}_{2j}|s_j=0} = \mathbb{E}[\lambda]. \quad (18)$$

Similarly, $q(\bar{w}_{2j}|s_j = 1)$ also follows a Gaussian given by

$$q(\bar{w}_{2j}|s_j = 1) = \mathcal{N} \left(\bar{w}_{2j}; \tilde{\mu}_{\bar{w}_{2j}|s_j=1}, \tilde{\lambda}_{\bar{w}_{2j}|s_j=1}^{-1} \right), \quad (19)$$

where

$$\begin{aligned}
\tilde{\mu}_{\bar{w}_{2j}|s_j=1} &= \tilde{\lambda}_{\bar{w}_{2j}|s_j=1}^{-1} \left(\sum_{i=1}^N \left(y_i - \frac{1}{2} \right) \mathbb{E}[X_{i1j}] \right. \\
&\quad \left. - 2\chi(\xi_i) \mathbb{E}[X_{i1j}] \mathbb{E} \left[\sum_{k \neq j} X_{i1k} w_{2k} \right] + \sum_{l=1}^L r_l^{-1} \mathbb{E} \left[\hat{w}_{2j}^{(l)} \right] \right), \quad (20)
\end{aligned}$$

$$\tilde{\lambda}_{\bar{w}_{2j}|s_j=1} = 2 \sum_{i=1}^N \chi(\xi_i) \mathbb{E} [X_{i1j}^2] + \mathbb{E}[\lambda] + \sum_{l=1}^L r_l^{-1}. \quad (21)$$

The stochastic updates of Eqs. 20 and 21 are shown in Section 1.3.

To derive $q(s_j)$, we use the Bayes' rule given by $q(s_j) = q(\bar{w}_{2j}, s_j)/q(\bar{w}_{2j}|s_j)$, yielding

$$q(s_j) = \text{Bern}(s_j; \tilde{\pi}_j), \quad (22)$$

where

$$\tilde{\pi}_j = \frac{\tilde{\rho}_{1j}}{\tilde{\rho}_{0j} + \tilde{\rho}_{1j}}, \quad (23)$$

and

$$\ln \tilde{\rho}_{0j} = \mathbb{E}[\ln(1 - \pi)] - \frac{1}{2} \ln \tilde{\lambda}_{\bar{w}_{2j}|s_j=0}, \quad (24)$$

$$\ln \tilde{\rho}_{1j} = \mathbb{E}[\ln \pi] + \frac{1}{2} \tilde{\lambda}_{\bar{w}_{2j}|s_j=1} \tilde{\mu}_{\bar{w}_{2j}|s_j=1}^2 - \frac{1}{2} \ln \tilde{\lambda}_{\bar{w}_{2j}|s_j=1}. \quad (25)$$

The posterior statistics of w_{2j} , including the expectation and variance, can be easily calculated based on Eqs. 12, 17, 18, 20 and 21. In particular, the posterior statistics of the marginal $q(\bar{w}_{2j})$ can be derived based on the laws of total expectation and variance, respectively.

1.3 Stochastic variational inference

As discussed in the Methods section in the main text, to scale up the inference algorithm to big data, we adopted SVI proposed in [3]. SVI updates variational parameters by summarizing data points based on stochastic gradient optimization, in which the natural gradient is used to account for measuring similarity between probability distributions. Thanks to the conditional conjugacy introduced in our model, the natural gradient enjoys a simple form without the calculation of the Hessian [3]. Then we can approximate the natural gradient by randomly sampling a single or a mini-batch of samples, greatly reducing the computational complexity per epoch. Here in our inference process, there are two steps where SVI needs to be applied.

(i) For the update of $q(\mathbf{w}_1)$ whose batch update is given by Eqs. 10 and 11, its stochastic update is given by

$$\phi_1^{(t)} = (1 - \epsilon_t) \phi_1^{(t-1)} + \epsilon_t \frac{N}{B} \sum_{i \in I} \left(y_i - \frac{1}{2} \right) \mathbf{X}_i \mathbb{E}[\mathbf{w}_2], \quad (26)$$

$$\phi_2^{(t)} = (1 - \epsilon_t) \phi_2^{(t-1)} + \epsilon_t \left(-\frac{1}{2} \mathbb{E}[\Lambda] - \frac{N}{B} \sum_{i \in I} \chi(\xi_i) \mathbf{X}_i \mathbb{E}[\mathbf{w}_2 \mathbf{w}_2^T] \mathbf{X}_i^T \right), \quad (27)$$

where ϕ_1 and ϕ_2 are natural parameters in the exponential family form for multivariate Gaussian, and I is a randomly sampled index set from $1 : N$ with size B . Then the distribution parameters in $q(\mathbf{w}_1)$ can be recovered by

$$\tilde{\boldsymbol{\mu}}_{w_1} = -\frac{1}{2} \phi_2^{-1} \phi_1, \quad (28)$$

$$\tilde{\Lambda}_{w_1} = -2\phi_2. \quad (29)$$

(ii) Similarly, for the update of $q(\bar{w}_{2j}|s_j = 1)$, its stochastic version is given by

$$\psi_{1j}^{(t)} = (1 - \epsilon_t)\psi_{1j}^{(t-1)} + \epsilon_t \left(\frac{N}{B} \sum_{i \in I} \left(y_i - \frac{1}{2} \right) \mathbb{E}[X_{i1j}] - 2\chi(\xi_i)\mathbb{E}[X_{i1j}]\mathbb{E} \left[\sum_{k \neq j} X_{i1k} w_{2k} \right] + \sum_{l=1}^L r_l^{-1} \mathbb{E} \left[\hat{w}_{2j}^{(l)} \right] \right), \quad (30)$$

$$\psi_{2j}^{(t)} = (1 - \epsilon_t)\psi_{2j}^{(t-1)} + \epsilon_t \left(-\frac{N}{B} \sum_{i \in I} \chi(\xi_i)\mathbb{E} [X_{i1j}^2] - \frac{1}{2}\mathbb{E}[\lambda] - \frac{1}{2} \sum_{l=1}^L r_l^{-1} \right), \quad (31)$$

where ψ_{1j} and ψ_{2j} are natural parameters in the exponential family form of Gaussian. In particular, the parameters in $q(\bar{w}_{2j}|s_j = 1)$ can be recovered by

$$\tilde{\mu}_{\bar{w}_{2j}|s_j=1} = -\frac{1}{2}\psi_{2j}^{-1}\psi_{1j}, \quad (32)$$

$$\tilde{\lambda}_{\bar{w}_{2j}|s_j=1} = -2\psi_{2j}. \quad (33)$$

1.4 Update rules for other global variational parameters

For other variational parameters, we perform standard MFVI and have

$$q(\Lambda; \tilde{\mathbf{W}}_\Lambda, \tilde{\nu}_\Lambda) = \mathcal{W}(\Lambda; \tilde{\mathbf{W}}_\Lambda, \tilde{\nu}_\Lambda), \quad (34)$$

$$q(\pi; \tilde{\alpha}_\pi, \tilde{\beta}_\pi) = \text{Beta}(\pi; \tilde{\alpha}_\pi, \tilde{\beta}_\pi), \quad (35)$$

$$q(\lambda; \tilde{a}_\lambda, \tilde{b}_\lambda) = \text{Gamma}(\lambda; \tilde{a}_\lambda, \tilde{b}_\lambda), \quad (36)$$

in which

$$\tilde{\mathbf{W}}_\Lambda^{-1} = \mathbf{W}_0^{-1} + \mathbb{E}[\mathbf{w}_1 \mathbf{w}_1^\top], \quad (37)$$

$$\tilde{\nu}_\Lambda = \nu_0 + 1, \quad (38)$$

$$\tilde{\alpha}_\pi = \alpha_0 + \sum_{j=1}^M \mathbb{E}[s_j], \quad (39)$$

$$\tilde{\beta}_\pi = \beta_0 + M - \sum_{j=1}^M \mathbb{E}[s_j], \quad (40)$$

$$\tilde{a}_\lambda = a_0 + \frac{M}{2}, \quad (41)$$

$$\tilde{b}_\lambda = b_0 + \frac{1}{2}\mathbb{E}[\bar{\mathbf{w}}_2^\top \bar{\mathbf{w}}_2], \quad (42)$$

1.5 Update rules for the local variational parameters

In addition to calculating posteriors, we also need to determine the local variational parameters ξ_i 's. According to our discussion in Section 1.1, we seek to optimizing ξ_i 's by

maximizing the lower bound $\mathcal{L}(\xi_{1:N})$ in Eq. 3. This corresponds to the M-step, in which the expected complete-data log-likelihood is maximized, i.e.,

$$\begin{aligned} Q(\boldsymbol{\xi}, \boldsymbol{\xi}^{\text{old}}) &\propto \mathbb{E} \left[\sum_{i=1}^N \ln \sigma(\xi_i) - \frac{1}{2} \xi_i - \chi(\xi_i) \left((\mathbf{w}_1^\top \mathbf{X}_i \mathbf{w}_2)^2 - \xi_i^2 \right) \right] \\ &= \sum_{i=1}^N \ln \sigma(\xi_i) - \frac{1}{2} \xi_i - \chi(\xi_i) \left(\text{Tr}(\mathbf{A}_i \text{Cov}[\mathbf{w}_1]) + \mathbb{E}[\mathbf{w}_1]^\top \mathbf{A}_i \mathbb{E}[\mathbf{w}_1] - \xi_i^2 \right), \end{aligned} \quad (43)$$

in which

$$\mathbf{A}_i = \mathbf{X}_i \mathbb{E}[\mathbf{w}_2 \mathbf{w}_2^\top] \mathbf{X}_i^\top. \quad (44)$$

By setting the derivate of Eq. 43 with respect to ξ_i to zero, we get

$$0 = \chi'(\xi_i) \left(\text{Tr}(\mathbf{A}_i \text{Cov}[\mathbf{w}_1]) + \mathbb{E}[\mathbf{w}_1]^\top \mathbf{A}_i \mathbb{E}[\mathbf{w}_1] - \xi_i^2 \right), \quad (45)$$

indicating that

$$\boxed{(\xi_i^{\text{new}})^2 = \text{Tr}(\mathbf{A}_i \text{Cov}[\mathbf{w}_1]) + \mathbb{E}[\mathbf{w}_1]^\top \mathbf{A}_i \mathbb{E}[\mathbf{w}_1]}. \quad (46)$$

Note that we can force ξ_i 's to be nonnegative without loss of generality due to the monotonicity of $\chi(\xi_i)$ when $\xi_i \geq 0$.

2 Update termination

To terminate the algorithm, we need to monitor the change of ELBO, whose computation is intense and undesirable. In this study, we followed the suggestions proposed in [2], in which we computed the average log predictive for a small held-out dataset to track ELBO evolution. We terminated the updates once the change of average log predictive fell below a threshold, indicating convergence. Here, we set $\text{tol} = 10^{-5}$ and terminate the algorithm when the proportion of change in ELBO is less than the tolerance. The inference algorithm is summarized in Algorithm 1.

References

- [1] C. M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006.
- [2] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- [3] M. D. Hoffman, D. M. Blei, C. Wang, and J. Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 14(4):1303–1347, 2013.
- [4] M. K. Titsias and M. Lázaro-Gredilla. Spike and slab variational inference for multi-task and multiple kernel learning. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 2339–2347. Curran Associates, Inc., 2011.

Algorithm 1: Stochastic MFVI for PULSE

Input : Model p , hyperparameters Θ and learning rate ϵ_t .

Output : Posteriors q and local variational parameters $\xi_{1:N}$.

```
1 Initialize variational parameters.
2 while not converged do
3   Randomly split the dataset into  $N/B$  mini-batches  $\mathcal{D}_{1:N/B}$ .
4   for  $i = 1 : N/B$  do
5     1. Update local variational parameters  $\xi_{1:N}$  based on Eq. 46.
6     2. Based on mini-batch  $\mathcal{D}_i$ , update  $\phi_1$ ,  $\phi_2$ ,  $\psi_{1j}$  and  $\psi_{2j}$  according to
7       Eqs. 26, 27, 30 and 31, respectively, and then update the corresponding
8       global variational parameters based on Eqs. 28, 29, 32 and 33.
9     3. Update other global variational parameters according to Eqs. 17, 18, 23,
10    37 to 40, successively.
11   end
12 Calculate average log predictive for the held-out dataset.
13 end
```
