

1 **Identification of Immune complement function as a determinant of adverse SARS-CoV-2**
2 **infection outcome**
3

4 Vijendra Ramlall^{1,2}, Phyllis M. Thangaraj^{1,3}, Cem Meydan^{5,6}, Jonathan Foon^{4,5}, Daniel Butler^{4,5}, Ben
5 May⁶, Jessica K. De Freitas^{7,8}, Benjamin S. Glicksberg^{7,8}, Christopher E. Mason^{4,5,9,10}, Nicholas P.
6 Tatonetti^{1,11*}, Sagi D. Shapira^{11*}
7

8 ¹ Department of Biomedical Informatics, Columbia University, New York, NY, USA.
9 USA.

10 ² Department of Physiology & Cellular Biophysics, Columbia University, New York, NY, USA.

11 ³ Vagelos College of Physicians and Surgeons, Columbia University, New York, NY, USA.

12 ⁴ Department of Physiology and Biophysics, Weill Cornell Medicine, New York, NY, USA

13 ⁵ The HRH Prince Alwaleed Bin Talal Bin Abdulaziz Alsaud Institute for Computational Biomedicine,
14 Weill Cornell Medicine, New York, NY, USA

15 ⁶ Herbert Irving Comprehensive Cancer Center, Columbia University, New York, NY, USA

16 ⁷ Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York,
17 NY, 10029

18 ⁸ Hasso Plattner Institute for Digital Health at Mount Sinai, Icahn School of Medicine at Mount Sinai,
19 New York, NY, 10065

20 ⁹ The WorldQuant Initiative for Quantitative Prediction, Weill Cornell Medicine, New York, NY, USA

21 ¹⁰ The Feil Family Brain and Mind Research Institute, Weill Cornell Medicine, New York, NY, USA

22 ¹¹ Department of Systems Biology, Columbia University, New York, NY, USA.

23 USA.

24 * address correspondence to nick.tatonetti@columbia.edu, ss4197@columbia.edu

25 **SUMMARY**

26 Understanding the pathophysiology of SARS-CoV-2 infection is critical for therapeutics and public
27 health intervention strategies. Viral-host interactions can guide discovery of regulators of disease
28 outcomes, and protein structure function analysis points to several immune pathways, including
29 complement and coagulation, as targets of the coronavirus proteome. To determine if conditions
30 associated with dysregulation of the complement or coagulation systems impact adverse clinical
31 outcomes, we performed a retrospective observational study of 11,116 patients who presented with
32 suspected SARS-CoV-2 infection. We found that history of macular degeneration (a proxy for
33 complement activation disorders) and history of coagulation disorders (thrombocytopenia, thrombosis,
34 and hemorrhage) are risk factors for morbidity and mortality in SARS-CoV-2 infected patients – effects
35 that could not be explained by age, sex, or history of smoking. Further, transcriptional profiling of
36 nasopharyngeal (NP) swabs from 650 control and SARS-CoV-2 infected patients demonstrated that in
37 addition to innate Type-I interferon and IL-6 dependent inflammatory immune responses, infection results
38 in robust engagement and activation of the complement and coagulation pathways. Finally, we conducted
39 a candidate driven genetic association study of severe SARS-CoV-2 disease. Among the findings, our
40 scan identified putative complement and coagulation associated loci including missense, eQTL and sQTL
41 variants of critical regulators of the complement and coagulation cascades. In addition to providing
42 evidence that complement function modulates SARS-CoV-2 infection outcome, the data point to putative
43 transcriptional genetic markers of susceptibility. The results highlight the value of using a multi-modal
44 analytical approach, combining molecular information from virus protein structure-function analysis with
45 clinical informatics, transcriptomics, and genomics to reveal determinants and predictors of immunity,
46 susceptibility, and clinical outcome associated with infection.

47 INTRODUCTION

48 The SARS-CoV-2 pandemic has had profound economic, social, and public health impact with over 6.1
49 million confirmed cases and over 370,000 deaths across the globe. The infection causes respiratory illness
50 with symptoms ranging from cough and fever to difficulty breathing. While highly variable age-
51 dependent mortality rates have been widely reported, the comorbidities that drive this dependence are not
52 fully understood. Further, with some notable exceptions¹⁻³, molecular studies have largely focused on
53 ACE-2, the receptor and determinant of cell entry and viral replication³. While ACE-2 expression is
54 critical, viruses employ a wide range of molecular strategies to infect cells, avoid detection, and
55 proliferate. In addition, viral replication and immune mediated pathology are the primary drivers of
56 morbidity and mortality associated with SARS-CoV-2 infection^{4,5}. Therefore, understanding how virus-
57 host interactions manifest as SARS-CoV-2 risk factors will facilitate clinical management, choice of
58 therapeutic interventions, and setting of appropriate social and public health measures.

59
60 Knowledge of the precise molecular interactions that control viral replicative cycles can delineate
61 regulatory programs that mediate immune pathology associated with infection and provide valuable clues
62 about disease determinants. For example, viruses, including SARS-CoV-2, deploy an array of genetically
63 encoded strategies to co-opt host machinery. Among the strategies, viruses encode multifunctional
64 proteins that harness or disrupt cellular functions, including nucleic acid metabolism and modulation of
65 immune responses, through protein-protein interactions and molecular mimicry – structural similarity
66 between viral and host proteins (for a full discussion please see accompanying paper). Recently, we
67 employed protein structure modeling to systematically chart interactions across all human infecting
68 viruses⁶ and in an accompanying paper, performed a virome-wide scan for molecular mimics. This
69 analysis points to broad diversification of strategies deployed by human infecting viruses and identifies
70 biological processes that underlie human disease. Of particular interest, we mapped over 140 cellular
71 proteins that are mimicked by coronaviruses (CoV). Among these, we identified components of the
72 complement and coagulation pathways as targets of structural mimicry across all CoV strains (see
73 companion paper).

74
75 Through activation of one of three cascades, (i) the classical pathway triggered by an antibody–antigen
76 complex, (ii) the alternative pathway triggered by binding to a host cell or pathogen surface, and (iii) the
77 lectin pathway triggered by polysaccharides on microbial surfaces, the complement system is a critical
78 regulator of host defense against pathogens including viruses⁷. When dysregulated by germline variants or
79 acquired through age-related effects or excessive acute and chronic tissue damage, complement activation
80 can contribute to pathologies mediated by inflammation⁷⁻⁹. Similarly, inflammation-induced coagulatory

81 programs -- which themselves can be regulated by the complement system -- as well as crosstalk between
82 pro-inflammatory cytokines and the coagulative and anticoagulant pathways play pivotal roles in
83 controlling pathogenesis associated with infections. Therefore, while the age-related differences in
84 susceptibility to SARS-CoV-2 are likely a consequence of multiple underlying variables, virally encoded
85 structural mimics of complement and coagulation pathway components may contribute to CoV associated
86 immune mediated pathology. Moreover, a corollary of these observations is that dysfunctions associated
87 with complement and/or coagulation may impact clinical outcome of SARS-CoV-2 infection. For
88 example, the companion study suggests that coagulation disorders, such as thrombocytopenia, thrombosis
89 and hemorrhage, may represent risk factors for SARS-CoV-2 clinical outcome. Among complement-
90 associated disorders, multiple genetic and experimental evidence (including animal models of disease,
91 histological examination of affected tissue, and germline mutational analysis) point to dysregulation of
92 the complement system as the major driver of both early-onset, and age-related macular degeneration
93 (AMD)⁸⁻¹¹. A hyperinflammatory phenotype mediated by complement leads to progressive immune-
94 mediated deterioration of the central retina. While AMD, the leading cause of blindness in elderly
95 individuals (affecting roughly 200 million people worldwide¹¹), is likely the result of multiple
96 pathological processes, dysregulation of complement activation has emerged as the most widely accepted
97 cause of disease⁹⁻¹².

98
99 To determine if conditions associated with dysregulation of the complement or coagulation systems
100 impact adverse clinical outcomes associated with SARS-CoV-2 infection, we conducted a retrospective
101 observational study of 11,116 patients at New York-Presbyterian/Columbia University Irving Medical
102 Center. In agreement with previous reports¹³, survival analysis identified significant risk of mechanical
103 respiration and mortality associated with age and sex, as well as history of hypertension, obesity, type 2
104 diabetes (T2D), and coronary artery disease (CAD). Moreover, we found that patients with history of
105 macular degeneration (a proxy for complement activation disorders) and coagulation disorders (i.e.
106 thrombocytopenia, thrombosis, and hemorrhage) were at significantly increased risk of adverse clinical
107 outcomes (including mechanical respiration and death) following SARS-CoV-2 infection. Importantly,
108 these effects could not be explained by either age or sex, nor did we find any evidence that history of
109 smoking contributes to risk of adverse clinical outcomes associated with SARS-CoV-2 infection.
110 Conversely, albeit in a small number of individuals, we observed that no patients with complement
111 deficiency disorders required mechanical respiration or succumbed to their illness. In addition,
112 transcriptional profiling of nasopharyngeal (NP) swabs from 650 control and SARS-CoV-2 infected
113 patients demonstrates that in addition to innate Type-I interferon and IL-6 dependent inflammatory

114 immune responses, infection results in robust engagement and activation of the complement and
115 coagulation pathways.

116

117 Finally, a focused analysis of proximal and distal variants of complement and coagulation components
118 using the April 2020 COVID data released by the UK Biobank revealed genetic markers associated with
119 severe SARS-CoV-2 infection. Among our findings, we identified variants in CD55 (a negative regulator
120 of complement activation¹⁴), CFH and C4BPA, which play central roles in complement activation and
121 innate immunity. Importantly, analysis of the May 2020 COVID data released by the UK Biobank
122 recapitulated these results and identified additional variants. For example, the scan revealed that variants
123 in Alpha-2-macroglobulin (A2M), a protease inhibitor and cytokine transporter which participates in the
124 formation of fibrin clots and regulates inflammatory cascades, were associated with adverse clinical
125 outcome. In addition to providing evidence that complement function modulates SARS-CoV-2 infection,
126 the data point to several putative genetic markers of susceptibility. The results highlight the value of using
127 a multi-modal analytical approach, combining molecular information from virus protein structure-
128 function analysis with clinical informatics, transcriptomics, and genomics to reveal determinants and
129 predictors of immunity, susceptibility, and clinical outcome associated with infection.

130

131 **RESULTS**

132 *Comorbidity statistics and covariances in a retrospective observational clinical cohort*

133 To explore if conditions associated with dysregulation of the complement or coagulation systems impact
134 adverse clinical outcomes associated with SARS-CoV-2, we conducted a retrospective observational
135 study of patients treated at New York-Presbyterian/Columbia University Irving Medical Center for
136 suspected infection (Table 1). Electronic health records (EHR) were used to define sex, age, and smoking
137 history status as well as histories of macular degeneration, coagulatory disorders (i.e. thrombocytopenia,
138 thrombosis, and hemorrhage), hypertension, type 2 diabetes, coronary artery disease, and obesity (see
139 Methods). As shown in Table 1, of the 11,116 patients that presented to the hospital between February 1,
140 2020 and April 25, 2020 with suspected SARS-CoV-2 infection, 6,398 tested positive for the virus.
141 Among these, 88 were patients with a history of macular degeneration, four were patients with
142 complement deficiency disorders, and 1,179 were patients with disorders associated with the coagulatory
143 system. In addition, hypertension, coronary artery disease, diabetes, obesity, and annotated cough were
144 represented by 1,922, 1,566, 847, 791, and 727 patients, respectively (Table 1). While CAD,
145 hypertension, T2D, obesity, and coagulation disorders represent a group with the highest covariance, we
146 find lower co-occurrence between these conditions and macular degeneration in both SARS-CoV-2
147 positive and negative individuals (Figure S1). In addition to these medical histories, smoking status, past

148 or present, was noted for 5,079 patients (of 1,359 smokers included in the study, 723 were SARS-CoV-2
149 positive). Finally, of patients who were put on mechanical ventilation, we observed a 35% mortality rate,
150 and 31% of deceased patients had been on mechanical respiration.

151

152 *Macular degeneration and coagulation disorders are associated with SARS-CoV-2 outcomes*

153 We estimated the univariate and age- and sex-corrected risk associated with baseline clinical history of
154 previously reported SARS-CoV-2 risk factors (including hypertension, obesity, type 2 diabetes, and
155 coronary artery disease) as well as coagulation and complement disorders using survival analysis and Cox
156 proportional hazards regression modeling. As shown in Figure 1 and Table 1, we identified significant
157 risk of mechanical respiration and mortality associated with age and sex, as well as history of
158 hypertension, obesity, and type 2 diabetes (T2D), coronary artery disease (CAD). Notably, we did not
159 find evidence that smoking status (past or present) is a significant risk factor for either mechanical
160 respiration or mortality. We found that those with a history of macular degeneration (a proxy for
161 complement activation disorders) and coagulation disorders (thrombocytopenia, thrombosis, and
162 hemorrhage) were at significantly increased risk of adverse clinical outcomes (including mechanical
163 respiration and death) following SARS-CoV-2 infection (Figure 1, Table 1). Specifically, we observed a
164 mechanical respiration rate of 15.9% (95% CI: 8.3-23.6; HR: 2.2, P value = 0.0046) and a mortality rate of
165 25% (95% CI: 16.0-34.0; HR 3.0, P value = 4.4×10^{-7}) among patients with a history of macular
166 degeneration, and rates of 9.4% (95% CI: 7.7-11.1; HR 1.5, P value = 9.6×10^{-5}) and 14.7% (95% CI: 12.7-
167 16.7; HR: 2.3, P value = 1.8×10^{-23}) for mechanical respiration and mortality, respectively, among patients
168 with coagulation disorders (Table 1). Moreover, as shown in Figure 1b, patients with a history of macular
169 degeneration appear to succumb to disease more rapidly than others. Critically, the contribution of age
170 and sex was not sufficient to explain the increased risks associated with history of macular degeneration
171 (Age/Sex-Corrected mechanical respiration HR=1.8 95% CI: 1.1-3.2, P value = 0.024; Age/Sex-Corrected
172 mortality HR=1.7 95% CI: 1.1-2.5, P value = 0.022) or coagulation disorders (Age/Sex-Corrected
173 mechanical respiration HR=1.5. 95% CI: 1.2-1.8, P value = 2.4×10^{-4} ; Age/Sex-Corrected mortality
174 HR=1.8 95% CI: 1.5-2.1, P value = 3.4×10^{-12}). Conversely, albeit in a small number of individuals, we
175 observed that among patients with complement deficiency disorders, who are normally at increased risk
176 of complications associated with infections, none required mechanical respiration or succumbed to their
177 illness (Table 1, Figure 1a and 1b). Importantly, while the correlation between macular degeneration or
178 coagulopathies and established covariates included in this study is low (as shown in Supplemental Figure
179 S1 and Supplemental Table S1, Tanimoto coefficients between 0.038 and 0.050 and 0.25 and 0.38,
180 respectively), further study, perhaps with larger patient cohorts, will be necessary to rule out
181 comorbidities that may be associated with macular degeneration and coagulopathies. Together, these data

182 suggest that hyper-active complement and coagulative states predispose individuals to adverse outcomes
183 associated with SARS-CoV-2 infection, and that deficiencies in complement components may be
184 protective. Importantly, given the low incidence rate of deficiencies in either complement or coagulation
185 pathways, further analysis with larger clinical cohorts is warranted.

186
187 *SARS-CoV-2 infection induces robust transcriptional regulation of complement and coagulation*
188 *components.*

189 Transcriptional responses of human NP epithelial cells during viral infection can provide critical
190 information about underlying immune programs. We leveraged whole genome RNA sequencing (RNA-
191 seq) profiles to identify differentially regulated genes and pathways in 650 NP swabs from control and
192 SARS-CoV-2 infected patients who presented to Weill-Cornell Medical Center. As shown in Figure 2a,
193 gene set enrichment analysis (GSEA) of HALLMARK gene sets found that SARS-CoV-2 infection (as
194 defined by presence of SARS-CoV-2 RNA and stratified into ‘positive’, ‘low’, ‘medium’ or ‘high’ based
195 on viral load; see Methods) induces genes related to pathways with known immune modulatory functions,
196 including ‘inflammatory_response’, ‘interferon_alpha_response’, and ‘IL6_JAK_STAT3_signaling (FDR
197 corrected P value < 0.001 ; Figure 2a). Moreover, we found that among the most enriched gene sets,
198 SARS-CoV-2 infection induces robust activation of the complement cascade (FDR corrected P value $<$
199 0.001), with increasing enrichment and significance with viral load (FDR corrected P value < 0.0001). We
200 extended the analysis to include all complement and coagulation associated gene sets in MsigDB and
201 identified ‘KEGG_Complement_and_Coagulation_Cascades’, ‘GO_Coagulation’, as well as
202 ‘Reactome_initial_triggering_of_complement’ to be enriched in expression profiles of SARS-CoV-2
203 infected samples (Q value < 0.05 ; representative GSEA profiles are shown in Figure 2b and a full list of
204 enriched pathways and gene sets can be found at <https://masonlab.shinyapps.io/CovidGenes/>). As
205 highlighted in Figure 2c-e, the pathway-level transcriptional regulation induced by SARS-CoV-2
206 identified by GSEA is also observed at the individual gene level for upregulated and downregulated
207 regulated transcripts as well as those that are particularly upregulated in the context of high viral load
208 (Figure 2d, e, f, respectively). Taken together, the data demonstrate that in addition to immune factors like
209 Type I interferons and dysregulation of IL6-dependent inflammatory responses which has been linked to
210 poor clinical outcome¹³, transcriptional control of complement and coagulation cascades is a feature of
211 SARS-CoV-2 infection.

212
213 *Genetic variation in complement and coagulation pathway components is associated with adverse SARS-*
214 *CoV-2 infection outcome*

215 The data highlighted above provide evidence that complement and coagulation disorders play a role in
216 SARS-CoV-2 infection outcome and that infection with this virus induces robust transcriptional
217 regulation of complement and coagulation pathway components. Moreover, dysfunction of complement
218 or coagulation cascades can be the result of either acquired dysregulation, genetically encoded variants, or
219 both. However, any genetic factors that may underlie the clinical trends we observed remain hidden due
220 to the retrospective nature of the study and the lack of available genetic data on these patients. On the
221 other hand, the UK Biobank, a prospective cohort study with deep genetic, physical, and health data
222 collected on ~500,000 individuals across the United Kingdom^{15,16}, recently released SARS-CoV-2
223 infection and outcome statuses for 1,474 patients, allowing for genetic and epidemiological associations
224 to be assessed. The release in April 2020 included 669 patients who tested positive for the virus, 572 of
225 whom required hospitalization.

226
227 We conducted a candidate driven study to evaluate if genetic variation in components of complement or
228 coagulation pathways are associated with poor SARS-CoV-2 clinical outcome. Briefly, we focused our
229 analysis on 337,147 (181,032 female) subjects of White British descent, excluding 3rd degree and above
230 relatedness and without aneuploidy¹⁵. Applying these restrictions to the April-2020 cohort resulted in 910
231 patients with suspected infection (388 positive, 332 positive and hospitalized; see *Methods*). As detailed
232 Supplemental Table S2, of the 805,426 genetic variants profiled in the UK Biobank, 2,888 are within a
233 60Kb window around 102 genes with known roles in regulating complement or coagulation cascades
234 (results that follow are robust to varying window size between 40Kb-80Kb; see *Methods*, Figure 3a-b).
235 We focused our analysis on single-nucleotide polymorphisms (SNP) with minor allele frequency (MAF)
236 above 1% and, as shown in Figure 3 and Supplemental Figure S2a-f, used an empirical permutation
237 analysis to set the study-wide significance alpha (α) thresholds for each analysis described below (see
238 *Methods*). As highlighted in Figure 3c and further detailed in Supplemental Table S2, we identified 11
239 loci representing 7 genes with study-wide significance ($\alpha = 0.001$) in the April-2020 cohort. Among
240 these, and proximal to coagulation factor III (F3), is variant rs72729504 which we find to be associated
241 with increased risk of adverse clinical outcome associated with SARS-CoV-2 infection (OR: 1.93). Fibrin
242 fragment D-dimer, one of several peptides produced when cross-linked fibrin is degraded by plasmin, is
243 the most widely used clinical marker of activated blood coagulation. Among the genetic loci that
244 influence D-dimer levels, GWAS studies have identified mutations in F3 as having the strongest
245 association¹⁷. Importantly, increased D-dimer levels were recently reported to correlate with poor clinical
246 outcome in SARS-CoV-2 infected patients¹³. So, while the functional role of rs72729504 remains to be
247 elucidated, our observations suggest that this locus may represent a genetic marker of SARS-CoV-2
248 susceptibility and outcomes.

249

250 In addition to the SNP highlighted above, we identified 4 variants (rs45574833, rs61821114, rs61821041,
251 and rs12064775) previously reported as risk alleles for AMD in the UKBB dataset¹⁸. Moreover, we find
252 that each of these variants predisposes carriers to adverse clinical outcome (i.e. hospitalization) following
253 SARS-CoV-2 infection (OR: 2.13-2.65). A fifth variant, rs2230199, which maps to complement C3, was
254 shown to be linked to AMD in an independent GWAS, however, this variant has not been associated with
255 increased AMD risk in the UK population. The three SNPs that map to C3 each appear to confer some
256 protection associated with SARS-CoV-2 infection (OR: 0.66-0.68). In addition, two of the identified
257 variants (rs61821114 and rs61821041) map to expression quantitative trait loci (eQTL) associated with
258 Complement Decay-Accelerating Factor (CD55)¹⁹. This protein negatively regulates complement
259 activation by accelerating the decay of complement proteins, thereby disrupting the cascade and
260 preventing immune-mediated damage⁷. As reported by GTex Consortium data¹⁹ and highlighted in Figure
261 3d, these eQTLs result in decreased expression of CD55, thereby relieving the restraining function of this
262 protein. In agreement with the functional role of CD55, we observe that these variants are associated with
263 increased risk of adverse clinical outcome associated with SARS-CoV-2 infection (OR: 2.34-2.4).

264

265 Genetic association studies performed on relatively small cohorts can be prone to false positives. While
266 permutation analyses to empirically determine statistical significance thresholds were implemented as
267 described in *Methods*, we also repeated the analysis using updated UKBB data released in May, 2020
268 which included 3,002 patients with suspected infection. Of the 1,073 that tested positive in the updated
269 cohort, 818 required hospitalization (651 and 500 respectively, after ancestry and relatedness filtering, see
270 *Methods*). Importantly, analysis of the May-2020 COVID data recapitulated 6 of 11 April-2020 findings
271 and identified 16 additional loci with study-wide significance ($\alpha = 0.0025$, Supplemental Table S2, Figure
272 3c). Among these, the scan revealed 5 variants proximal to Alpha-2-macroglobulin (A2M), a protease
273 inhibitor and cytokine transporter which participates in the formation of fibrin clots and regulates
274 inflammatory cascades²⁰. Of these, 3 (rs10842898, rs669, and rs4883215) are eQTLs associated with
275 significant downregulation of A2M (and concomitant upregulation of A2M-AS1, the antisense RNA of
276 A2M; data available on gtexportal.org) in multiple tissues including mucosa of the esophagus (P value =
277 1.9×10^{-15}) as highlighted in Figure 3e. In addition to A2M, rs10842898 and rs669 are splicing quantitative
278 trait loci (sQTLs) for Mannose-6-Phosphate Receptor (M6PR) a P-type lectin that regulates lysosomal
279 cargo loading and participates in cellular responses to wound healing, cell growth and viral infection²¹ -
280 suggesting that the SNPs identified may contribute to complex regulation of transcripts with
281 immunological and antiviral roles.

282

283 As detailed in Supplemental Table S2, 936 of the variants that were part of the study are within haplotype
284 blocks of analyzed genes (see *Methods*). Analysis focused on SNPs in complement and coagulation
285 haplotype blocks (based on linkage disequilibrium; LD, See *Methods*) resulted in 16 study-wide
286 significant SNPs ($\alpha = 0.01$, Figure S3) using the April-2020 cohort, of which 8 repeated at study-wide
287 significance ($\alpha = 0.0075$, Figure S3) using the May-2020 dataset. These include rs45574833, a variant
288 highlighted above that results in a missense mutation in C4BPA, a protein that controls activation of the
289 classical complement pathway by mediating hydrolysis of complement factor C4b and degradation of the
290 C3 convertase²² (see Supplemental Table S2). In addition, the haplotype-based analysis identified a link
291 between rs731034 (an eQTL in Collectin Subfamily Member 11; COLEC11) and poor clinical outcome
292 in both April-2020 (OR: 1.27) and May-2020 (OR: 1.33) cohorts. COLEC11, a member of the collectin
293 family of C-type lectins, plays an important role in the innate immune system by binding to carbohydrate
294 antigens (with a preference for fucose and mannose) on microorganisms including viruses, facilitating
295 their recognition and removal. This eQTL variant results in significant upregulation of COLEC11 across
296 multiple tissues including lung (P value = 1×10^{-11}) and suggests that sugar moieties on viral proteins may
297 serve as antigenic targets of immunological responses to SARS-CoV-2 infection. Though experimental
298 validation and functional interrogation of the variants we have identified is required to elucidate their
299 precise pathophysiology, taken together, our observations point to genetic variation in complement and
300 coagulation components as a contributing factor in SARS-CoV-2 mediated disease.

301

302

303 **DISCUSSION**

304 Zoonotic infections like the SARS-CoV-2 pandemic pose tremendous risk to public health and
305 socioeconomic factors on a global scale. While the innate and adaptive arms of the immune system are
306 exquisitely equipped to deal with noxious agents including viruses, interactions between emerging
307 pathogens and their human hosts can manifest in unpredictable ways. In the case of SARS-CoV-2
308 infection a combination of viral replication and immune mediated pathology are the primary drivers of
309 morbidity and mortality. While recent analysis of coronavirus patients in China, suggests that high serum
310 levels of interleukin-6 (IL-6), a proinflammatory cytokine, is associated with poor prognosis¹³ (and as
311 shown in Figure 2, found to be transcriptionally regulated in SARS-CoV-2 patients) further delineation of
312 the regulatory programs that mediate immune pathology associated with SARS-CoV-2 infection is
313 necessary. As illustrated in the accompanying paper and by the results presented herein, knowledge of
314 molecular interactions between virus and host can refine hypothesis-driven discovery of disease
315 determinants.

316

317 Our scan for virus-encoded structural mimics across Earth's virome pointed to molecular mimicry as a
318 pervasive strategy employed by viruses and indicated that the protein structure space used by a given
319 virus is dictated by the host proteome (see accompanying paper). Moreover, observations about how
320 coronaviruses exploit this strategy provided clues about the cellular processes driving pathogenesis.
321 Together with knowledge that CoV infections, including the SARS-CoV outbreak in 2002-2003 and the
322 current SARS-CoV-2 outbreak¹³, result in hyper-coagulative phenotypes²³, our protein structure-function
323 analysis led us to hypothesize that conditions associated with complement or coagulatory dysfunction
324 may influence outcomes of SARS-CoV-2 infections. Of these, among the most common are AMD (which
325 is associated with hyper-activation of the complement pathway) and hyper-coagulative disorders. Their
326 relatively high incidence rates together with SARS-CoV-2 prevalence in and around New York City made
327 them reasonable candidates for a retrospective clinical study.

328
329 As presented above, in addition to rediscovering previously identified risk factors including age, sex,
330 hypertension, and CAD we found that history of macular degeneration or coagulatory dysfunctions
331 predispose patients to poor clinical outcomes (including increased risk of mechanical ventilation and
332 death) following SARS-CoV-2 infection. Complement deficiencies on the other hand, appear to be
333 protective. Their low incidence rates, however, make for a small sample size and invite further
334 investigation. Moreover, retrospective studies of observational data have notable limitations in their data
335 completeness, selection biases, and methods of data capture. As a result, claims on causality cannot be
336 made - nor can we definitively rule out other clinical factors as possible drivers. Nevertheless, in an
337 orthogonal analysis of 650 transcriptional profiles of NP swabs, we demonstrate that in addition to
338 immune factors like Type I interferons and dysregulation of IL-6-dependent inflammatory responses,
339 SARS-CoV-2 infection results in engagement and robust activation of complement and coagulation
340 cascades. Dysregulation of complement and coagulation pathways leading to pathology resulting from
341 viral infection is not without precedent. Indeed, it has been associated with Dengue virus infection where
342 immune mediated pathology and dysregulation of complement is correlated with disease severity and
343 mirrors that of acute SARS-CoV-2 disease²⁴. Moreover, though different from the variants identified in
344 this study, polymorphisms and haplotypes in CFH have been associated with severity of Dengue
345 infection²⁵, suggesting that complement and coagulatory dysfunctions may represent risk factors for a
346 broader range of pathogens.

347
348 Finally, since complement and coagulative dysfunctions can have both acquired and congenital etiologies,
349 we implemented a focused, candidate-driven analysis of UK Biobank data to evaluate linkage between
350 severe SARS-CoV-2 disease and genetic variation associated with complement and coagulation

351 pathways. Our analysis identified putative complement and coagulation associated loci including
352 missense, eQTL and sQTL variants of critical regulators of the complement and coagulation cascades.
353 Though interpretation of these findings may be limited by sample size, site-specific biases in clinical care
354 decisions, ancestral homogeneity and population stratification in the biobank data, and socioeconomic
355 status of affected populations, to our knowledge, this is the first study to identify complement and
356 coagulation functions as underlying risk-factors of SARS-CoV-2 disease outcome. In addition, given an
357 existing menu of immune-modulatory therapies that target complement and coagulation pathways, the
358 discovery provides a rationale to investigate these options for the treatment of SARS-CoV-2 associated
359 pathology. Indeed, the therapeutic potential of complement modulation was recently introduced and
360 further shown to be of significant benefit in a cohort of SARS-CoV-2 patients^{26,27}.

361
362 Our study highlights the value of combining molecular information from virus protein structure-function
363 analysis with orthogonal clinical data analysis to reveal determinants and/or predictors of immunity,
364 susceptibility, and clinical outcome associated with infection. Such a framework can help refine large-
365 scale genomics efforts and help power genomics studies based on informed biological and clinical
366 conjectures. While identification of CoV encoded structural mimics guided the retrospective clinical
367 studies, a molecular and functional link between those observations and our discovery of complement and
368 coagulation functions as risk factors for SARS-CoV-2 pathogenesis remains to be elucidated.
369 Nevertheless, the findings advance our understanding of how SARS-CoV-2 infection leads to disease and
370 can help explain variability in clinical outcomes. Among the implications, the data warrant heightened
371 public health awareness for individuals most vulnerable to developing adverse SARS-CoV-2 mediated
372 pathology.

373
374 **ACKNOWLEDGEMENTS**
375 This work was funded by NIH grants 5R01GM109018 and 5U54CA209997 to SS,
376 R35GM131905 to NPT, F30HL140946 to PT, and equipment grants S10OD012351 and S10OD021764
377 to the Columbia University Department of Systems Biology. CEM would like to thank the Scientific
378 Computing Unit (SCU), XSEDE Supercomputing Resources, the Starr Cancer Consortium (I13-0052),
379 and funding from the WorldQuant Foundation, The Pershing Square Sohn Cancer Research Alliance,
380 NASA (NNX14AH50G, NNX17AB26G), the National Institutes of Health (R21AI129851,
381 R01MH117406, R01AI151059

382
383 **DECLARATION OF INTERESTS**
384 The authors declare no competing interests

385 **FIGURE LEGENDS**

386 **Figure 1|** History of macular degeneration and coagulation disorders are associated with adverse
387 outcomes after confirmed SARS-CoV-2 infection. **a**, Kaplan-Meier curves for 10 binary conditions: age
388 over 65, male sex, macular degeneration (Macula), complement deficiency disorders (CD), coagulation,
389 hypertension, type 2 diabetes (T2DM), obesity, coronary artery disease (CAD), and cough. The survival
390 for the patients with the named condition are shown in orange. The shaded region indicates the 95%
391 confidence interval. The blue survival line is for patients without the named condition. Note that none of
392 the four patients with CD required mechanical ventilation. **b**, Kaplan-Meier curves for the same 10
393 conditions as in **(a)**. All four patients with CD survived (not statistically significant). **c**, Intubation rates
394 across the binary conditions. Mortality (N=88) was highest in patients with a history of macular
395 degeneration, followed by Type 2 Diabetes and Hypertension. **d**, Mortality rates across the binary
396 conditions. Patients with a history of macular degeneration saw the highest mortality rates, followed by
397 Age ≥ 65 and Type 2 Diabetes. **e**, Hazard ratios, estimated using a Cox proportional hazards model, for
398 risk if intubation (as a validated proxy for requiring mechanical respiration). **f**, Similarly, hazard ratios for
399 mortality, estimated using a Cox proportional hazards model. Hazard ratios and statistical significances
400 are shown in Table 1.

401
402 **Figure 2|** SARS-CoV-2 infection engages robust transcriptional regulation of complement and
403 coagulation cascades. **a**, GSEA of HALLMARK gene sets was applied to RNA-seq profiles of NP swabs
404 from 650 control and SARS-CoV-2 infected patients stratified by SARS-CoV-2 positive (green) or low
405 (yellow), medium (orange), high (red) viral load (significantly enriched gene sets highlighted in blue; **b**,
406 Leading edge enrichment plots from GSEA analysis of MsigDB-wide gene sets are shown for
407 HALLMARK_Complement and KEGG_Complement_and_Coagulation_Cascade gene sets with SARS-
408 CoV-2 stratification indicated by color. **c**, Hierarchical clustering of Z-score normalized mRNA profiles
409 of complement and coagulation components that undergo significant (FDR corrected P value < 0.01)
410 transcriptional regulation in response to SARS-CoV-2 infection (cold and hot color scale reflects down,
411 or up regulated expression, respectively). **d-f**, Violin plots (transcripts per million; TPM shown on y-axis)
412 of highlighted differentially regulated genes are shown for upregulated (**d**), downregulated (**e**), or
413 particularly upregulated in the context of high viral load (**f**). Normalized enrichment scores (NES) and
414 FDR-corrected P values are shown.

415
416 **Figure 3|** Targeted genetic association study identifies SNPs in complement and coagulation pathway
417 components associated with clinical outcome of SARS-CoV-2 infection. **a-b**, P values from a Negative
418 Binomial distribution fit to permutation of SNPs sampled (left) and case:control phenotypes (center)

419 generated under the null hypothesis are shown for the April-2020 (**a**) or May-2020 (**b**) cohort (α and
420 distance pairs as indicated; for more information see *Methods*). Also shown are the number of hits that
421 pass the corresponding alpha study-wide significance threshold by distance (right) for April-2020 (**a**) or
422 May-2020 (**b**) cohorts. **c**, Manhattan plots of 2,888 variants within 60kb of complement and coagulation
423 pathway genes for analyses using the April-2020 cohort (top) and May-2020 cohort (bottom). Study-wide
424 significance threshold shown as dashed green lines, nominal significance threshold shown as black
425 dashed line, and SNPs color alternates by chromosome. Significant SNPs are shown as colored markers
426 and annotated with the nearest gene by base-pair distance. SNPs shown in green are study-wide
427 significant in both April-2020 and May-2020. SNPs shown as diamonds are also study-wide significant in
428 haplotype-based analysis (see *Methods*). eQTLs are further highlighted in (**d**) and (**e**). **d**, eQTL
429 relationship for rs61821114 and *CD55* in thyroid¹⁹. The T allele of rs61821114 is associated with
430 significantly lower expression of *CD55*. **e**, eQTL relationship for rs669 and *A2M*¹⁹. The C allele of rs669
431 is associated with significant lower expression of *A2M* in 17 tissues, including the esophageal mucosa
432 (shown) and lung.

433
434 **Figure S1**| Covariate correlations in EHR clinical data. **a**, Spearman correlation between modeled
435 covariates in patients were diagnosed or tested positive for SARS-CoV-2: age, sex, macular degeneration
436 (macula), complement deficiency disorders (CD), coagulation disorders (coagulation), hypertension, Type
437 2 Diabetes, obesity, and coronary artery disease (CAD). **b**, Spearman correlations, as in (**a**), for all
438 patients (includes patients who tested negative for SARS-CoV-2). **c**, Tanimoto coefficients as in (**a**), for
439 patients who tested positive for SARS-CoV-2 infection. Age was binarized as “Age over 65” to compute
440 the score. **d**, Tanimoto coefficients as in (**c**) for all patients.

441
442 **Figure S2**| Results of permutation testing and fits to negative binomial distributions for (**a**) April-2020
443 phenotype permutations, (**b**) April-2020 SNP permutations, (**c**) May-2020 phenotype permutations, (**d**)
444 May-2020 SNP permutations, (**e**) Haplotype SNPs-only April-2020 phenotype permutations, and (**f**)
445 Haplotype SNPs-only May-2020 phenotype permutations. Histograms indicate the number of
446 permutations with X significant hits (black/grey bars). Negative binomial fits are shown in red (see
447 *Methods*). Chi-squared goodness-of-fit tests were performed for each distribution. Distributions which
448 passed the goodness-of-fit test ($p > 0.05$) are shown in black and those that failed ($p \leq 0.05$) are shown in
449 grey. Results are visualized for 5 distances (columns) and 9 alpha thresholds (rows). All fits are available
450 as supplement data.

451

452 **Figure S3** | *P*values from a Negative Binomial distribution fit to permutation of case:control phenotypes
453 generated under the null hypothesis are shown for the Haplotype SNPs-only analyses using the April-
454 2020 (a) or May-2020 (b) cohort. α and distance pairs as indicated; for more information see *Methods*.

455
456 **Figure S4** | Percent of significant eQTLs within a given distance of the gene body. Significant eQTLs
457 were downloaded from the GTEx Portal website for Esophagus, Lung, and Heart tissues (9 tissues total)
458 and used the provided significance thresholds to determine significance. Shown is the percent of
459 significant eQTLs that are within X base pairs of their target gene aggregated over 9 tissues. Over 70% of
460 significant eQTLs are within 60 Kb of their target gene. Black dashed line represents 60 Kb, grey lines
461 represent 40 and 80 Kb.

462
463 **Figure S5** | Comparison of MAF distributions across sampled SNP sets. The medians, means, interquartile
464 range, 95% confidence interval, minimum, and maximum are shown for each of the 100 samples of SNP
465 sets (see *Empirical Permutation Evaluation to set Study-wide Alpha Thresholds* for details). Also shown
466 are the same distribution statistics for the SNP set within 60Kb of complement and coagulation gene
467 bodies (red). Each of the 100 sampled SNP sets MAF distributions were compared to the study SNP set
468 and tested for differences using a two-sample Mann-Whitney U test. Those that were not significantly
469 different ($p > 0.05$) are shown in black. Those that are significantly different ($p \leq 0.05$) are shown in grey
470 and were dropped from the analysis.

471

472 **METHODS**

473

474 *Ethics and Data Governance Approval*

475 The study is approved by the Columbia University Irving Medical Center Institutional Review Board
476 (IRB# AAAL0601) and the requirement for an informed consent was waived. A data request associated
477 with this protocol was submitted to the Tri-Institutional Request Assessment Committee (TRAC) of New-
478 York Presbyterian, Columbia, and Cornell and approved. The research on the UK Biobank data has been
479 conducted using the UK Biobank Resource under Application Number 41039. The transcriptomics
480 analysis samples were collected and processed through the Weill Cornell Medicine Institutional Review
481 Board (IRB) Protocol 19-11021069.

482

483 **Retrospective Clinical Study**

484 *Cohort and Study Description*

485 In this observational cohort study, we used a data warehouse derived from electronic health records
486 (EHRs) from 11,116 patients treated at New York-Presbyterian/Columbia University Irving Medical
487 Center for suspected cases of SARS-CoV-2 infection. For these patients we collected contemporary data
488 from their current encounter (i.e. the encounter associated with their suspected SARS-CoV-2 infection) as
489 well as historical data, if available, from their previous encounters. Contemporary data (data collected
490 between February 1, 2020 and April 12, 2020) included insurance billing information, laboratory
491 measurements, procedures, and SARS-CoV-2 diagnostic test results. These data were derived from the
492 data warehouse tables in Epic. 6,927 patients have historical data (data collected prior to September 24,
493 2019) available from an OMOP v5 instance stored using MySQL, which included all of the standard
494 tables for recording condition, procedure, medication, and measurement data (among others). Of these we
495 used the insurance billing information from the condition occurrence table and demographics from the
496 person table. See *Preparation of data for modeling* for further details on data preparation.

497

498 We used the contemporary data to define inclusion criteria and outcomes (requiring mechanical
499 respiration and mortality) and used historical data to define patient comorbidities. We defined three
500 hypothesized comorbidity covariates, macular degeneration, complement deficiency disorders, and
501 disorders of coagulation. We used historical data to define these comorbidities, age, and sex. We did not
502 include race and ethnicity data in the modeling as we have previously found issues with the data quality²⁸.
503 The race/ethnicity data we do have is included in the tables for reference. We also modeled other
504 comorbidities previously associated with morbidity and mortality (Zhou et al and others), including
505 history of cardiovascular disease, hypertension, obesity, and diabetes (Table 1, Table S1) -- all derived

506 from the historical data. Coded covariate definitions, as well as lists of which diagnosis codes are most
507 common in each group, are available in the supplemental materials and methods. We used established
508 institutional procedures and an institutional clinical data warehouse to extract all data from the EHR.

509

510 *Defining patient outcomes*

511 Outcome definitions were defined by data derived from the electronic health record between February 1,
512 2020 and April 12, 2020. Mortality is derived from a death note filed by a resident or primary provider
513 that records the date and time of death. Intubation was used as an intermediary endpoint and is a proxy for
514 a patient requiring mechanical respiration. We used note types that were developed for patients with
515 SARS-CoV-2 infection to record that this procedure was completed. We validated outcome data derived
516 from notes against the patient's medical record using manual review.

517

518 *Preparation of data for modeling*

519 We used MySQL and python libraries (pymysql, pandas) to extract and prepare the data for modeling.
520 The code for data preparation is available in the github ([https://github.com/tatonetti-](https://github.com/tatonetti-lab/complementcovid)
521 [lab/complementcovid](https://github.com/tatonetti-lab/complementcovid)) as a Jupyter Notebook titled Data Setup. We begin by creating a master list of
522 suspected covid patients. These are patients that are either diagnosed with the disease, as indicated by a
523 ICD10 code for SARS-CoV-2 infection, in their billing data or a patient that was tested for the presence
524 of the virus using RT-PCR as indicated by a "lab" order for the test. We found 2,821 using the former
525 method and 11,116 patients using the latter. We then extracted birthdates, death dates (if the patient had
526 died or a null value otherwise), and sex codes (1 for female, 2 for male). Patients which had sex codes for
527 non-binary genders were excluded from our analysis. We then define a "first diagnosis date" for each
528 patient as either their first diagnosis date (by billing code) or the first date that they tested positive for
529 SARS-CoV-2, whichever comes first. Next, we calculate each patient's age at the time of this "first
530 diagnosis date." Each of the outcomes and covariates are extracted from their respective tables as detailed
531 in the github. Whenever possible, we use the highest-level ancestor code (from the structured vocabulary
532 in OMOP) that represents the concept we want to model. We then use the concept ancestor tables to grab
533 all the descendant codes. Note that diabetic kidney disease was considered for inclusion and so is
534 represented in the data preparation script, however, it was never modeled. Cough is included as a
535 covariate as a reference symptom for comparison. The last step in the preparation process was to compute
536 the censor dates. To do, we iterated through each patient in our master list and computed their time (in
537 days) to intubation (if they required mechanical respiration) or death (if they died). If not, then the study
538 end date (April 25, 2020) was used as the patient's censored time (in days). Finally, for any patients that
539 were not SARS-CoV-2 positive, their time-to-event values were set to a null indicator to be dropped from

540 the dataset later. Finally, the data are all combined in a pandas (version 1.0.3) dataframe and saved to disk
541 as a pickle file for efficient loading.

542

543 *Statistical Model*

544 Our patient timelines may be censored since our study cohort included patients that were being treated at
545 the time of analysis. We performed survival analysis on the intubation orders and death using a Cox
546 proportional-hazards model and visualized the risk using Kaplan-Meier curves using the lifelines python
547 package (version 0.24.4). Error estimates on the Kaplan-Meier curves are estimated using Greenwood's
548 Exponential Formula²⁹. We fit both univariate models and models fit on the covariate, age, and sex and
549 used log-likelihood to assess significance. We reported Cox proportional hazards coefficients and their
550 95% confidence intervals (Table 1). We modeled whether or not a patient had macular degeneration, a
551 complement deficiency disorder, or a coagulation disorder as binary variables (1=yes, 0=no). Code
552 definitions provided in Table S1. We also included other significant comorbidities suggested by previous
553 studies, CAD, hypertension, T2DM, obesity, or smoking status as binary variables (1=yes, 0=no), sex as a
554 binary variable (0=female, 1=male), age as quantitative variable, older age over 65 (note that age over 65
555 is used *only* for illustrative purposes and is not used in multivariate modeling -- in the multivariate model
556 age as a quantitative variable is used), and outcome as a binary variable (1=yes, 0=no). The outcome of
557 interest was coded as 0 until the day it occurred (the date of the first intubation order following admission
558 or the death date) or the date of analysis, whichever occurred first. Survival curves are generated for the
559 indicated variables by setting all other variables to their respected averages within the training data. Note
560 that we dropped patients who experienced the outcome before their initial diagnosis. This is either due to
561 patients being hospitalized prior to infection (in the case of intubation) or errors in the coded data. We
562 dropped 121 patients for intubation prior to infection and 12 patients for prior death. We also restricted
563 the study to 90 days from the start date. One patient was removed for having an event outside of this
564 range.

565

566 *Covariate Correlations*

567 Using the data prepared as discussed above, we computed pairwise statistical correlations between age,
568 sex as well as history of macular degeneration, complement deficiency disorders, coagulation disorders,
569 HTN, T2DM, obesity, and CAD. We computed them using data from all suspected patients (tested both
570 positive and negative) as well as only those patients who tested positive. We used spearman rho and the
571 tanimoto coefficients (1-Jaccard distance) as our measures of correlation. For the comparison using the
572 tanimoto coefficient we binarized age as greater than or equal to 65.

573

574 *Statistical Software*

575 We used Jupyter Notebooks (jupyter-client version 5.3.4 and jupyter-core version 4.6.1) running Python
576 3.7 and all fit models using the python lifelines package (version 0.24.4).

577

578 **Transcriptomic Analysis of NP swabs**

579 *Sample Collection and Processing*

580 Patient specimens were collected with patients' consent at New York Presbyterian Hospital (NYPH) and
581 then processed for RT-PCR as described previously³⁰. Nasopharyngeal (NP) swab specimens were
582 collected using the BD Universal Viral Transport Media system (Becton, Dickinson and Company,
583 Franklin Lakes, NJ) from symptomatic patients.

584

585 *Extraction of Viral RNA and RT-PCR detection*

586 Total viral RNA was extracted from deactivated samples using automated nucleic acid extraction on the
587 QIA Symphony and the DSP Virus/Pathogen Mini Kit (QIAGEN). One step reverse transcription to
588 cDNA and real-time PCR (RT-PCR) amplification of viral targets, E (envelope) and S (spike) genes and
589 internal control, was performed using the Rotor-Gene Q thermocycler (QIAGEN).

590

591 *Human Transcriptome Analysis*

592 RNA-seq reads that mapped unambiguously to the human reference genome via Kraken2 were used to
593 detect transcriptional responses to SARS-CoV-2 infection as described previously³⁰. Briefly, reads were
594 trimmed with TrimGalore, aligned with STAR (v2.6.1d) to the human reference build GRCh38 and the
595 GENCODE v33 transcriptome reference, gene expression was quantified using featureCounts, stringTie
596 and salmon using the nf-core RNAseq pipeline. Sample QC was reported using fastqc, RSeQC, qualimap,
597 dupradar, Preseq and MultiQC. Reads, as reported by featureCounts, were normalized using variance-
598 stabilizing transform (vst) in DESeq2 package in R and DESeq2 was used to call differential expression
599 with either Positive cases vs Negative, or viral load (High/Medium/Low/None) as reported by RT-PCR
600 cycle threshold (Ct) values. Transcript counts (per million) were used to rank genes and perform gene set
601 enrichment analysis (GSEA).

602

603 *Reverse Transcriptase, quantitative real-time PCR (RT-PCR)*

604 The presence of SARS-CoV-2 in clinical samples was determined by RT-PCR. Briefly, primers for the E
605 (envelope) gene (which detects all members of the lineage B of beta-CoVs), and the S (spike) gene
606 (which specifically detect SARS-CoV-2). Samples were annotated using RT-PCR cycle threshold (Ct)
607 value for SARS-CoV-2 primers as follows: Ct \leq 18 were assigned "high viral load"; Ct 18 - 24 were

608 assigned "medium viral load"; and Ct 24 - 40 were assigned "low viral load" stratifications; Ct > 40 was
609 classified as negative (-).

610

611 **Genetic Analysis of UK Biobank**

612 *Data Source*

613 UK Biobank subjects that were of White British descent, in the UK Biobank PCA calculations and
614 therefore without 3rd degree and above relatedness and without aneuploidy, were used in this study,
615 totaling 337,147 subjects (181,032 females and 156,115 males) (Bycroft 2018). Of the nearly 500,000
616 participants, approximately 50,000 subjects were genotyped on the UK BiLEVE Array by Affymetrix
617 while the rest were genotyped using the Applied Biosystems UK Biobank Axiom Array, with over
618 800,000 markers using build GRCh37 (hg19). The arrays share 95% marker coverage. We extracted
619 markers with a minor allele frequency greater than 0.005, INFO score greater than 0.3, and Hardy-
620 Weinberg equilibrium test mid-p value greater than 10⁻¹⁰ using PLINK2³¹. UKBB version 3 Imputation
621 combined the Haplotype Research Consortium with the UK10K haplotype resource using the software
622 IMPUTE4 (UK Biobank White paper). Association analyses were performed using a logistic regression
623 model with additive gene dosage and covariates including age at 2018, sex, first 10 principal components
624 (provided by the UK Biobank), and the genotyping array the sample was carried out on. We determined
625 the alpha threshold for study-wide significance using an empirical permutation analysis (see *Empirical*
626 *Permutation Evaluation to set Alpha Thresholds*). We performed a study-wide association analysis
627 comparing variants for subjects that were SARS-CoV-2 positive and required hospitalization against the
628 entire population of 337,147 subjects

629

630 *Targeted Gene Set Definition*

631 The union of coagulation and complement related gene sets (with immunoglobulin genes removed) that
632 are part of MsigDB was used to define the set of 102 genes used in this study. For each gene, we used the
633 transcriptional start and stop site from the hg19 build of the human genome to define a catchment window
634 of 80kbp. From the 805,426 variants profiled in the UK Biobank genotyping data after quality control and
635 QC filters using PLINK2 (see above), 3,540 variants within the transcribed region of the genes of interest
636 or within 80kbp flanking the transcribed region, 2,888 are within 60kbp, 2,292 are within 40kbp, and 936
637 are located in haplotype blocks with study genes.

638

639 *Empirical Permutation Evaluation to set Study-wide Alpha Thresholds*

640 We used permutation to estimate null distributions of the number of hits expected at 9 alpha thresholds
641 varying from (5x10⁻⁵ to 0.05) and by varying the distance threshold from 40kb to 80kb. As shown

642 previously, 80% of GWAS hits are within 60Kb of the nearest gene³². Further, as shown in Supplemental
643 Figure S4, we empirically determined that the majority of eQTLs (>70%) are within 60kb of gene bodies.
644 We performed two sets of permutation analyses: (i) permuted the initial set of genes on which the
645 included variant loci were chosen and (ii) permuted the case/control labels. We repeated each 100 times
646 and used the resulting data to fit a negative binomial distribution as our estimate of the null. Additionally,
647 we evaluated each of the sampled SNP variant sets from (i) and compared their MAF distribution with the
648 MAF distribution of the Complement and Coagulation set. We removed any sets that were significantly
649 different (nominal p-value < 0.05) according to a Mann-Whitney U test (52 of 100 sets were removed due
650 to this criterion; see Supplemental Figure S5). We found that the negative binomial fit the data the best
651 according to a goodness of fit test (Supplemental Figure S2). We used this distribution to assess statistical
652 significance for each combination of alpha and distance values. The result is two estimates of the
653 significance for each alpha (α), distance (d) pair, $P^{(i)}_{\alpha,d}$ and $P^{(ii)}_{\alpha,d}$, from permutation analyses (i) and (ii)
654 above, respectively. For example:

$$\begin{aligned} 655 & \\ 656 & X^{(i)}_{\alpha,d} \sim \text{NB}(r, p) \\ 657 & P^{(i)}_{\alpha,d} = 1 - \text{CDF}_{\text{NB}(r,p)}(k_{\alpha,d}) \end{aligned}$$

658
659 where $X^{(i)}_{\alpha,d}$ is the number of permutation loci with a p-value under the threshold, α . The parameters r and
660 p of the negative binomial represent the number of successes/failures and the probability of success,
661 respectively. Both r and p are fit using non-linear least squares (the `curve_fit` function in `scipy.optimize`)
662 on $X^{(i)}_{\alpha,d}$, the count data from the permutation analyses for the given α and d. The P is then calculated
663 using the CDF of the fitted negative binomial distribution.

664
665 For the gene set permutation analysis (i.e. (i) above) we evaluated each of the 100 replicates to confirm
666 that the minor allele frequency distribution was statistically indistinguishable from that of the complement
667 and coagulation gene set variants. We did so by performing a Mann-Whitney U test between the two
668 distributions and excluded any replicates that showed a significant difference (nominal p-value < 0.05).
669 52 replicates were excluded because of this requirement (Figure SX). This MAF distribution analysis is
670 not necessary for the case/control permutation analysis (i.e. (ii) above) as the loci are the same in each
671 replicate and it is the case/control labels that are permuted.

672
673 Finally, to set the study-wide alpha for each study we chose the greatest threshold value that was gave a P
674 of 0.05 or less for both permutation analysis method:

675

676
$$\max \alpha \text{ s.t. } P_{\alpha,d}^{(i)} < 0.05 \text{ and } P_{\alpha,d}^{(ii)} < 0.05.$$

677
678 Finally, this entire process was repeated for two cohorts of patients, (a) the initial COVID cohort released
679 by the UK Biobank in April 2020 and (b) the updated COVID cohort released in May 2020. The chosen α
680 for April was 0.001 and the chosen α for May was 0.0025. A data file of all of the distribution fit results
681 and their resulting chi-squared goodness-of-fit statistics is made available in the supplemental materials.

682
683 We also performed this permutation significance estimation for the haplotype-derived SNP sets although
684 the distances for all loci chosen using that method are below the minimum in this analysis of 40Kb so
685 those results are constant with regards to distance (Figures S3a-b). The chosen α for the LD-derived SNP
686 sets is 0.01 and 0.0075 for April and May, respectively.

687
688 *Haplotype block-based selection of SNPs*

689 We identified haplotype blocks based on linkage disequilibrium within the UK Biobank data genotype
690 data of the 337,147 subjects using PLINK1.9, where the lower 90% confidence interval is greater than
691 0.70 and the upper 90% confidence interval is at least 0.98. We identified blocks of interests, and
692 subsequently the variants within those blocks, as those that contain any part of the genes of interest as
693 denoted by the transcriptional start and end sites from the hg19 build of the human genome. From the
694 805,426 variants profiles in the UK Biobank genotype data, we identified 7,281 variants within the genes
695 of interest. After applying additional QC filters using PLINK2, 936 variants remained for analysis.

696
697 *Software*

698 We used PLINK v2.00a2LM 64-bit Intel (26 Aug 2019) to run the genetic association analysis. We used
699 PLINK v1.90b6.10 64-bit (17 Jun 2019) to identify haplotype blocks based on linkage disequilibrium. We
700 used Jupyter Notebooks (jupyter-client version 5.3.4 and jupyter-core version 4.6.1) running Python 3.7,
701 numpy 1.18.1, and scipy 1.4.1 for the permutation analyses.

702 **REFERENCES**

703

- 704 1 Zhang, L. *et al.* Crystal structure of SARS-CoV-2 main protease provides a basis for design
705 of improved alpha-ketoamide inhibitors. *Science* **368**, 409-412,
706 doi:10.1126/science.abb3405 (2020).
- 707 2 Dai, W. *et al.* Structure-based design of antiviral drug candidates targeting the SARS-
708 CoV-2 main protease. *Science*, doi:10.1126/science.abb4489 (2020).
- 709 3 Gordon, D. E. *et al.* A SARS-CoV-2-Human Protein-Protein Interaction Map Reveals Drug
710 Targets and Potential Drug-Repurposing. *bioRxiv*, 2020.2003.2022.002386,
711 doi:10.1101/2020.03.22.002386 (2020).
- 712 4 Chen, G. *et al.* Clinical and immunological features of severe and moderate coronavirus
713 disease 2019. *J Clin Invest*, doi:10.1172/JCI137244 (2020).
- 714 5 Moore, B. J. B. & June, C. H. Cytokine release syndrome in severe COVID-19. *Science*,
715 doi:10.1126/science.abb8925 (2020).
- 716 6 Lasso, G. *et al.* A Structure-Informed Atlas of Human-Virus Interactions. *Cell* **178**, 1526-
717 1541 e1516, doi:10.1016/j.cell.2019.08.005 (2019).
- 718 7 Merle, N. S., Church, S. E., Fremeaux-Bacchi, V. & Roumenina, L. T. Complement System
719 Part I - Molecular Mechanisms of Activation and Regulation. *Front Immunol* **6**, 262,
720 doi:10.3389/fimmu.2015.00262 (2015).
- 721 8 Holers, V. M. Complement and its receptors: new insights into human disease. *Annu Rev*
722 *Immunol* **32**, 433-459, doi:10.1146/annurev-immunol-032713-120154 (2014).
- 723 9 Liszewski, M. K., Java, A., Schramm, E. C. & Atkinson, J. P. Complement Dysregulation
724 and Disease: Insights from Contemporary Genetics. *Annu Rev Pathol* **12**, 25-52,
725 doi:10.1146/annurev-pathol-012615-044145 (2017).
- 726 10 Wu, J. & Sun, X. Complement system and age-related macular degeneration: drugs and
727 challenges. *Drug Des Devel Ther* **13**, 2413-2425, doi:10.2147/DDDT.S206355 (2019).
- 728 11 Ambati, J., Atkinson, J. P. & Gelfand, B. D. Immunology of age-related macular
729 degeneration. *Nat Rev Immunol* **13**, 438-451, doi:10.1038/nri3459 (2013).
- 730 12 Degn, S. E., Jensenius, J. C. & Thiel, S. Disease-causing mutations in genes of the
731 complement system. *Am J Hum Genet* **88**, 689-705, doi:10.1016/j.ajhg.2011.05.011
732 (2011).
- 733 13 Zhou, F. *et al.* Clinical course and risk factors for mortality of adult inpatients with
734 COVID-19 in Wuhan, China: a retrospective cohort study. *Lancet* **395**, 1054-1062,
735 doi:10.1016/S0140-6736(20)30566-3 (2020).
- 736 14 Nicholson-Weller, A. & Wang, C. E. Structure and function of decay accelerating factor
737 CD55. *J Lab Clin Med* **123**, 485-491 (1994).
- 738 15 Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
739 *Nature* **562**, 203-209, doi:10.1038/s41586-018-0579-z (2018).
- 740 16 Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes of a
741 wide range of complex diseases of middle and old age. *PLoS Med* **12**, e1001779,
742 doi:10.1371/journal.pmed.1001779 (2015).
- 743 17 Smith, N. L. *et al.* Genetic predictors of fibrin D-dimer levels in healthy adults. *Circulation*
744 **123**, 1864-1872, doi:10.1161/CIRCULATIONAHA.110.009480 (2011).

- 745 18 Han, X. *et al.* Genome-wide meta-analysis identifies novel loci associated with age-
746 related macular degeneration. *J Hum Genet*, doi:10.1038/s10038-020-0750-x (2020).
- 747 19 Consortium, G. T. *et al.* Genetic effects on gene expression across human tissues. *Nature*
748 **550**, 204-213, doi:10.1038/nature24277 (2017).
- 749 20 Rehman, A. A., Ahsan, H. & Khan, F. H. alpha-2-Macroglobulin: a physiological guardian.
750 *J Cell Physiol* **228**, 1665-1675, doi:10.1002/jcp.24266 (2013).
- 751 21 Gary-Bobo, M., Nirde, P., Jeanjean, A., Morere, A. & Garcia, M. Mannose 6-phosphate
752 receptor targeting and its applications in human diseases. *Curr Med Chem* **14**, 2945-
753 2953, doi:10.2174/092986707782794005 (2007).
- 754 22 Ermert, D. & Blom, A. M. C4b-binding protein: The good, the bad and the deadly. Novel
755 functions of an old friend. *Immunol Lett* **169**, 82-92, doi:10.1016/j.imlet.2015.11.014
756 (2016).
- 757 23 Goeijenbier, M. *et al.* Review: Viral infections and mechanisms of thrombosis and
758 bleeding. *J Med Virol* **84**, 1680-1696, doi:10.1002/jmv.23354 (2012).
- 759 24 Nascimento, E. J. *et al.* Alternative complement pathway deregulation is correlated with
760 dengue severity. *PLoS One* **4**, e6782, doi:10.1371/journal.pone.0006782 (2009).
- 761 25 Pastor, A. F. *et al.* Complement factor H gene (CFH) polymorphisms C-257T, G257A and
762 haplotypes are associated with protection against severe dengue phenotype, possible
763 related with high CFH expression. *Hum Immunol* **74**, 1225-1230,
764 doi:10.1016/j.humimm.2013.05.005 (2013).
- 765 26 Risitano, A. M. *et al.* Complement as a target in COVID-19? *Nat Rev Immunol*,
766 doi:10.1038/s41577-020-0320-7 (2020).
- 767 27 Mastaglio, S. *et al.* The first case of COVID-19 treated with the complement C3 inhibitor
768 AMY-101. *Clin Immunol* **215**, 108450, doi:10.1016/j.clim.2020.108450 (2020).
- 769 28 Polubriaginof, F. C. G. *et al.* Challenges with quality of race and ethnicity data in
770 observational databases. *J Am Med Inform Assoc* **26**, 730-736,
771 doi:10.1093/jamia/ocz113 (2019).
- 772 29 Hosmer, D. W., Lemeshow, S. & May, S. *Applied survival analysis : regression modeling*
773 *of time-to-event data*. 2nd edn, (Wiley-Interscience, 2008).
- 774 30 Butler, D. J. *et al.* Shotgun Transcriptome and Isothermal Profiling of SARS-CoV-2
775 Infection Reveals Unique Host Responses, Viral Diversification, and Drug Interactions.
776 *bioRxiv*, 2020.2004.2020.048066, doi:10.1101/2020.04.20.048066 (2020).
- 777 31 Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer
778 datasets. *Gigascience* **4**, 7, doi:10.1186/s13742-015-0047-8 (2015).
- 779 32 Brodie, A., Azaria, J. R. & Ofran, Y. How far from the SNP may the causative genes be?
780 *Nucleic Acids Res* **44**, 6046-6054, doi:10.1093/nar/gkw500 (2016).
- 781

Figure 1

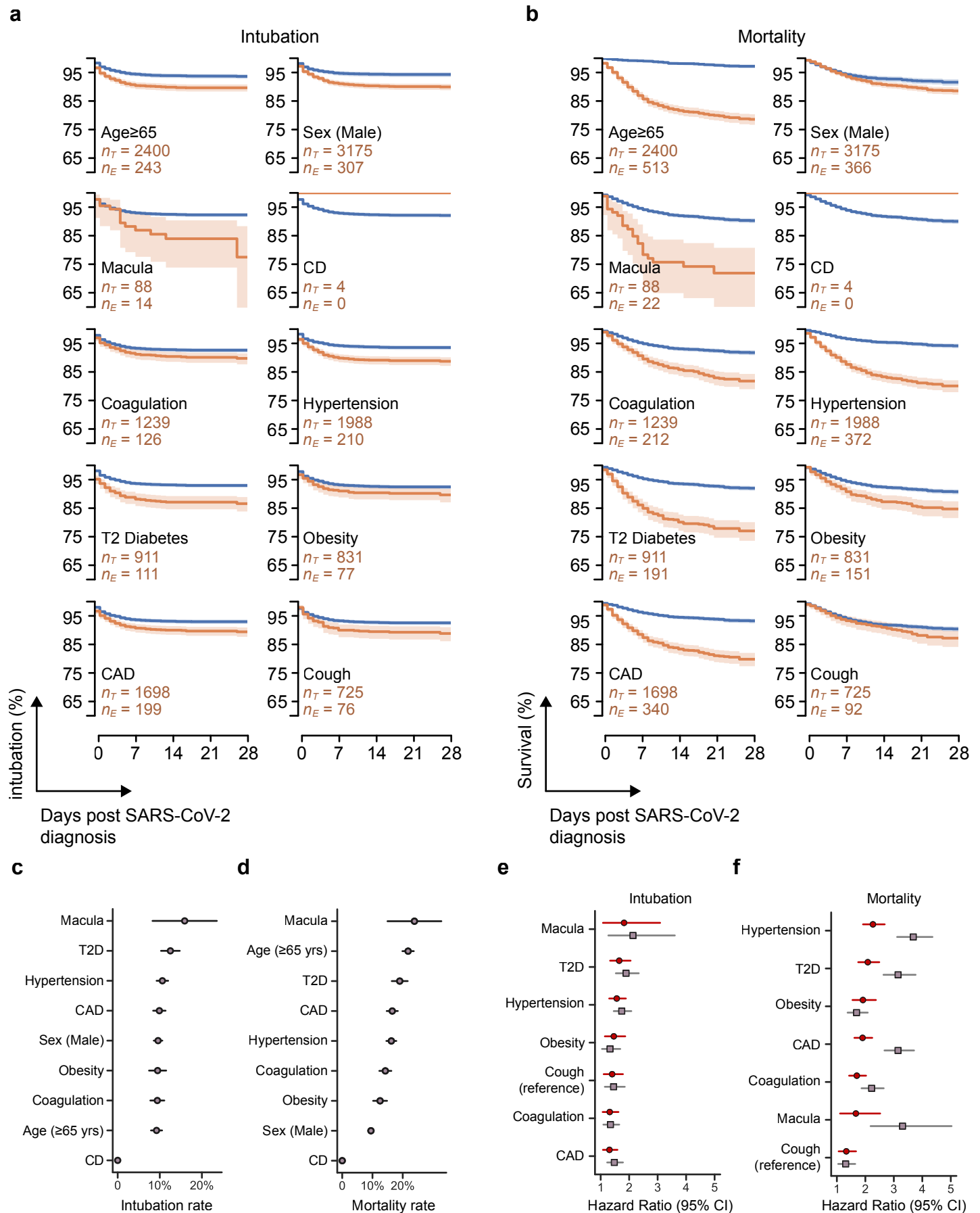


Figure 2

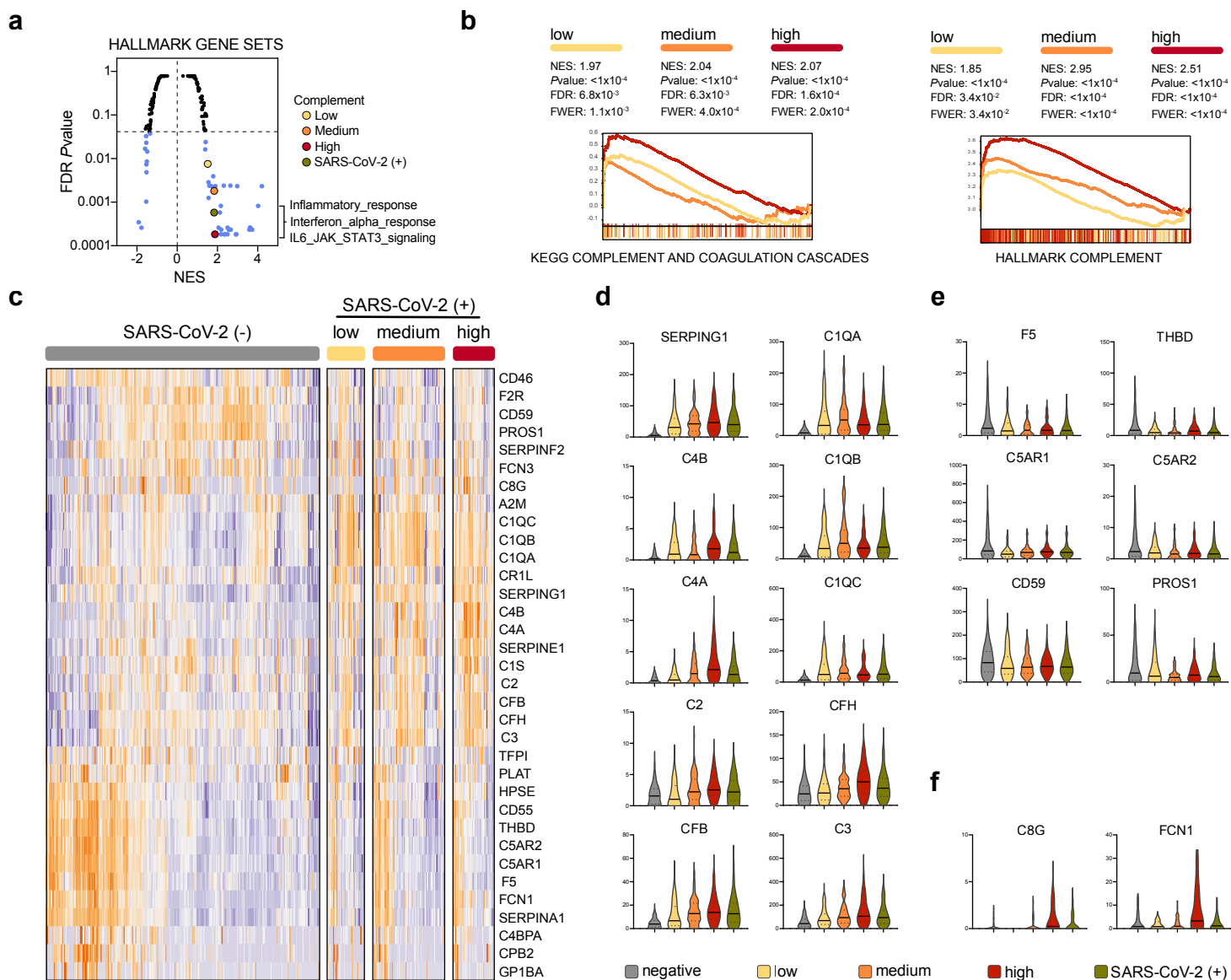


Figure 3

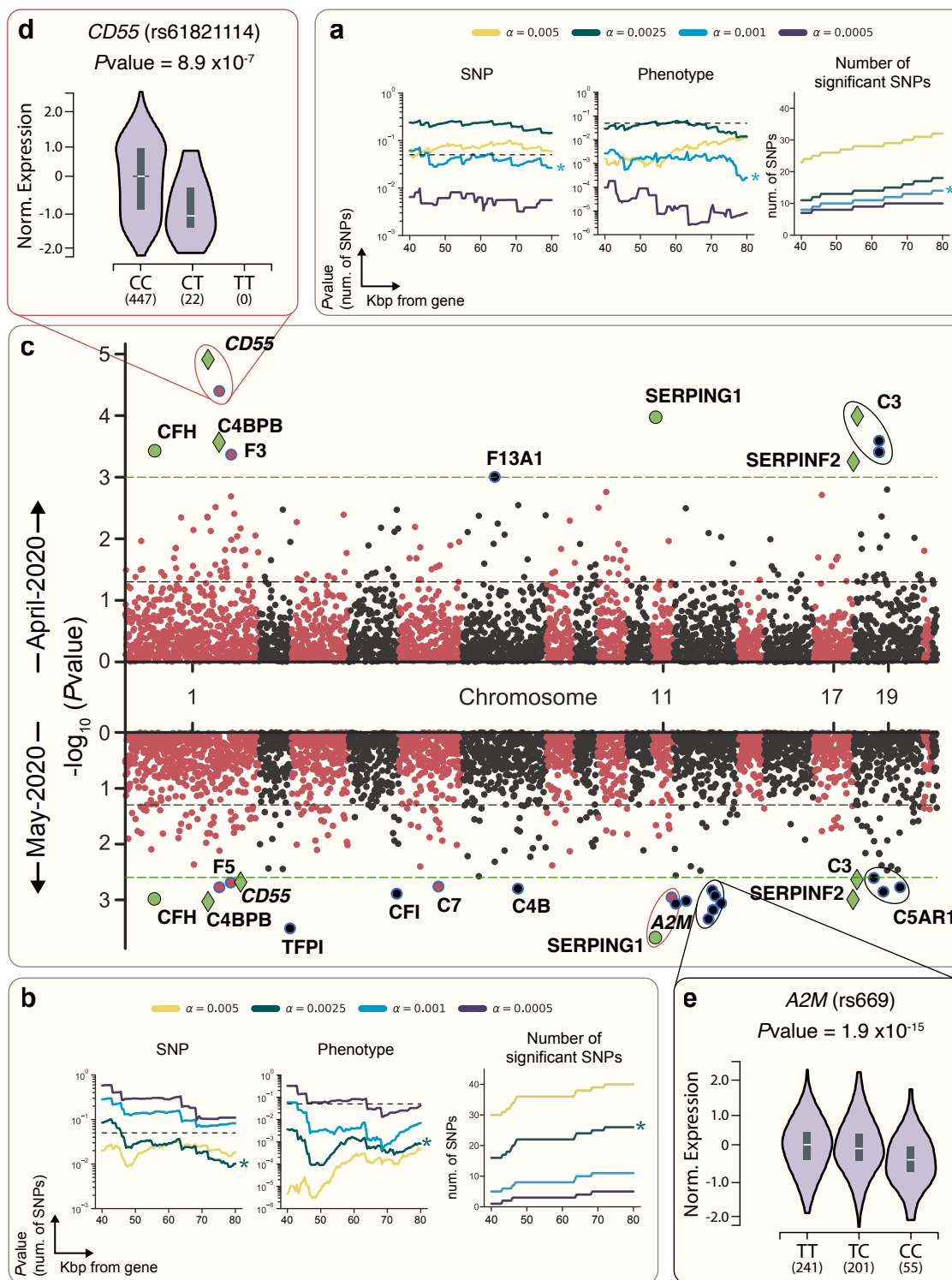
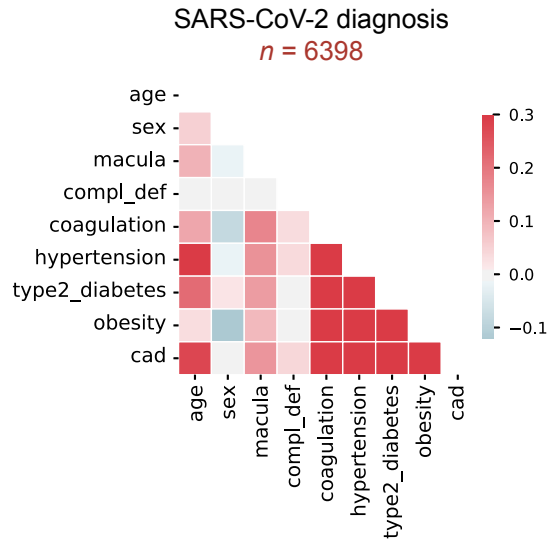
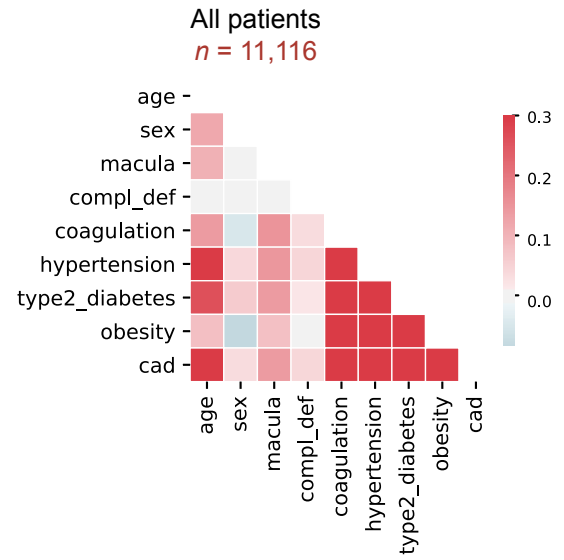


Figure S1

a



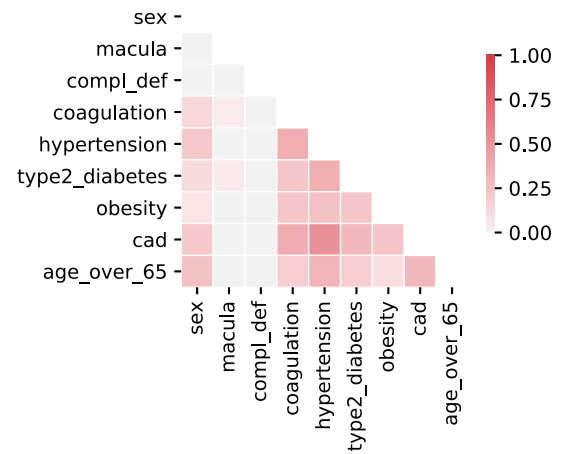
b



c

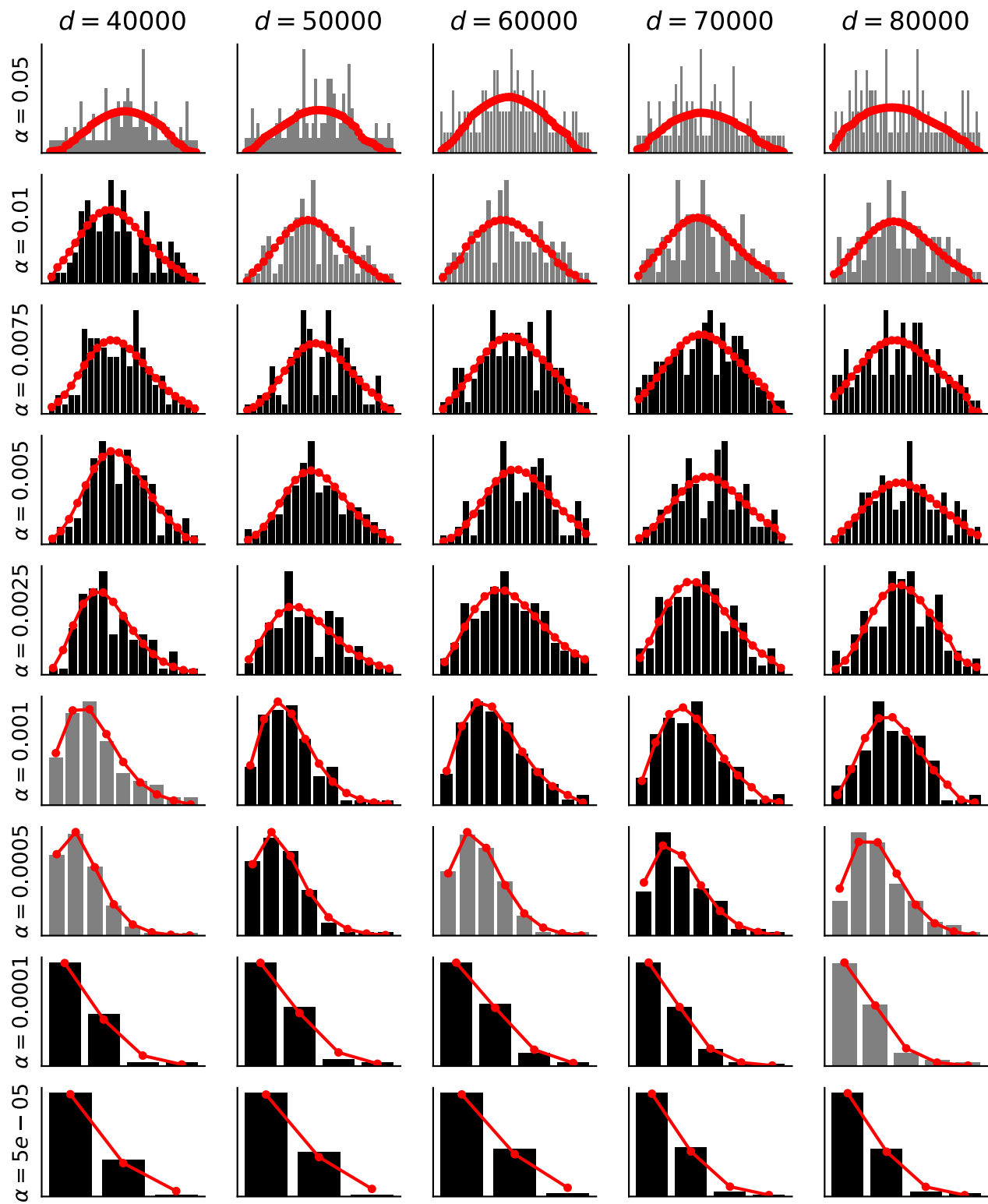


d



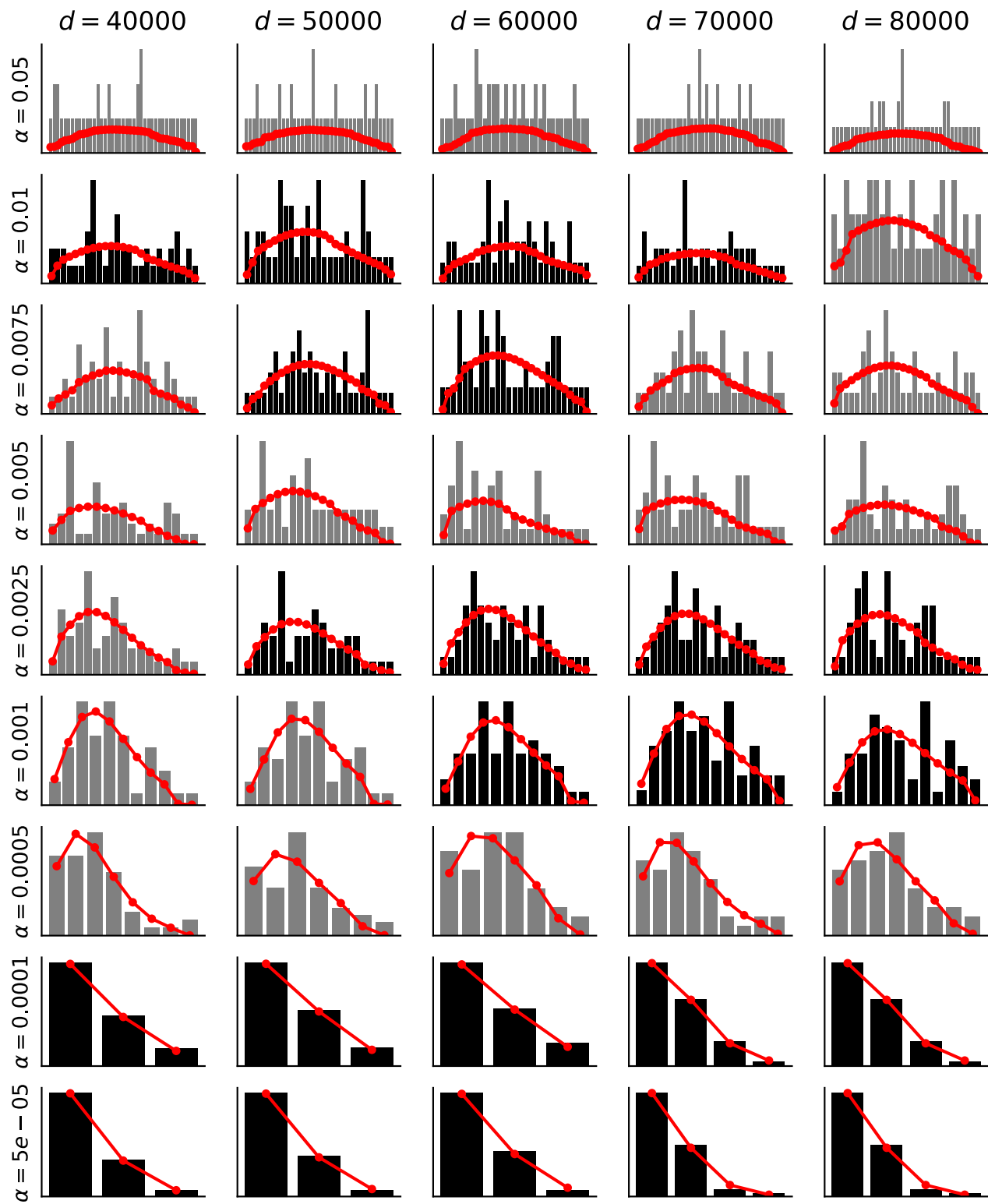
Supplemental Figure S2a

April 2020 Phenotype Permutation



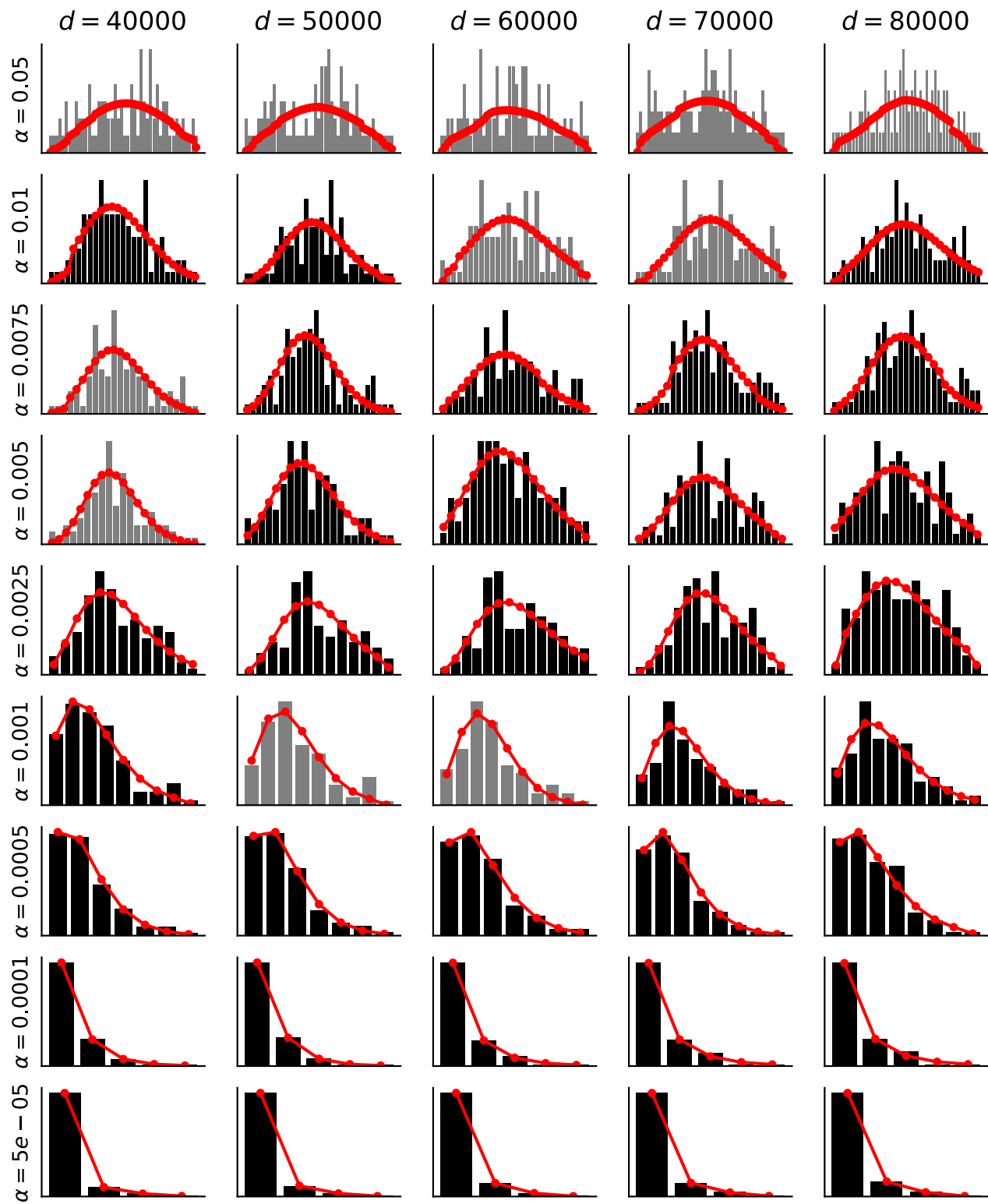
Supplemental Figure S2b

April 2020 SNP Permutation



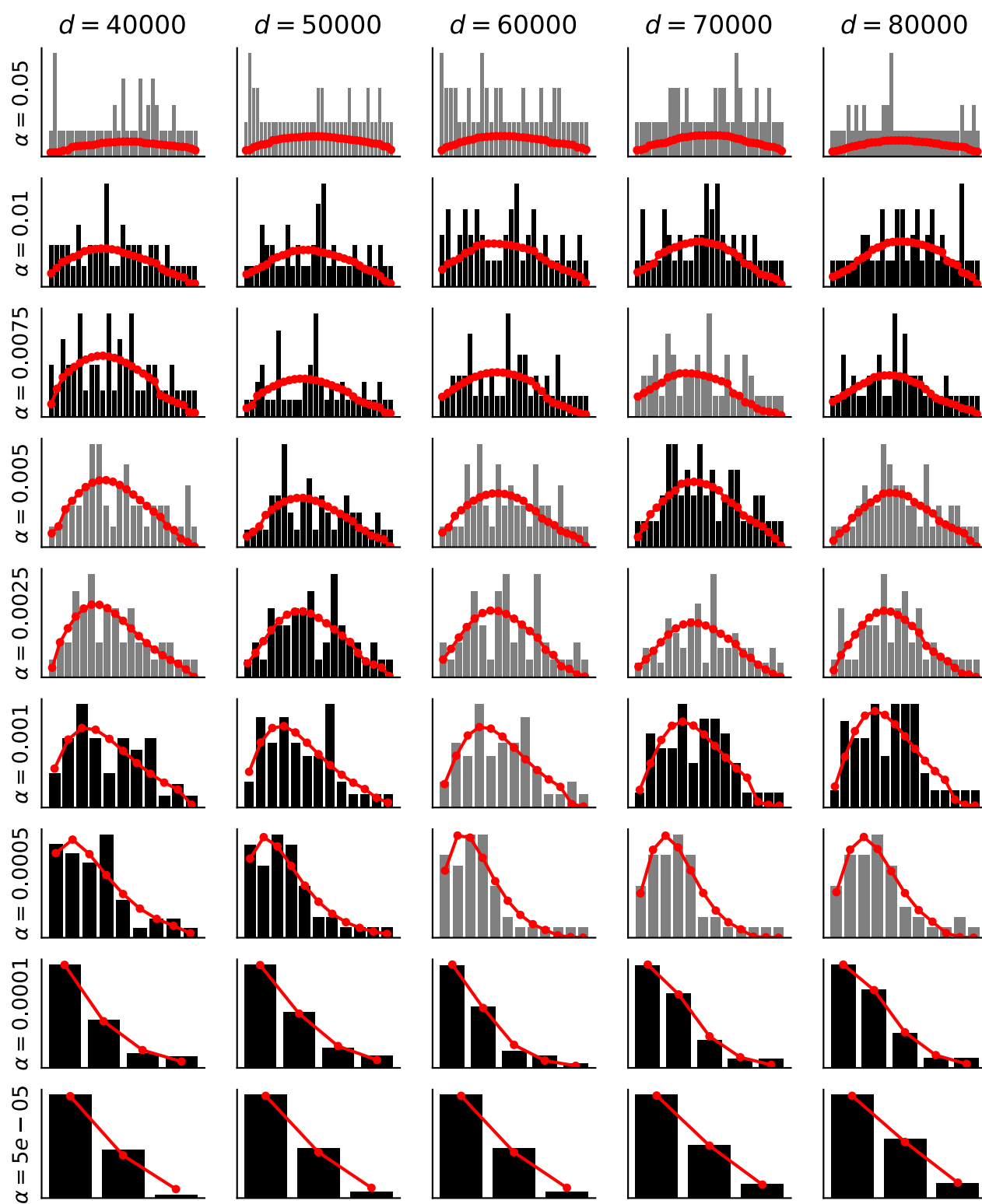
Supplemental Figure S2c

May 2020 Phenotype Permutation



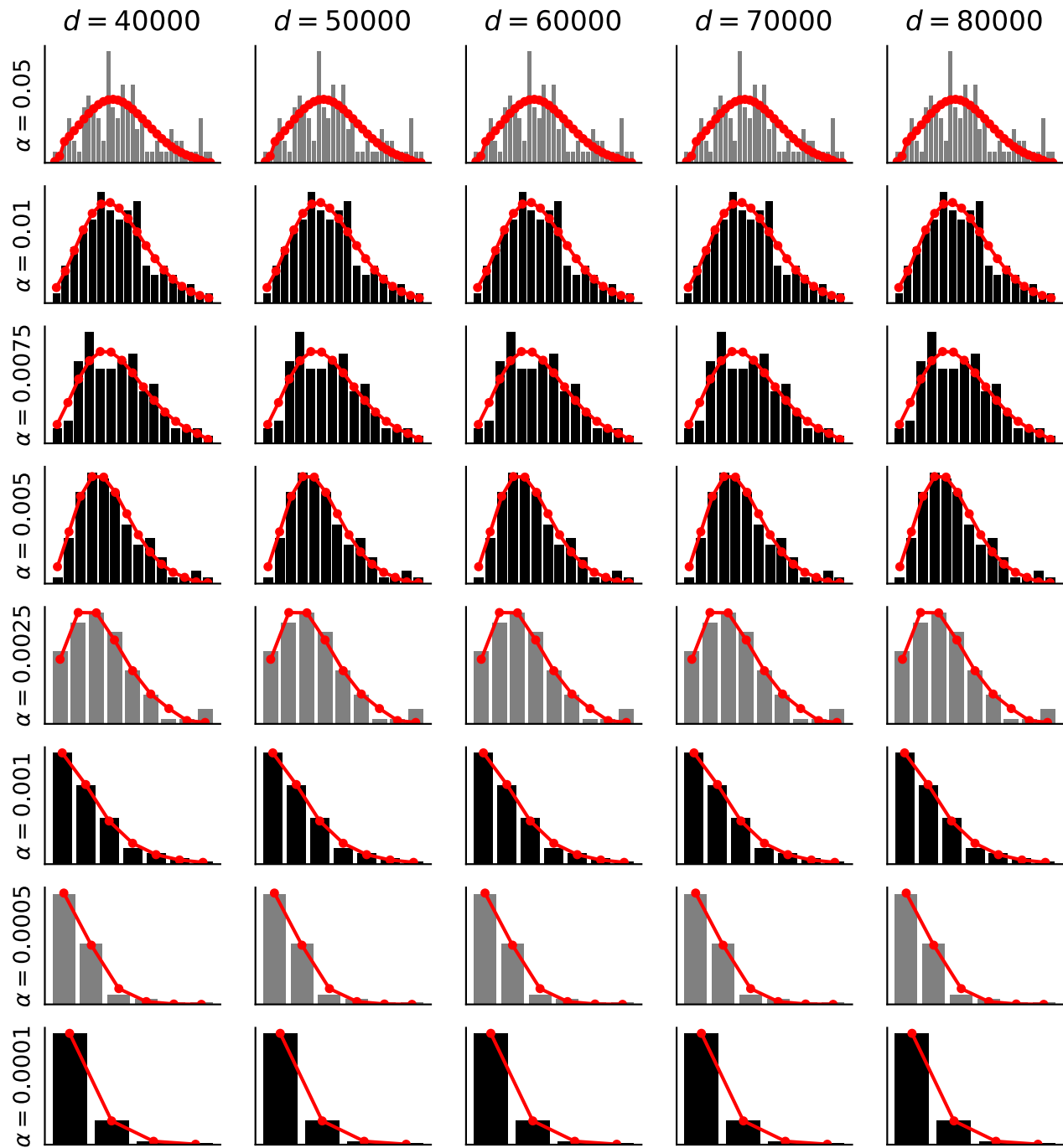
Supplemental Figure S2d

May 2020 SNP Permutation



Supplemental Figure S2e

April 2020 Haplotype Phenotype Permutation



Supplemental Figure S2f

May 2020 Haplotype Phenotype Permutation

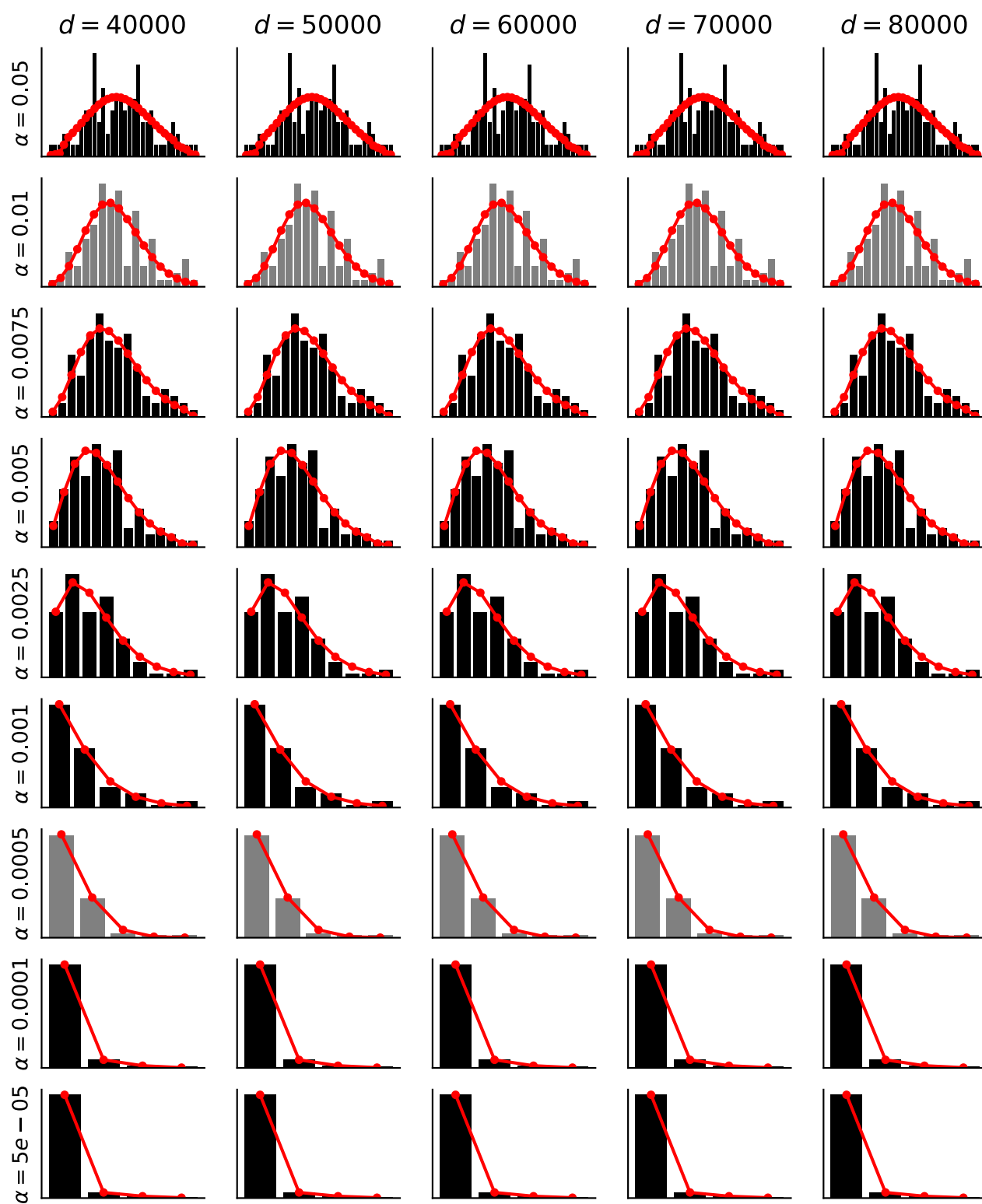
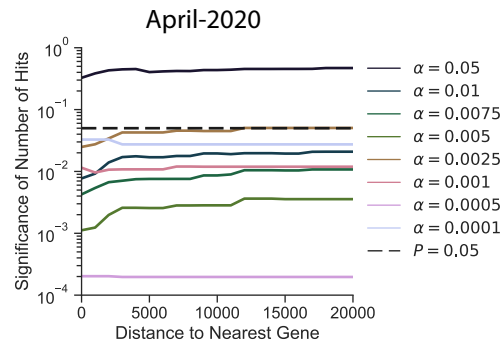


Figure S3

a



b

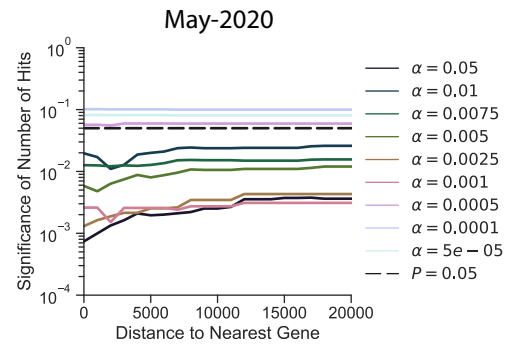


Figure S4

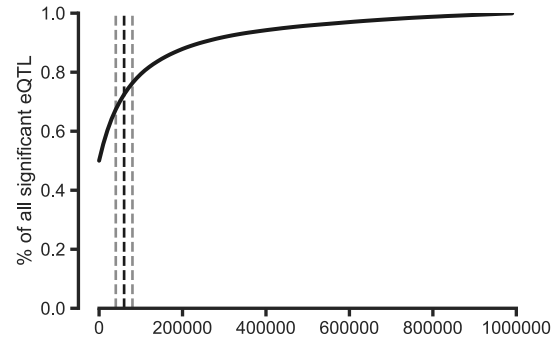


Figure S5

