



## 18 SUMMARY

19 Understanding the pathophysiology of SARS-CoV-2 infection is critical for therapeutics and public  
20 health intervention strategies. Viral-host interactions can guide discovery of regulators of disease  
21 outcomes, and protein structure function analysis points to several immune pathways, including  
22 complement and coagulation, as targets of the coronavirus proteome. To determine if conditions  
23 associated with dysregulation of the complement or coagulation systems impact adverse clinical outcomes  
24 associated with SARS-CoV-2 infection, we performed a retrospective observational study of 11,116  
25 patients suspected of SARS-CoV-2 infection. We found that history of macular degeneration (a proxy for  
26 complement activation disorders) and history of coagulation disorders (thrombocytopenia, thrombosis,  
27 and hemorrhage) are risk factors for morbidity and mortality in SARS-CoV-2 infected patients – effects  
28 that could not be explained by age or sex. In addition, using data from the UK Biobank, we implemented  
29 a candidate driven approach to evaluate linkage between severe SARS-CoV-2 disease and genetic  
30 variation associated with complement and coagulation pathways. Among our findings, our scan identified  
31 an eQTL for CD55 (a negative regulator of complement activation) and SNPs in Complement Factor H  
32 (CFH) and Complement Component 4 Binding Protein Alpha (C4BPA), which play central roles in  
33 complement activation and innate immunity and were previously linked to Age Related Macular  
34 Degeneration (AMD) in a Genome-Wide Association Study (GWAS). In addition to providing evidence  
35 that complement function modulates SARS-CoV-2 infection outcome, the data point to several putative  
36 genetic markers of susceptibility. The results highlight the value of using a multi-modal analytical  
37 approach, combining molecular information from virus protein structure-function analysis with clinical  
38 informatics and genomics to reveal determinants and predictors of immunity, susceptibility, and clinical  
39 outcome associated with infection.

40

## 41 INTRODUCTION

42 The SARS-CoV-2 pandemic has had profound economic, social, and public health impact with over 3  
43 million confirmed cases and over 210,000 deaths across the globe. The infection causes respiratory illness  
44 with symptoms ranging from cough and fever to difficulty breathing. While highly variable age-  
45 dependent mortality rates have been widely reported, the comorbidities that drive this dependence are not  
46 fully understood. Further, with some notable exceptions<sup>1-3</sup>, molecular studies have largely focused on  
47 ACE-2, the receptor and determinant of cell entry and viral replication<sup>3</sup>. While ACE-2 expression is  
48 critical, viruses employ a wide range of molecular strategies to infect cells, avoid detection, and  
49 proliferate. In addition, viral replication and immune mediated pathology are the primary drivers of  
50 morbidity and mortality associated with SARS-CoV-2 infection<sup>4,5</sup>. Therefore, understanding how virus-

51 host interactions manifest as SARS-CoV-2 risk factors will facilitate clinical management, choice of  
52 therapeutic interventions, and setting of appropriate social and public health measures.

53

54 Knowledge of the precise molecular interactions that control viral replicative cycles can delineate  
55 regulatory programs that mediate immune pathology associated with infection and provide valuable clues  
56 about disease determinants. For example, viruses, including SARS-CoV-2, deploy an array of genetically  
57 encoded strategies to co-opt host machinery. Among the strategies, viruses encode multifunctional  
58 proteins that harness or disrupt cellular functions, including nucleic acid metabolism and modulation of  
59 immune responses, through protein-protein interactions and molecular mimicry – structural similarity  
60 between viral and host proteins (for a full discussion please see accompanying paper). Recently, we  
61 employed protein structure modeling to systematically chart interactions across all human infecting  
62 viruses<sup>6</sup> and in an accompanying paper, performed a virome-wide scan for molecular mimics. This  
63 analysis points to broad diversification of strategies deployed by human infecting viruses and identifies  
64 biological processes that underlie human disease. Of particular interest, we mapped over 140 cellular  
65 proteins that are mimicked by coronaviruses (CoV). Among these, we identified components of the  
66 complement and coagulation pathways as targets of structural mimicry across all CoV strains (see  
67 companion paper).

68

69 Through activation of one of three cascades, (i) the classical pathway triggered by an antibody–antigen  
70 complex, (ii) the alternative pathway triggered by binding to a host cell or pathogen surface, and (iii) the  
71 lectin pathway triggered by polysaccharides on microbial surfaces, the complement system is a critical  
72 regulator of host defense against pathogens including viruses<sup>7</sup>. When dysregulated by age-related effects  
73 or excessive acute and chronic tissue damage, complement activation can contribute to pathologies  
74 mediated by inflammation<sup>7,8</sup>. Similarly, inflammation-induced coagulatory programs as well as crosstalk  
75 between pro-inflammatory cytokines and the coagulative and anticoagulant pathways play pivotal roles in  
76 controlling pathogenesis associated with infections. Therefore, while the age-related differences in  
77 susceptibility to SARS-CoV-2 are likely a consequence of multiple underlying variables, virally encoded  
78 structural mimics of complement and coagulation pathway components may contribute to CoV associated  
79 immune mediated pathology. Moreover, a corollary of these observations is that dysfunctions associated  
80 with complement and/or coagulation may impact clinical outcome of SARS-CoV-2 infection. For  
81 example, the companion study suggests that coagulation disorders, such as thrombocytopenia, thrombosis  
82 and hemorrhage, may represent risk factors for SARS-CoV-2 clinical outcome. Among complement-  
83 associated disorders, multiple genetic and experimental evidence (including animal models of disease,  
84 histological examination of affected tissue, and germline mutational analysis) point to dysregulation of

85 the complement system as the major driver of both early-onset, and age-related macular degeneration  
86 (AMD)<sup>9,10</sup>. A hyperinflammatory phenotype mediated by complement leads to progressive immune-  
87 mediated deterioration of the central retina. While AMD, the leading cause of blindness in elderly  
88 individuals (affecting roughly 200 million people worldwide<sup>11</sup>), is likely the result of multiple  
89 pathological processes, dysregulation of complement activation has emerged as the most widely accepted  
90 cause of disease<sup>11-13</sup>.

91  
92 To determine if conditions associated with dysregulation of the complement or coagulation systems  
93 impact adverse clinical outcomes associated with SARS-CoV-2 infection, we conducted a retrospective  
94 observational study of 11,116 patients at New York-Presbyterian/Columbia University Irving Medical  
95 Center. In agreement with previous reports<sup>14</sup>, survival analysis identified significant risk of mechanical  
96 respiration and mortality associated with age and sex, as well as history of hypertension, obesity, and type  
97 2 diabetes (T2D), coronary artery disease (CAD). Moreover, we found that history of macular  
98 degeneration (a proxy for complement activation disorders) and coagulation disorders (thrombocytopenia,  
99 thrombosis, and hemorrhage) were at significantly increased risk of adverse clinical outcomes (including  
100 mechanical respiration and death) following SARS-CoV-2 infection. Importantly, these effects could not  
101 be explained by either age or sex. Conversely, albeit in a small number of individuals, we observed that  
102 no patients with complement deficiency disorders required mechanical respiration or succumbed to their  
103 illness. Finally, in an independent analysis of data from the UK Biobank that focused on variants  
104 associated with the complement and coagulation pathways, we found significant genetic markers in  
105 patients presenting with severe SARS-CoV-2 infection. In particular, we identified variants in CD55 (a  
106 negative regulator of complement activation<sup>15</sup>), CFH and C4BPA, which play central roles in complement  
107 activation and innate immunity, to be associated with adverse clinical outcome. In addition to providing  
108 evidence that complement function modulates SARS-CoV-2 infection, the data point to several putative  
109 genetic markers of susceptibility. The results highlight the value of using a multi-modal analytical  
110 approach, combining molecular information from virus protein structure-function analysis with clinical  
111 informatics and genomics to reveal determinants and predictors of immunity, susceptibility, and clinical  
112 outcome associated with infection.

113

## 114 **RESULTS**

### 115 *Comorbidity statistics and covariances in our retrospective observational clinical cohort*

116 To explore if conditions associated with dysregulation of the complement or coagulation systems impact  
117 adverse clinical outcomes associated with SARS-CoV-2, we conducted a retrospective observational  
118 study of patients treated at New York-Presbyterian/Columbia University Irving Medical Center for

119 suspected infection (Table 1). Electronic health records (EHR) were used to define sex and age as well as  
120 histories of macular degeneration, thrombocytopenia, thrombosis, and hemorrhage, hypertension, type 2  
121 diabetes, coronary artery disease, and obesity (see Methods). As shown in Table 1, of the 11,116 patients  
122 that presented to the hospital between February 1, 2020 and April 25, 2020 with suspected SARS-CoV-2  
123 infection, 6,398 tested positive for the virus. Among these, 88 were patients with history of macular  
124 degeneration, four patients with complement deficiency disorders, and 1,179 patients with disorders  
125 associated with the coagulatory system. In addition, hypertension, coronary artery disease, diabetes,  
126 obesity, and annotated cough were represented by 1,922, 1,566, 847, 791, and 727 patients, respectively  
127 (Table 1). While CAD, hypertension, T2D, obesity, and coagulation disorders represent a group with the  
128 highest covariance, we find lower co-occurrence between these conditions and macular degeneration in  
129 both SARS-CoV-2 positive and negative individuals (Figure S1). Finally, of patients who are put on  
130 mechanical ventilation, we observed a 35% mortality rate, and 31% of deceased patients had been on  
131 mechanical respiration.

132

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

134 We estimated the univariate and age- and sex-corrected risk associated with baseline clinical history of  
135 previously reported SARS-CoV-2 risk factors (including hypertension, obesity, type 2 diabetes, and  
136 coronary artery disease) as well as coagulation and complement disorders using survival analysis and Cox  
137 proportional hazards regression modeling. As shown in Figure 1 and Table 1, we identified significant  
138 risk of mechanical respiration and mortality associated with age and sex, as well as history of  
139 hypertension, obesity, and type 2 diabetes (T2D), coronary artery disease (CAD). Moreover, we found  
140 that history of macular degeneration (a proxy for complement activation disorders) and coagulation  
141 disorders (thrombocytopenia, thrombosis, and hemorrhage) were at significantly increased risk of adverse  
142 clinical outcomes (including mechanical respiration and death) following SARS-CoV-2 infection (Figure  
143 1, Table 1). Specifically, we observed a mechanical respiration rate of 15.9% (95% CI: 8.3-23.6) and a  
144 mortality rate of 25% (95% CI: 16.0-34.0) among patients with a history of macular degeneration, and  
145 rates of 9.4% (95% CI: 7.7-11.1) and 14.7% (95% CI: 12.7-16.7) for mechanical respiration and  
146 mortality, respectively, among patients with coagulation disorders (Table 1). Moreover, as shown in  
147 Figure 1b, patients with a history of macular degeneration appear to succumb to disease more rapidly than  
148 others. Critically, the contribution of age and sex was not sufficient to explain the increased risks  
149 associated with history of macular degeneration (Age/Sex-Corrected mechanical respiration HR=1.8 95%  
150 CI: 1.1-3.2,  $P$ value = 0.024; Age/Sex-Corrected mortality HR=1.7 95% CI: 1.1-2.5,  $P$ value = 0.022).  
151 Conversely, albeit in a small number of individuals, we observed that among patients with complement  
152 deficiency disorders, who are normally at increased risk of complications associated with infections, none

153 required mechanical respiration or succumbed to their illness (Figure 1a and 1b). While we cannot rule  
154 out comorbidities that may be associated with macular degeneration, as shown in Figure S1, the  
155 correlation between macular degeneration and established covariates included in this study is low  
156 (correlation coefficients between 0.09 and 0.15). Together, these data suggest that hyper-active  
157 complement and coagulative states predispose individuals to adverse outcomes associated with SARS-  
158 CoV-2 infection, and that deficiencies in complement components may be protective. Importantly, given  
159 the low incidence rate of deficiencies in either complement or coagulation pathways, further analysis with  
160 larger clinical cohorts is warranted.

161

162 *Genetic variation in complement and coagulation pathway components is associated with adverse SARS-*  
163 *CoV-2 infection outcome*

164 The data highlighted above provide evidence that macular degeneration and coagulation disorders play a  
165 role in SARS-CoV-2 infection outcome. Importantly, macular degeneration and coagulation disorders  
166 have established genetic markers associated with regulators of these functions. However, any genetic  
167 components that may underlie the clinical trends we observed remain hidden due to the retrospective  
168 nature of the study and the lack of available genetic data on these patients. On the other hand, the UK  
169 Biobank, a prospective cohort study with deep genetic, physical, and health data collected on ~500,000  
170 individuals across the United Kingdom<sup>16</sup>, allows for genetic and epidemiological associations to be made.  
171 Among UK Biobank participants, recently released data include SARS-CoV-2-related clinical  
172 information on 1,474 suspected cases, including 669 patients who tested positive and 572 who required  
173 hospitalization. In a candidate driven approach, we leveraged this resource to evaluate if SNPs associated  
174 with components of complement or coagulation pathways are associated with SARS-CoV-2 infection or  
175 hospitalization. Briefly, we focused our analysis on 337,147 (181,032 female) subjects of White British  
176 descent, excluding 3rd degree and above relatedness and without aneuploidy<sup>16</sup>. Applying these  
177 restrictions to the UK Biobank SARS-CoV-2 cohort resulted in 957 patients with suspected infection (388  
178 positive, 332 positive and hospitalized; see Methods).

179

180 Of the 805,426 genetic variants profiled in the UK Biobank, 4,248 are associated with 67 genes with  
181 known roles in regulating complement or coagulation pathways (see Methods). As highlighted in Figure 2  
182 and further delineated in Table 2, we identified 10 loci ( $P$ value =  $3 \times 10^{-6}$ ; see Methods) representing 7  
183 genes with study-wide significance at a minor allele frequency of 0.005 with multiple-hypothesis adjusted  
184 p-values less than 0.05. Among these and proximal to coagulation factor III (F3) is variant rs72729504,  
185 which we find to be associated with increased risk of adverse clinical outcome associated with SARS-  
186 CoV-2 infection (OR: 1.93 95% CI 1.34-2.79). Fibrin fragment D-dimer, one of several peptides

187 produced when cross-linked fibrin is degraded by plasmin, is the most widely used clinical marker of  
188 activated blood coagulation. Among the genetic loci that influence D-dimer levels, GWAS studies have  
189 identified mutations in F3 as having the strongest association<sup>17</sup>. Importantly, increased D-dimer levels  
190 were recently reported to correlate with poor clinical outcome in SARS-CoV-2 infected patients<sup>14</sup>. So,  
191 while the functional role of rs72729504 remains to be elucidated, our observations suggest that this locus  
192 may represent a genetic marker of SARS-CoV-2 susceptibility and outcomes.

193  
194 In addition to the SNP highlighted above, we identified 4 variants (rs45574833, rs61821114, rs61821041,  
195 and rs12064775) previously identified as risk alleles for AMD in the UKBB dataset<sup>18</sup>. Moreover, we find  
196 that each of these variants predisposes carriers to adverse clinical outcome (i.e. hospitalization) following  
197 SARS-CoV-2 infection (OR: 2.13-2.65; see Table 3 for variant specific 95% CI). A fifth variant,  
198 rs2230199, which maps to complement C3, was shown to be linked to AMD in an independent GWAS,  
199 however, this variant has not been associated with increased AMD risk in the UK population. The three  
200 SNPs that map to C3 each appear to confer some protection associated with SARS-CoV-2 infection (OR:  
201 0.66-0.68 see Table 3 for variant specific 95% CI). In addition, two of the identified variants (rs61821114  
202 and rs61821041) map to expression quantitative trait loci (eQTL) associated with Complement Decay-  
203 Accelerating Factor (CD55). This protein negatively regulates complement activation by accelerating the  
204 decay of complement proteins, thereby disrupting the cascade and preventing immune-mediated damage<sup>7</sup>.  
205 As shown in Figure 2b, these eQTLs result in decreased expression of CD55, thereby relieving the  
206 restraining function of this protein. In agreement with the functional role of CD55, we observe that these  
207 variants are associated with increased risk of adverse clinical outcome associated with SARS-CoV-2  
208 infection (OR: 2.34-2.4 see Table 3 for variant specific 95% CI). Together, our observations point to  
209 genetic variation in complement and coagulation components as a contributing factor in SARS-CoV-2  
210 mediated disease.

## 211 212 **DISCUSSION**

213 Zoonotic infections like the SARS-CoV-2 pandemic pose tremendous risk to public health and  
214 socioeconomic factors on a global scale. While the innate and adaptive arms of the immune system are  
215 exquisitely equipped to deal with noxious agents including viruses, interactions between emerging  
216 pathogens and their human hosts can manifest in unpredictable ways. In the case of SARS-CoV-2  
217 infection a combination of viral replication and immune mediated pathology are the primary drivers of  
218 morbidity and mortality. While recent analysis of coronavirus patients in China, suggests that high serum  
219 levels of interleukin-6 (IL-6), a proinflammatory cytokine, is associated with poor prognosis<sup>14</sup>, further  
220 delineation of the regulatory programs that mediate immune pathology associated with SARS-CoV-2

221 infection is necessary. As illustrated in the accompanying paper and by the results presented herein,  
222 knowledge of molecular interactions between virus and host can refine hypothesis-driven discovery of  
223 disease determinants.

224

225 Our scan for virus-encoded structural mimics across Earth's virome pointed to molecular mimicry as a  
226 pervasive strategy employed by viruses and indicated that the protein structure space used by a given  
227 virus is dictated by the host proteome (see accompanying paper). Moreover, observations about how  
228 coronaviruses exploit this strategy provided clues about the cellular processes driving pathogenesis.  
229 Together with knowledge that CoV infections, including the SARS-CoV outbreak in 2002-2003 and the  
230 current SARS-CoV-2 outbreak<sup>14</sup>, result in hyper-coagulative phenotypes<sup>19</sup>, our protein structure-function  
231 analysis led us to hypothesize that conditions associated with complement or coagulatory dysfunction  
232 may influence outcomes of SARS-CoV-2 infections. Of these, among the most common are AMD (which  
233 is associated with hyper-activation of the complement pathway) and hyper-coagulative disorders. Their  
234 relatively high incidence rates together with SARS-CoV-2 prevalence in and around New York City made  
235 them reasonable candidates for a retrospective clinical study.

236

237 As presented above, in addition to rediscovering previously identified risk factors including age, sex,  
238 hypertension, and CAD we found that history of macular degeneration or coagulatory dysfunctions  
239 predispose patients to poor clinical outcomes (including increased risk of mechanical ventilation and  
240 death) following SARS-CoV-2 infection. Complement deficiencies on the other hand, appear to be  
241 protective. Their low incidence rates, however, make for a small sample size and invite further  
242 investigation. Further, retrospective studies of observational data have notable limitations in their data  
243 completeness, selection biases, and methods of data capture. As a result, claims on causality cannot be  
244 made - nor can we definitively rule out other clinical factors as possible drivers. Recognizing these  
245 limitations and that AMD and coagulative dysfunctions can have acquired and congenital etiologies, we  
246 implemented a focused, candidate-driven analysis of UK Biobank data to evaluate linkage between severe  
247 SARS-CoV-2 disease and genetic variation associated with complement and coagulation pathways. Our  
248 analysis identified 10 complement and coagulation associated loci including 4 that have been associated  
249 with AMD and 2 eQTLs that negatively impact expression of CD55, a critical negative regulator of the  
250 complement cascade. Though interpretation of our results may be limited by sample size, site-specific  
251 biases in clinical care decisions, ancestral homogeneity in the biobank data, and socioeconomic status of  
252 affected populations, to our knowledge, this is the first study to identify complement and coagulation  
253 functions as an underlying risk-factors of SARS-CoV-2 disease outcome. In addition, given an existing

254 menu of immune-modulatory therapies that target complement and coagulation pathways, the discovery  
255 provides a rationale to investigate these options for the treatment of SARS-CoV-2 associated pathology.

256

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

268

#### 269 **ACKNOWLEDGEMENTS**

270 This work was funded by NIH grants 5R01GM109018 and 5U54CA209997 to SS,  
271 R35GM131905 to NPT, F30HL140946 to PT, and equipment grants S10OD012351 and S10OD021764  
272 to the Columbia University Department of Systems Biology.

273

#### 274 **DECLARATION OF INTERESTS**

275 The authors declare no competing interests

276

277

278 **FIGURE LEGENDS**

279 **Figure 1|** History of macular degeneration and coagulation disorders are associated with adverse  
280 outcomes after confirmed SARS-CoV-2 infection. **a**, Kaplan-Meier curves for 10 binary conditions: age  
281 over 70, male sex, macular degeneration (Macula), complement deficiency disorders (CD), coagulation,  
282 hypertension, type 2 diabetes (T2DM), obesity, coronary artery disease (CAD), and cough. The survival  
283 for the patients with the named condition are shown in orange. The shaded region indicates the 95%  
284 confidence interval. The blue survival line is for patients without the named condition. Note that none of  
285 the four patients with CD required mechanical ventilation. **b**, Kaplan-Meier curves for the same 10  
286 conditions as in **(a)**. All four patients with CD survived (not statistically significant). **c**, Intubation rates  
287 across the binary conditions. Mortality (N=88) was highest in patients with a history of macular  
288 degeneration, followed by Type 2 Diabetes and Hypertension. **d**, Mortality rates across the binary  
289 conditions. Patients with a history of macular degeneration saw the highest mortality rates, followed by  
290 Age  $\geq 65$  and Type 2 Diabetes. **e**, Hazard ratios, estimated using a Cox proportional hazards model, for  
291 risk if intubation (as a validated proxy for requiring mechanical respiration). **f**, Similarly, hazard ratios for  
292 mortality, estimated using a Cox proportional hazards model. Hazard ratios and statistical significances  
293 are shown in Table 1.

294  
295 **Figure 2|** Genetic association study of 332 SARS-CoV-2 infected patients who required hospitalization in  
296 the UK Biobank. Shown is a Manhattan plot for 2,097 single nucleotide polymorphisms (SNPs)  
297 associated with the complement and coagulation pathway genes. Study-wide significance was determined  
298 using a Benjamini-Hochberg adjusted  $P$ value  $< 0.05$ . The 10 study-wide significant loci are annotated  
299 with related complement genes, full details in Table 2. The panels above show expression quantitative  
300 trait relationships between two study-wide significant SNPs and CD55. The minor allele is associated  
301 with significantly reduced CD55 expression in omentum and thyroid.

302  
303 **Figure S1|** Covariate correlation in clinical data. **a**, Spearman correlation between modeled covariates in  
304 patients were diagnosed or tested positive for SARS-CoV-2: age, sex, macular degeneration (macula),  
305 complement deficiency disorders (CD), coagulation disorders (coagulation), hypertension, Type 2  
306 Diabetes, obesity, and coronary artery disease (CAD). **b**, Spearman correlations, as in **(a)**, for all patients  
307 (includes patients who tested negative for SARS-CoV-2).

308  
309  
310  
311

## 312 **METHODS**

### 313 Retrospective Clinical Study

#### 314 *Cohort and Study Description*

315 In this observational cohort study, we used a data warehouse derived from electronic health records  
316 (EHRs) from 11,116 patients treated at New York-Presbyterian/Columbia University Irving Medical  
317 Center for suspected cases of SARS-CoV-2 infection. For these patients we collected contemporary data  
318 from their current encounter (i.e. the encounter associated with their suspected SARS-CoV-2 infection) as  
319 well as historical data, if available, from their previous encounters. Contemporary data (data collected  
320 between February 1, 2020 and April 12, 2020) included insurance billing information, laboratory  
321 measurements, procedures, and SARS-CoV-2 diagnostic test results. These data were derived from the  
322 data warehouse tables in Epic. 6,927 patients have historical data (data collected prior to September 24,  
323 2019) available from an OMOP v5 instance stored using MySQL, which included all of the standard  
324 tables for recording condition, procedure, medication, and measurement data (among others). Of these we  
325 used the insurance billing information from the condition occurrence table and demographics from the  
326 person table. See *Preparation of data for modeling* for further details on data preparation.

327  
328 We used the contemporary data to define inclusion criteria and outcomes (requiring mechanical  
329 respiration and mortality) and used historical data to define patient comorbidities. We defined three  
330 hypothesized comorbidity covariates, macular degeneration, complement deficiency disorders, and  
331 disorders of coagulation. We used historical data to define these comorbidities, age, and sex. We did not  
332 include race and ethnicity data in the modeling as we have previously found issues with the data quality<sup>20</sup>.  
333 The race/ethnicity data we do have is included in the tables for reference. We also modeled other  
334 comorbidities previously associated with morbidity and mortality (Zhou et al and others), including  
335 history of cardiovascular disease, hypertension, obesity, and diabetes (Table 1, Table S1) -- all derived  
336 from the historical data. Coded covariate definitions, as well as lists of which diagnosis codes are most  
337 common in each group, are available in the supplemental materials and methods. We used established  
338 institutional procedures and an institutional clinical data warehouse to extract all data from the EHR.

#### 339 340 *Defining patient outcomes*

341 Outcome definitions were defined by data derived from the electronic health record between February 1,  
342 2020 and April 12, 2020. Mortality is derived from a death note filed by a resident or primary provider  
343 that records the date and time of death. Intubation was used as an intermediary endpoint and is a proxy for  
344 a patient requiring mechanical respiration. We used note types that were developed for patients with

345 SARS-CoV-2 infection to record that this procedure was completed. We validated outcome data derived  
346 from notes against the patient’s medical record using manual review.

347

#### 348 *Ethics and Data Governance Approval*

349 The study is approved by the Columbia University Irving Medical Center Institutional Review Board  
350 (IRB# AAAL0601) and the requirement for an informed consent was waived. A data request associated  
351 with this protocol was submitted to the Tri-Institutional Request Assessment Committee (TRAC) of New-  
352 York Presbyterian, Columbia, and Cornell and approved. The research on the UK Biobank data has been  
353 conducted using the UK Biobank Resource under Application Number 41039.

354

#### 355 *Preparation of data for modeling*

356 We used MySQL and python libraries (pymysql, pandas) to extract and prepare the data for modeling.  
357 The code for data preparation is available in the github ([https://github.com/tatonetti-](https://github.com/tatonetti-lab/complementcovid)  
358 [lab/complementcovid](https://github.com/tatonetti-lab/complementcovid)) as a Jupyter Notebook titled Data Setup. We begin by creating a master list of  
359 suspected covid patients. These are patients that are either diagnosed with the disease, as indicated by a  
360 ICD10 code for SARS-CoV-2 infection, in their billing data or a patient that was tested for the presence  
361 of the virus using RT-PCR as indicated by a “lab” order for the test. We found 2,821 using the former  
362 method and 11,116 patients using the latter. We then extracted birthdates, death dates (if the patient had  
363 died or a null value otherwise), and sex codes (1 for female, 2 for male). Patients which had sex codes for  
364 non-binary genders were excluded from our analysis. We then define a “first diagnosis date” for each  
365 patient as either their first diagnosis date (by billing code) or the first date that they tested positive for  
366 SARS-CoV-2, whichever comes first. Next, we calculate each patient’s age at the time of this “first  
367 diagnosis date.” Each of the outcomes and covariates are extracted from their respective tables as detailed  
368 in the github. Whenever possible, we use the highest-level ancestor code (from the structured vocabulary  
369 in OMOP) that represents the concept we want to model. We then use the concept ancestor tables to grab  
370 all the descendant codes. Note that diabetic kidney disease was considered for inclusion and so is  
371 represented in the data preparation script, however, it was never modeled. Cough is included as a  
372 covariate as a reference symptom for comparison. The last step in the preparation process was to compute  
373 the censor dates. To do, we iterated through each patient in our master list and computed their time (in  
374 days) to intubation (if they required mechanical respiration) or death (if they died). If not, then the study  
375 end date (April 25, 2020) was used as the patient’s censored time (in days). Finally, for any patients that  
376 were not SARS-CoV-2 positive, their time-to-event values were set to a null indicator to be dropped from  
377 the dataset later. Finally, the data are all combined in a pandas (version 1.0.3) dataframe and saved to disk  
378 as a pickle file for efficient loading.

379

### 380 *Statistical Model*

381 Our patient timelines may be censored since our study cohort included patients that were being treated at  
382 the time of analysis. We performed survival analysis on the intubation orders and death using a Cox  
383 proportional-hazards model and visualized the risk using Kaplan-Meier curves using the lifelines python  
384 package (version 0.24.4). Error estimates on the Kaplan-Meier curves are estimated using Greenwood's  
385 Exponential Formula<sup>21</sup>. We fit both univariate models and models fit on the covariate, age, and sex and  
386 used log-likelihood to assess significance. We reported Cox proportional hazards coefficients and their  
387 95% confidence intervals (Table 1). We modeled whether or not a patient had macular degeneration, a  
388 complement deficiency disorder, or a coagulation disorder as binary variables (1=yes, 0=no). Code  
389 definitions provided in Table S1. We also included other significant comorbidities suggested by previous  
390 studies, CAD, hypertension, T2DM, or obesity as binary variables (1=yes, 0=no), sex as a binary variable  
391 (0=female, 1=male), age as quantitative variable, older age (65+), and outcome as a binary variable  
392 (1=yes, 0=no). The outcome of interest was coded as 0 until the day it occurred (the date of the first  
393 intubation order following admission or the death date) or the date of analysis, whichever occurred first.  
394 Survival curves are generated for the indicated variables by setting all other variables to their respected  
395 averages within the training data. Note that we dropped patients who experienced the outcome before  
396 their initial diagnosis. This is either due to patients being hospitalized prior to infection (in the case of  
397 intubation) or errors in the coded data. We dropped 121 patients for intubation prior to infection and 12  
398 patients for prior death. We also restricted the study to 90 days from the start date. One patient was  
399 removed for having an event outside of this range.

400

### 401 *Covariate Correlations*

402 Using the data prepared as discussed above, we computed pairwise statistical correlations between age,  
403 sex as well as history of macular degeneration, complement deficiency disorders, coagulation disorders,  
404 HTN, T2DM, obesity, and CAD. We computed them using data from all suspected patients (tested both  
405 positive and negative) as well as only those patients who tested positive. We chose spearman rho as our  
406 measure of correlation.

407

### 408 *Statistical Software*

409 Models were generated each day that data was available beginning on March 23rd, 2019 with data from  
410 patients available through that day. We used Jupyter Notebooks (jupyter-client version 5.3.4 and jupyter-  
411 core version 4.6.1) running Python 3.7 and all fit models using the python lifelines package (version  
412 0.24.4).

413

## 414 UK Biobank Genetic Analysis

### 415 *Data Source*

416 UK Biobank subjects that were of White British descent, in the UK Biobank PCA calculations and  
417 therefore without 3rd degree and above relatedness and without aneuploidy, were used in this study,  
418 totaling 337,147 subjects (181,032 females and 156,115 males) (Bycroft 2018). Of the nearly 500,000  
419 participants, approximately 50,000 subjects were genotyped on the UK BiLEVE Array by Affymetrix  
420 while the rest were genotyped using the Applied Biosystems UK Biobank Axiom Array, with over  
421 800,000 markers using build GRCh37 (hg19). The arrays share 95% marker coverage. We extracted  
422 markers with a minor allele frequency greater than 0.005, INFO score greater than 0.3, and Hardy-  
423 Weinberg equilibrium test mid-p value greater than 10<sup>-10</sup> using PLINK2<sup>22</sup>. UKBB version 3 Imputation  
424 combined the Haplotype Research Consortium with the UK10K haplotype resource using the software  
425 IMPUTE4 (UK Biobank White paper). Association analyses were performed using a logistic regression  
426 model with additive gene dosage and covariates including age at 2018, sex, first 10 principal components  
427 (provided by the UK Biobank), and the genotyping array the sample was carried out on. We adjusted for  
428 multiple testing with FDR-corrected p-values using the Benjamini-Hochberg method.

429

### 430 *Genetic Association Studies*

431 We performed three study-wide association analyses: (i) comparing variants for SARS-CoV-2 positive  
432 patients against the entire population of 337,147 subjects, (ii) comparing positive patients who required  
433 hospitalization against the entire population, and (iii) comparing patients who tested negative versus the  
434 entire population.

435

### 436 *Targeted Gene Set Definition*

437 We identified a set of 69 genes relating to the complement and coagulation cascades from the Kyoto  
438 Encyclopedia of Genes and Genomes (KEGG accession id: hsa04610). For each gene, we used the  
439 transcriptional start and stop site from the hg19 build of the human genome to define a catchment window  
440 of 60kbp. From the 805,426 variants profiled in the UK Biobank genotyping data after quality control,  
441 4,248 variants within the transcribed region of the genes of interest or within 60kbp flanking the  
442 transcribed region. After applying additional QC filters using PLINK2 (see *Data Source* above), 2,097  
443 SNPs remained for analysis. We calculated counts for each variant for each of our groups of interest listed  
444 in *Genetic Association Studies* above.

445

### 446 *SNP Set Empirical Statistical Evaluation*

447 To assess the probability of getting 10 study-wide significant hits (using BH corrected p-value < 0.05),  
448 we used empirical sampling to generate 100 sets of randomly chosen SNPs. In each sample, 69 genes  
449 were chosen at random from the genome and mapped to nearby SNPs (within a 60kbp flanking region),  
450 resulting in sets of SNPs sized 1712 to 2945 – similar to the 2097 that resulted from our complement and  
451 coagulation cascade set. We then repeated the association study and counted the number of significant  
452 hits (using BH corrected p-value < 0.05). We fit the empirical data to a Poisson distribution and used the  
453 derived lambda to compute p-values for our observations of 10, 4, and 1 hit (corresponding to the number  
454 of significant results from our severe analysis, positive analysis, and negative, respectively). We  
455 performed a chi-squared goodness-of-fit test to verify the data were consistent with a Poisson.

456

#### 457 *Software*

458 We used PLINK v2.00a2LM 64-bit Intel (26 Aug 2019) to run the genetic association analysis.

459

460

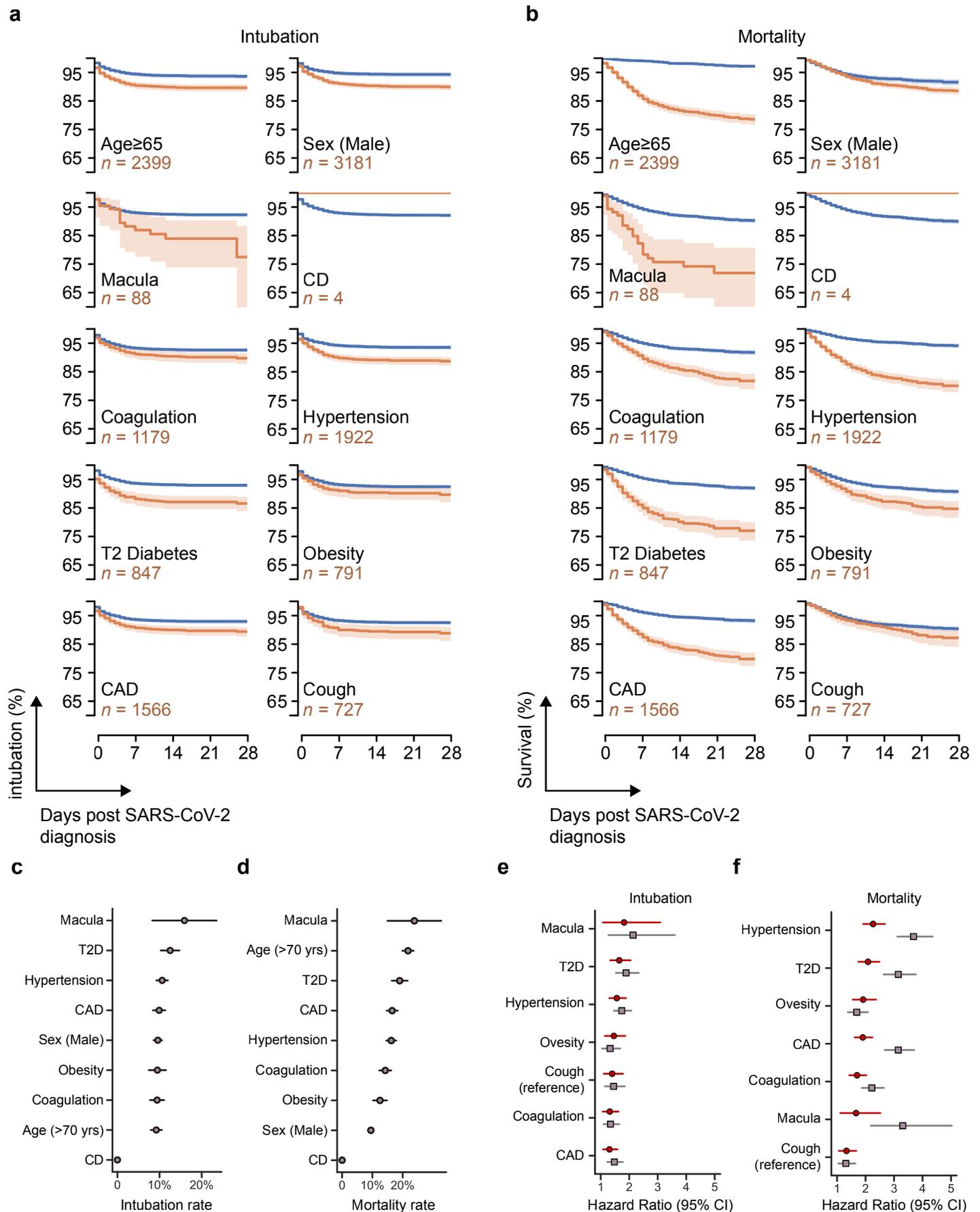
461 **REFERENCES**

462

- 463 1 Zhang, L. *et al.* Crystal structure of SARS-CoV-2 main protease provides a basis for design  
464 of improved alpha-ketoamide inhibitors. *Science* **368**, 409-412,  
465 doi:10.1126/science.abb3405 (2020).
- 466 2 Dai, W. *et al.* Structure-based design of antiviral drug candidates targeting the SARS-  
467 CoV-2 main protease. *Science*, doi:10.1126/science.abb4489 (2020).
- 468 3 Gordon, D. E. *et al.* A SARS-CoV-2-Human Protein-Protein Interaction Map Reveals Drug  
469 Targets and Potential Drug-Repurposing. *bioRxiv*, 2020.2003.2022.002386,  
470 doi:10.1101/2020.03.22.002386 (2020).
- 471 4 Chen, G. *et al.* Clinical and immunological features of severe and moderate coronavirus  
472 disease 2019. *J Clin Invest*, doi:10.1172/JCI137244 (2020).
- 473 5 Moore, B. J. B. & June, C. H. Cytokine release syndrome in severe COVID-19. *Science*,  
474 doi:10.1126/science.abb8925 (2020).
- 475 6 Lasso, G. *et al.* A Structure-Informed Atlas of Human-Virus Interactions. *Cell* **178**, 1526-  
476 1541 e1516, doi:10.1016/j.cell.2019.08.005 (2019).
- 477 7 Merle, N. S., Church, S. E., Fremeaux-Bacchi, V. & Roumenina, L. T. Complement System  
478 Part I - Molecular Mechanisms of Activation and Regulation. *Front Immunol* **6**, 262,  
479 doi:10.3389/fimmu.2015.00262 (2015).
- 480 8 Holers, V. M. Complement and its receptors: new insights into human disease. *Annu Rev*  
481 *Immunol* **32**, 433-459, doi:10.1146/annurev-immunol-032713-120154 (2014).
- 482 9 Wu, J. & Sun, X. Complement system and age-related macular degeneration: drugs and  
483 challenges. *Drug Des Devel Ther* **13**, 2413-2425, doi:10.2147/DDDT.S206355 (2019).
- 484 10 Ambati, J., Atkinson, J. P. & Gelfand, B. D. Immunology of age-related macular  
485 degeneration. *Nat Rev Immunol* **13**, 438-451, doi:10.1038/nri3459 (2013).
- 486 11 Khandhadia, S., Cipriani, V., Yates, J. R. & Lotery, A. J. Age-related macular degeneration  
487 and the complement system. *Immunobiology* **217**, 127-146,  
488 doi:10.1016/j.imbio.2011.07.019 (2012).
- 489 12 Degen, S. E., Jensenius, J. C. & Thiel, S. Disease-causing mutations in genes of the  
490 complement system. *Am J Hum Genet* **88**, 689-705, doi:10.1016/j.ajhg.2011.05.011  
491 (2011).
- 492 13 Liszewski, M. K., Java, A., Schramm, E. C. & Atkinson, J. P. Complement Dysregulation  
493 and Disease: Insights from Contemporary Genetics. *Annu Rev Pathol* **12**, 25-52,  
494 doi:10.1146/annurev-pathol-012615-044145 (2017).
- 495 14 Zhou, F. *et al.* Clinical course and risk factors for mortality of adult inpatients with  
496 COVID-19 in Wuhan, China: a retrospective cohort study. *Lancet* **395**, 1054-1062,  
497 doi:10.1016/S0140-6736(20)30566-3 (2020).
- 498 15 Nicholson-Weller, A. & Wang, C. E. Structure and function of decay accelerating factor  
499 CD55. *J Lab Clin Med* **123**, 485-491 (1994).
- 500 16 Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.  
501 *Nature* **562**, 203-209, doi:10.1038/s41586-018-0579-z (2018).
- 502 17 Smith, N. L. *et al.* Genetic predictors of fibrin D-dimer levels in healthy adults. *Circulation*  
503 **123**, 1864-1872, doi:10.1161/CIRCULATIONAHA.110.009480 (2011).

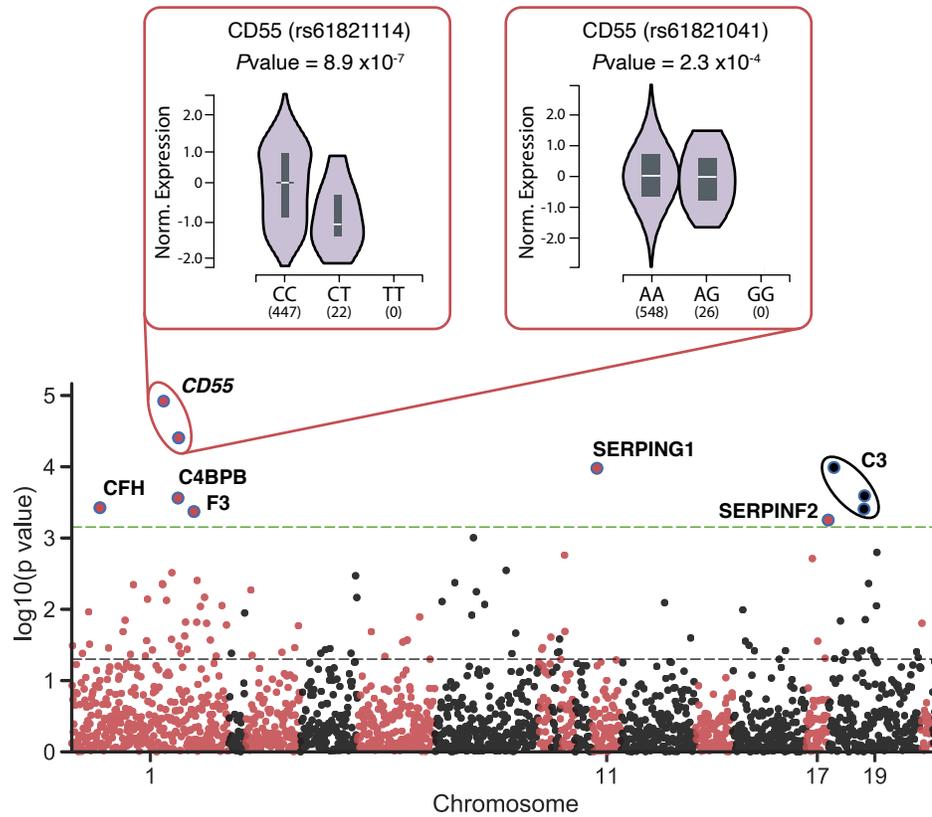
- 504 18 Han, X. *et al.* Genome-wide meta-analysis identifies novel loci associated with age-  
505 related macular degeneration. *J Hum Genet*, doi:10.1038/s10038-020-0750-x (2020).
- 506 19 Goeijenbier, M. *et al.* Review: Viral infections and mechanisms of thrombosis and  
507 bleeding. *J Med Virol* **84**, 1680-1696, doi:10.1002/jmv.23354 (2012).
- 508 20 Polubriaginof, F. C. G. *et al.* Challenges with quality of race and ethnicity data in  
509 observational databases. *J Am Med Inform Assoc* **26**, 730-736,  
510 doi:10.1093/jamia/ocz113 (2019).
- 511 21 Hosmer, D. W., Lemeshow, S. & May, S. *Applied survival analysis : regression modeling*  
512 *of time-to-event data*. 2nd edn, (Wiley-Interscience, 2008).
- 513 22 Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer  
514 datasets. *Gigascience* **4**, 7, doi:10.1186/s13742-015-0047-8 (2015).  
515

**Figure 1**



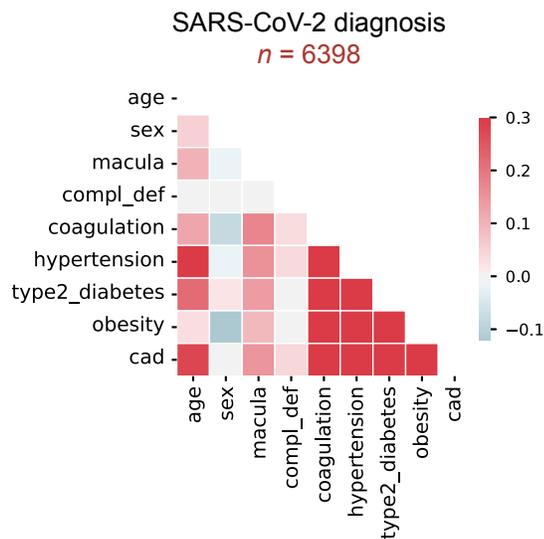
**This report has not been certified by peer review. This should not be relied on to guide clinical practice or health-related behavior and should not be reported in news media as established information.**

**Figure 2**

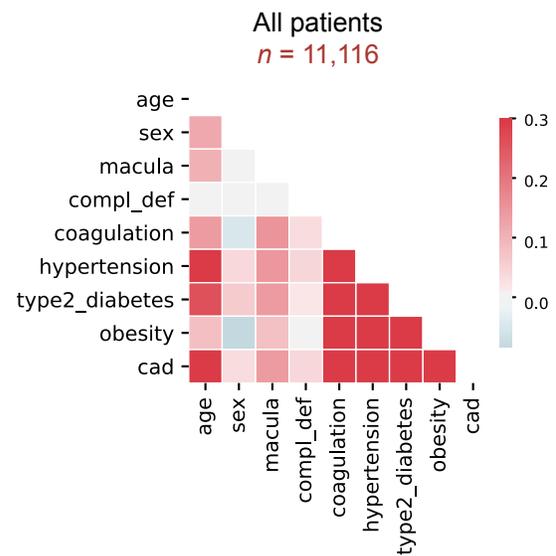


## Figure S1

**a**



**b**





Variant	Chr	Position	Nearby Gene(s)	eQTL Gene(s)	Locus is GWAS hit for AMD in UKBiobank	SARS-CoV-2 Positive and Hospitalized (N=332)						Others (N=337,147)			Odds Ratio (95%CI)	Unadjusted p-value	Adjusted p-value (BH)
						A	B	AA	AB	BB	AA	AB	BB				
rs45574833	1	207300070	C4BPB, C4BPA	--	yes	A	G	1	19	368	45	8144	328570	2.65 (1.71-4.09)	1.20x10 <sup>-5</sup>	4.55x10 <sup>-3</sup>	
rs61821114	1	207610967	CR2, CR1	CD55 (Pvalue=9x10 <sup>-7</sup> )	yes	T	C	1	21	366	72	9816	326871	2.40 (1.58-3.63)	3.94x10 <sup>-5</sup>	1.50x10 <sup>-2</sup>	
rs61821041	1	207352581	C4BPA	CD55 (Pvalue=2x10 <sup>-4</sup> )	yes	G	A	1	17	370	47	8831	328381	2.34 (1.48-3.69)	2.74x10 <sup>-4</sup>	3.22x10 <sup>-2</sup>	
rs12064775	1	196600605	CFH	--	yes	G	A	0	24	364	95	10988	325677	2.13 (1.41-3.23)	3.71x10 <sup>-4</sup>	3.22x10 <sup>-2</sup>	
rs72729504	1	94940206	F3	--	no	T	C	0	33	355	191	15646	320922	1.93 (1.34-2.79)	4.24x10 <sup>-4</sup>	3.22x10 <sup>-2</sup>	
rs117284601	11	57425228	SERPING1	--	no	A	G	1	49	338	544	26315	309900	1.80 (1.34-2.42)	1.60x10 <sup>-4</sup>	7.42x10 <sup>-3</sup>	
rs9913923	17	1703982	SERPINF2	--	no	T	C	6	87	295	3096	57767	275896	1.48 (1.19-1.85)	5.59x10 <sup>-4</sup>	3.02x10 <sup>-2</sup>	
rs1047286	19	6713262	C3	--	no	A	G	7	109	272	15403	113286	208070	0.66 (0.53-0.81)	1.02x10 <sup>-4</sup>	2.28x10 <sup>-2</sup>	
rs2230203	19	6710782	C3	--	no	T	G	7	98	283	12134	104411	220214	0.66 (0.53-0.82)	2.57x10 <sup>-4</sup>	2.87x10 <sup>-2</sup>	
rs2230199	19	6718387	C3	--	no	C	G	7	112	260	14577	108566	201225	0.68 (0.55-0.84)	3.92x10 <sup>-4</sup>	2.93x10 <sup>-2</sup>	

**This report has not been certified by peer review. This should not be relied on to guide clinical practice or health-related behavior and should not be reported in news media as established information.**