Initial Whole Genome Sequencing and Analysis of the Host Genetic Contribution to COVID-19 Severity and Susceptibility

- 3
- 4 Fang Wang^{1*}, Shujia Huang^{2,3*}, Rongsui Gao^{1*}, Yuwen Zhou^{2,4*}, Changxiang Lai^{1*}, Zhichao
- 5 Li^{2,4}*, Wenjie Xian¹, Xiaobo Qian^{2,4}, Zhiyu Li¹, Yushan Huang^{2,4}, Qiyuan Tang¹, Panhong
- 6 Liu^{2,4}, Ruikun Chen¹, Rong Liu², Xuan Li¹, Xin Tong², Xuan Zhou¹, Yong Bai², Gang Duan¹,
- 7 Tao Zhang², Xun Xu^{2,5}, Jian Wang^{2,6}, Huanming Yang^{2,6}, Siyang Liu^{2#}, Qing He^{1#}, Xin Jin^{2,3#},
- 8 Lei Liu^{1#}
- 9 1. The Third People's Hospital of Shenzhen, National Clinical Research Center for Infectious
- 10 Disease, The Second Affiliated Hospital of Southern University of Science and Technology,
- 11 Shenzhen 518112, Guangdong, China
- 12 2. BGI-Shenzhen, Shenzhen 518083, Guangdong, China
- 13 3. School of Medicine, South China University of Technology, Guangzhou 510006,14 Guangdong, China
- 15 4. BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083,
- 16 Guangdong, China
- 17 5. Guangdong Provincial Key Laboratory of Genome Read and Write, BGI-Shenzhen,
- 18 Shenzhen, 518120, China
- 19 6. James D. Watson Institute of Genome Science, 310008 Hangzhou, China
- 20 *Those authors contribute equally
- 21 Corresponding to any of the followings:
- 22 Lei Liu liuleiszsdsrmyy@163.com
- 23 Xin Jin jinxin@genomics.cn
- 24 Qing He heqingjoe@163.com
- 25 Siyang Liu liusiyang@genomics.cn
- 26
- 27 NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

28 Abstract

The COVID-19 pandemic has accounted for more than five million infections 29 and hundreds of thousand deaths worldwide in the past six months. The 30 patients demonstrate a great diversity in clinical and laboratory manifestations 31 32 and disease severity. Nonetheless, little is known about the host genetic contribution to the observed inter-individual phenotypic variability. Here, we 33 report the first host genetic study in China by deeply sequencing and analyzing 34 35 332 COVID-19 patients categorized by varying levels of severity from the Shenzhen Third People's Hospital. Upon a total of 22.2 million genetic variants, 36 37 we conducted both single-variant and gene-based association tests among five 38 severity groups including asymptomatic, mild, moderate, severe and critical ill patients after the correction of potential confounding factors. The most 39 40 significant gene locus associated with severity is located in TMEM189-UBE2V1 41 involved in the IL-1 signaling pathway. The p.Val197Met missense variant that 42 affects the stability of the TMPRSS2 protein displays a decreasing allele 43 frequency among the severe patients compared to the mild and the general 44 population. We also identified that the HLA-A*11:01, B*51:01 and C*14:02 alleles significantly predispose the worst outcome of the patients. This initial 45 46 study of Chinese patients provides a comprehensive view of the genetic difference among the COVID-19 patient groups and highlighted genes and 47 variants that may help guide targeted efforts in containing the outbreak. 48 49 Limitations and advantages of the study were also reviewed to guide future 50 international efforts on elucidating the genetic architecture of host-pathogen interaction for COVID-19 and other infectious and complex diseases. 51

52

53 Introduction

54 It has been more than 100 years since the 1918 influenza outbreak killed at least fifty million people worldwide¹. Now we are facing another pandemic. 55 Since the late December of 2019, the 2019 novel coronavirus diseases 56 57 (COVID-19) has spread rapidly throughout the world, resulting in more than five million confirmed cases and hundreds of thousands deaths in less than six 58 months^{2,3}. The disease was caused by the infection of a novel enveloped RNA 59 60 betacoronavirus that has been named severe acute respiratory syndrome 61 coronavirus 2 (SARS-CoV-2), which is the seventh coronavirus species that causes respiratory disease in humans^{4,5}. The virus causes serious respiratory 62 illnesses such as pneumonia, lung failure and even death⁶. Until now, there is 63 no specific therapeutics and vaccine available for its control. Continuing 64 epidemiological and molecular biological study to better understand, treat and 65 prevent COVID-19 are urgently needed. 66

A characteristic feature of many human infections is that only a proportion of exposed individuals develop clinical disease and for the infected persons, severity varies from person to person⁷. In the COVID-19 outbreak, a high level of inter-individual variability was observed in terms of disease severity and symptomatic presentation. Around 80%-85% of the laboratory confirmed patients were classified as mild (i.e. nonpneumonia and mild pneumonia) while 15%-20% would progress to severe or critical stage with a high probability of

respiratory failure⁸⁻¹¹. Patients with severe disease had more prominent 74 75 laboratory abnormalities including lymphocytopenia and leukopenia than those with non-severe disease^{12,13}. In addition, not all people exposed to SARS-CoV-76 2 were infected according to the epidemiological observation of the patients' 77 close contacts^{14,15}. Notably, previous studies have indicated that genetic 78 background plays an essential role in determining the host responses to 79 HIV¹⁶⁻¹⁸, HBV¹⁹, HCV²⁰, influenza²¹⁻²⁴, SARS-CoV^{25,26} and 80 infections by numerous common viruses²⁷ etc. Those studies highlighted the HLA alleles and 81 several genes involved in the interferon production and viral replication pathway 82 and indicates that genetic factors may also play an important role to explain the 83 inter-individual clinical variability among patients infected by SARS-CoV-2. 84

Till now, the global genetic community has been actively investigating in 85 86 the genetic contribution to COVID-19. A recent twin study in UK suggests a 30% 87 -50% genetic heritability for self-reported symptoms of COVID-19 and the predictive disease onset²⁸, indicating a very strong genetic background 88 predisposing the COVID-19 patients' clinical manifestation and susceptibility. 89 90 An earlier studies comparing the distribution of ABO blood group from 1,775 patients infected with SARS-CoV-2 with 3,694 normal people from Wuhan city 91 and 23, 386 people from Shenzhen city suggested that blood group A had a 92 significantly higher risk for COVID-19 (OR=1.20, p=0.02) while blood group O 93 had the lower risk²⁹. Using allele frequency and expression quantitative loci 94 (eQTL) information of general healthy population from 1000 genome project 95 and others, a few studies investigate the mutation frequency spectrum in 96 different populations in candidate genes such as ACE2 and TRMPSS2³⁰⁻³². 97 Genome-wide association test on array data from the UK Biobank participants 98 with a positive and negative PCR-tests also reveals a few suggestive genes²⁷. 99 The COVID-19 host genetics initiative was established to encourage generation, 100 sharing and meta- analysis of the genome-wide association summary statistics 101 data around the world³³. International collaborative efforts are necessary to 102 elucidate the role of host genetic factors defining the severity and susceptibility 103 104 of the SARS-CoV-2 virus pandemic.

Herein, we report the first genetic study of COVID-19 disease severity in 105 China by deeply analyzing the association between the genetic variants present 106 in the patients' genome and their disease progression. We have recruited 332 107 hospitalized patients from a designated infectious disease hospital in Shenzhen 108 City³⁴. The patients display varying clinical and laboratory features and were 109 categorized as asymptomatic, mild, moderate, severe and critical cases 110 according to the criteria made by the Chinese Center for Disease Control and 111 Prevention⁶. To maximize the statistical power given the relatively 112

small hospitalized sample size and for accurate detection of extremely rare variants, we conducted deep whole genome sequencing (average 46x) for the patients. Given a fixed samples size, this protocol facilitates the estimation of genetic effects of rare and loss of function variants in addition to the common variants that may be potentially contributing to the COVID-19 clinical variability³⁵. Based on the 22.2 million variation detected from the patients, we investigated host factors by conducting both single variant and gene-based genome-wide

association study and by evaluating the difference of allele frequency of the 120 protein truncating variants and HLA alleles among the patient groups. In 121 addition, we performed joint-calling of the genetic variants of the unrelated 122 COVID-19 patients (N=284) and the publicly available Chinese genomes from 123 the 1000 genome project³⁶ (N=301, \sim 7x) and 665 selected Chinese genomes 124 from the Chinese Reference Panel Population (manuscript in preparation, ~30x) 125 126 to explore potential genetic factors that may contribute genetic susceptibility of 127 SARS-CoV-2 infection.

128

129

130 **Results**

131 Clinical and laboratory features of the 332 hospitalized COVID-19 patients

The 332 recruited patients with laboratory-confirmation of SARS-COV-2 infection were being quarantined and treated in the Shenzhen Third Hospital. We extracted and analyzed the clinical symptoms, laboratory assessment, recent exposure history of the patients from the hospital's electronic medical records. The 332 patients consist of 48 family members and 284 unrelated individuals.

25 (7.5%), 12 (3.6%), 225 (67.8%), 53 (16.0%) and 17 (5.1%) patients 138 were defined as asymptomatic, mild, moderate, severe and critically severe 139 140 according to the most severe stage they encountered during the disease 141 course following the Chinese CDC criteria⁶ (Figure 1A). The asymptomatic, 142 mild and the moderate groups of patients had positive RT-PCR test result but 143 did not have or only had mild pneumonia. The severe patients met any one of 144 the following criteria: respiratory rate (RR) \geq 30/min, blood oxygen saturation \leq 93%, partial pressure of arterial oxygen to fraction of inspired oxygen ratio 145 146 (PaO2/FiO2) < 300 mmHg and/or lung infiltrates > 50% within 24-48 hours. A 147 severe patient was classified as critical ill if he/she experienced any one of the following situations: respiratory failure, septic shock and/or multiple organ 148 149 dysfunction or failure. A broader definition of the mild group includes the 150 asymptomatic, mild and moderate patients, and of the severe group, includes the severe and critically severe patients. 151

The patients displayed several clinical presentations typical to COVID-152 153 19, which mainly involved fever (70.8%), cough (54.2%), fatigue (23.9%), hoarse voice (17.6%), loss of appetite (16.2%), delirium (15.1%) (Figure S1). 154 155 Less than 10% had also experienced diarrhea, chest and abdominal pain, shortness of breath and anosmia. More than 50% of the patients had at least 156 one medical comorbidities (e.g., hypertension). Consistent with previous 157 158 report, the broadly defined severe patients tend to be older (severe average 45 years old vs mild average 58 years old, t-test p=0.03, Figure 1C), suffer 159 from a longer course of disease between the onset and the first negative RT-160 PCR test outcome (Figure 1D) and shorter exposure time (Figure S2) . In 161 addition, the severe patient group consist of more males than females (severe 162 66.7% vs mild 41.3%, χ^2 test p=4.3e-4, **Figure 1E**) and tend to undergo 163 medical comorbidities more frequently (severe 58.8% vs mild 45.1%, x^2 test 164 p=0.07) (Figure 1F) than the mild patients. 165

During hospitalization, a series of sixty four laboratory assessments 166 including a complete blood count and blood chemical analysis, assessment of 167 liver function, assessment of renal functions, test of humoral immunity, test of 168 coagulation, measure of electrolyte and measure of blood gas electrolyte 169 (Figure S3) and a time-series evaluation of T lymphocyte subgroups (Figure 170 S4) were performed for each of the patients to monitor their disease status and 171 172 progression. Using a tree-based machine learning prediction model³⁷, we computed the local interaction effects of the sixty four laboratory assessment 173 features as well as three demographic features including age, gender and w/o 174 175 medical comorbidities for classification of the patient severity category (Figure **S5**). The top ten features of greatest importance that contribute to a severer 176 disease outcome include decreased lymphocyte counts (Tc-Count, T-CELL, 177 178 LYMPH#) and platelet counts, evaluated interleukin 6, C-reactive protein and 179 D-dimer, increased age and decreased A/G and CO2 (Figure 1B), consistent with previous reports³⁸. We applied the top twenty features of importance to 180 181 assign a severity score for each patient to reflect their disease status (Figure 182 S6).

183

184 Deep whole genome sequencing and genetic variation

We obtained the whole blood and performed deep whole genome sequencing for the recruited patients. There is no significant difference for sequencing depth between the broadly defined mild and severe group (mild 46.26x vs severe 46.71x) (**Figure 2A**). We conducted variation detection and genotyping using the GATK joint genotyping framework to avoid any potential batch effect derived from individual variant calling. Bioinformatics analysis and the data quality control process were described in details in the **Online methods**.

Among the 332 patients, we identified a total of 22.2 million variants 192 including 17.9 million bi-allelic single nucleotide polymorphism, 1.75 million bi-193 allelic small insertions and deletions and 2.49 million multi-allelic variants 194 (Figure 2B). The average transition/transversion (ts/tv) ratio is 2.12 and the 195 proportion of heterozygous versus homozygous variants among all the samples 196 is 1.29, consistent with our expectation³⁹ and indicates good quality of the 197 variant calls (Figure S7). Particularly, we have identified 398K variants that 198 result in an alteration of the protein coding sequence (Figure 2C). The QC 199 metrics were detailed in Table S1. 200

Our first question was whether the most vulnerable severe and critical 201 patients may have a monogenic basis for their demonstration. We investigated 202 the burden of loss of function variants predicted by the ensemble variant effect 203 predictor among the patient groups⁴⁰. In total, we have identified 4,891 204 predicted loss of function variants including 1,860 frameshift, 1,447 stop gained, 205 505 splice donor and 380 splice acceptor variants among the 332 patients. On 206 207 average, each patient possessed 201 predicted loss of function variants in their genome (Figure S8). 261 of those variants were uniquely presented in the 208 COVID-19 patients (18.6%) and have not been previously reported in the 1000 209 genome and the gnomAD studies^{36,41,42}. Interestingly, the severe and the critical 210 patients tend to have more loss of function insertions than the asymptomatic, 211

mild and the moderate groups in a logistic regression taking the number of loss 212 213 of function variants as variable and the patients' age, gender, the twenty principle components and effective sequencing depth as covariates (p=0.004) 214 (Figure 2D). When performing a mutation burden test for each of the 16,801 215 genes that have more than one variant among the 284 unrelated patients, we 216 didn't identify genes that were enriched in loss of function variants in the severe 217 and critical patients (Figure 2E). On the other hand, we found two heterozygous 218 219 loss of function variants located in MST1R and RASA2 that were only present in the asymptomatic patients (Figure S9). The MST1R encodes the 220 221 macrophage stimulating 1 receptor expressed on the ciliated epithelia of the mucociliary transport apparatus of the lung and follows an autosomal dominant 222 inheritance mode for susceptibility to nasopharyngeal carcinoma⁴³. Because 223 224 those loss of function variants were only present in one patient, we didn't build 225 up links to the COVID-19 severity.

226 Particularly, we have inspected the missense and loss of function variants 227 present in the SARS-CoV-2 S protein host cellular receptor gene ACE2 and the 228 S protein primer gene *TMPRSS2* that plays a critical role in controlling the viral entry into the host cell, as well as a few other genes that were predicted to play 229 a role in the host pathogen interaction network like SLC6A19, ADAM17, RPS6, 230 HNRNPA1, SUMO1, NACA and BTF3⁴⁴. The majority of the functional variants 231 have minor allele frequency less than 1% except for the p.Val197Met missense 232 variant in TMPRSS2 (Figure 2F). Although not statistically significant, the 233 p.Val197Met variant (rs12329760) displays a higher allele frequency in the 234 asymptomatic and mild group compared to the rest of the group (asymptomatic: 235 236 0.46, mild: 0.50, moderate: 0.38, severe: 0.39, critical severe: 0.26). 237 p.Val197Met was previously found to have higher allele frequency in East Asian (0.31-0.41) and Finnish (0.36) but is less frequently seen in South Asians (0.14-238 0.29) and the Europeans (0.17-0.23) (Figure S10). By computational protein 239 modelling, the p.Val197Met TMPRSS2 isoform could decrease the stability of 240 the TMPRSS2 protein, promote the binding to S-protein and inhibit its binding 241 with ACE2⁴⁴. The decreasing allele frequency in the severe patient groups 242 supports that the p.Val197Met is related to the disease outcomes of COVID-19. 243 The other genes didn't contain significant allele frequency difference among the 244 patient groups (Figure S11). 245

246

247 Genetic association of common and rare variants with COVID-19 severity

To further investigate genetic effects for the patient severity, we performed 248 genome-wide single variant association test and sequence kernel association 249 250 test (SKAT) analysis of three traits implicating patient severity. We defined the first trait as a dichotomous classification of the broadly defined "severe group" 251 that consists of the severe and critical ill patients (N=70) and the "mild group" 252 (N=262) that consists of the asymptomatic, mild and moderate patients. We 253 defined the second trait as a quantitative measurement of the severity level 254 trained from the demographic features such as age, gender and the sixty-four 255 laboratory assessments (N=332) (Figure S5-6). We used the disease duration 256 from the electronic health records as the third trait which corresponds to the 257

duration of time between the complained disease onset and the first laboratory 258 259 confirmed PCR-test negative outcome (N=233) (Figure 1D). Power analysis indicates that given 80% statistical power, we will be able to identify 260 associations between genotypes and phenotypes for variants with minor allele 261 frequency greater than 0.2 and with a relative genetic risk contribution greater 262 than 2 given the current sample size for dichotomous trait and similarly for the 263 264 quantitative trait (Figure S12). Principal component analysis of the patients suggests little genetic differentiation (Figure S13-14). 265

We tested all the QC-passed 19.6 million bi-allelic variants for association 266 267 with each of three traits in a logistic or linear regression model that includes gender, age, and the top 20 PC axes as covariates. The global distribution of 268 resulting p-values was very close to the null expectation (λ =0.996~1.1, Figure 269 270 **S15**) indicating that stratification was adequately controlled. The most 271 significant SNP rs6020298 is located in the intron of a read-through transcript 272 *TMEM189-UBE2V1* in the 20q13.13 region. (Figure 3A-B). The rs6020298 273 (hg38 chr20:50152518, A allele frequency severe vs non-severe: 0.59 vs 0.45) 274 marks a suggestive significant association signal for both the Severe and Mild binary trait (logistic regression p=4.1e-6, OR=1.2) and the quantitative 275 measurement of the severity score (linear regression p=1.1e-6, beta=0.35). 276 SNPs in linkage disequilibrium with rs6020298 (r²>0.8) also affect the gene 277 UBE2V1 and TMEM189 (Figure 4A). The UBE2V1 gene encodes the ubiquitin-278 conjugating enzyme E2 variant 1. Both the UBE2V1 and TMEM189-UBE2V1 279 have been involved in the interleukin-1 (IL-1) signaling pathway⁴⁵ and 280 suggested to work together with TRIM5 to promotes innate immune signaling⁴⁶. 281 IL-1 is elevated in COVID-19 patients especially the severe and critical patients 282 who suffer from the cytokine storm and severe inflation⁴⁷. Clinical trial using IL-283 1 blockade on critical patients results in an improvement in respiratory function 284 in 72% of the patients⁴⁸. The lead SNP rs6020284 has a minor allele frequency 285 close to 0.5 among the worldwide populations except for the African population 286 (AF=0.13) (Figure 4B). It is also an eQTL for LINC01273, TMEM189 among 287 several tissues including the lung where the risk A allele increases the 288 TMEM189 and LINC01273 expression in several tissues (Figure S16). This 289 may indicate that an inborn evaluated TMEM189 expression in the patients may 290 promote IL-1 signaling and predisposes the patients towards a poorer outcome 291 against the COVID-19 infection. However, given the limited sample size in this 292 study and that the intermediate pathways between TMEM189 and IL-1 293 production is still unclear, more replication and functional validation efforts 294 should be made to re-evaluate this association signal. Notably, the TMEM189-295 UBE2V1 locus has been associated with monocyte percentage of leukocytes 296 and granulocyte percentage of myeloid white cells⁴⁹. Nonetheless, we didn't 297 observe nominal association (p<0.05) at the lead SNP rs6020298 with all the 298 sixty-four laboratory assessments among the patients (Figure 4C). Therefore, 299 the observed signal is not supposed to be confounded by individual variability 300 on blood cell types. There is no strong genetic association with the disease 301 durations (Figure 3C). 302

We further performed optimal SKAT gene-based association test on the 303 304 functional variants including a total of 99,166 missense and loss of function variants that were predicted to have high or moderate impacts by variant effect 305 predictor among the patients. The *NOA1* gene tend to higher mutation burden 306 307 in the severe group (P= 8.1e-07) (**Figure 3D**). This gene encodes the GTPase that functions in the mitochondrion and has been associated with platelet count 308 309 and leukocyte count⁴⁵. We didn't identify other genes that are genome-wide significantly associated with the severity score or the disease duration (Figure3 310 **E-F**). 311

312

313 HLA gene alleles associated with severity in the COVID-19 patients

Manifestation of numerous infectious diseases are closely related to the genetic 314 variants across the major histocompatibility complex (MHC) genes, i.e. the 315 human leukocyte antigen (HLA) genes, which play an essential role in 316 presenting the antigen determinant epitopes from the pathogens to the T cell or 317 B cell to activate the host immune response^{50,51}. In the 2003 severe acute 318 respiratory (SARS) outbreak, caused by the SARS coronavirus (SARS-CoV) 319 related to SARS-CoV-2, the HLA-B*46:01 was reported to be associated with 320 infection severity in East Asian patients²⁵. Herein, we investigated the genetic 321 effect from HLA genes on the COVID-19 patient severity. We re-aligned all the 322 reads mapped to the eight HLA haplotypes in the human reference genome 323 (GRCh38) and all the unaligned reads and typed the three class I HLA genes 324 (A, B, C) and four class II HLA genes (DPB1, DQA1, DQB1, DRB1) using the 325 xHLA⁵² and the SOAP-HLA approach⁵³. 4-digit haplotyping resolution was 326 achieved for 99% of the patients for all the genes except for DQA1 where three 327 patients were only typed to the 2-digit resolution. We observed zero mendelian 328 error rate for the typing results using the family members involved in the study. 329 We investigated whether some HLA alleles may significantly differ between the 330 broadly defined severe (severe and critical, N=69) and mild (asymptomatic, mild 331 and moderate, N=215) groups of unrelated patients using a logistic regression 332 333 with age, gender and the top 20 principal components as covariates. The frequency comparison between the severe and mild groups for the total 30 HLA-334 A, 51 HLA-B, 28 HLA-C, 20 DPB1, 21 DQA1, 16 DQB1 and 32 DRB1 alleles 335 were displayed in Figure 4 and Table S2. Among the class I HLA genes, 336 C*14:02 (severe 8.7% vs mild 4.6%, OR=4.7, P=3e-3), B*51:01 (severe 10.1%) 337 vs mild 5.8%, OR=3.3, P=7e-3), A*11:01 (severe 29.7% vs 26.2%, OR=2.3, 338 P=8.5e-3) are the top three most significant alleles between the two groups that 339 predispose the patients entering the severe stage (Table 1). The HLA-A*11:01, 340 B*51:01 and C*14:02 is in strong linkage equilibrium with each other and thus 341 represents one haplotype. This haplotype has an average allele frequency 2.4% 342 - 3.6% among the Chinese populations according to the HLA Allele Frequency 343 Net Database⁵⁴. In our study, we find that this haplotype is more prevalent in 344 the severe patients compared with the mild patients. 345

Notably, although B*46:01 has been suggested to present the fewest SARS-CoV and SARS-CoV-2 peptides in an in silico analysis⁵⁵ and has been associated with the SARS-CoV in a small sample size association analysis

without correcting demographic and geographic covariates²⁵, our analysis 349 350 doesn't support this allele is associated with the disease severity (OR=0.5, P=0.15). On the contrary, allele frequency of B*46:01 is less frequent in the 351 severe patients (10.1%) than among the mild patients (12.8%). Class II HLA 352 genes is less significantly associated with the disease severity compared to the 353 Class I genes (Table 1). DRB1*14:04 (severe 2% vs mild 0.5%, P=0.01), 354 355 DRB1*01:01 (severe 2.2% vs 0.5%), DQA1*01:01 (severe 2.9% vs 0.9%) are the top three risk alleles while DPB1*03:01 (severe 0.7% vs mild 4.5%) and 356 DRB1*12:01 (severe 2.2% vs mild 3.7%) might display a protective effect. 357

358

Comparison with general population for potential genetic contribution to SARS-CoV-2 infection susceptibility

Our study till now has been restricted in the infected patients to understand 361 genetic contribution to patient severity. Mapping genes related to infection 362 susceptibility is more difficult. The ideal design commands a comparison 363 between people who are exposed or not exposed to the pathogen. This is 364 hard to meet because early detection and isolation of infected patients are the 365 primary containment strategies against an outbreak⁵⁶. Therefore, we choose 366 another approach to investigate genetic susceptibility by comparing the 284 367 368 unrelated hospitalized patients (the Case) with two general populations including 301 Chinese individuals in 1000 genome project³⁶ (the Control I) and 369 665 individuals recruited from the Chinese Reference Panel program (CNPR, 370 371 manuscript in preparation, the Control II). Control I and Control II differ in 372 terms of the similarity of the adopted sequencing protocol compared to the 373 Case. All the technical components are almost the same between the Case and Control 2 except for sequencing depth (case 46x versus control 2 30x). 374 375 On the other hand, various factors are different between the Case and Control 1. including types of sample (case fresh blood versus control 1 cell line). 376 sequencing technology (case MGI's nanoball sequencing versus control 1 377 378 Illumina sequencing), sequencing read cycles (case 100bp pair-end versus 379 control 1 150bp pair-end) and the sequencing depth (case average 46x versus average 7x). Study like this can reveal genetic difference between the infected 380 381 population and the general population if any and if not, instruct on what cautions should be taken when comparing the disease cohorts versus the 382 general in the whole genome sequencing context. 383

We analyzed the data carefully by jointly genotype the samples from their 384 individual gvcf files using the GATK best practices³⁹ instead of simply merging 385 the population vcf files of the case and the control. Principle component analysis 386 indicates that population structure is the dominant confounding factor and 387 sequencing induced batch effects were difficult to identify in the PCs (Figure 388 **S17**, **Figure S18**). Similarly, we conducted both single variant and gene-based 389 association tests for the two case-control data sets using the top 20 PCs, 390 gender and age (age was not available for 1KGP samples and was used for the 391 CNPR alone) as covariates. Surprisingly, in the single association test for the 392 high and moderate impact variants, many variants in the HLA region displayed 393 significant associations between the COVID-19 patients and the 1KGP Chinese 394

(Figure 6A) even though the inflation was seemingly adequately controlled 395 (Figure S19). In the gene-based association test, we observed significantly 396 different mutation burdens in the immunoglobulin loci (Figure 6B). However, 397 this was not replicated when we compared the COVID-19 patients with the 665 398 399 CNRP individuals (Figure 6C-D). Therefore, we inferred that the association signals between the 1KGP and the COVID-19 patients were probably due to 400 401 sequencing batch effects. As the fresh blood of an infected individual contains 402 numerous somatic mutated B-cells, patients tend to accumulate more mutations in the immunoglobulin genes⁵⁷. As many studies try to directly compare the 403 allele frequency between the general population and the COVID-19 patients^{30,32}, 404 our discoveries remind us of the necessity for re-evaluation of the significant 405 hits given distinct experimental protocol for case and control. 406

407 In the single variant association test between the COVID-19 patients and 408 the CNPR who were sequenced using the same experimental protocol and 409 were laboratory PCR tested negative, we identified genome-wide significant 410 associated signals tagged by a novel missense variant (Patient T allele frequency=0.34, CNPR T AF=0.14, OR=18, P=4,7e-17) in *MUC2*; a missense 411 variant rs200584390 (Patient G allele frequency=0.31, CNPR G AF=0.09, 412 413 OR=9.29, P=1.5e-13) in *RIMBP3* and a missense variant rs200975425 (Patient T allele frequency=0.24, CNPR T AF=0.39, OR=5.4, P=9.4e-10) in GOLGA8B 414 (Figure 6C). Gene-based association test also indicates that RIMBP3 and 415 416 GOLGA8B were different between the patients and the CNPR (Figure 6D). 417 Those discoveries require further replication and interpretation when more sequencing data for patients and for general populations become available 418 worldwide³³. 419

420

421 Discussion

We have conducted the first genetic association study for the COVID-19 422 severity and SARS-CoV-2 infection susceptibility by studying the genome and 423 clinical outcome of 332 patients in a designated infectious disease hospital in 424 425 the Shenzhen City. Instead of using the microarray or the exome genome sequencing, we have carried out high-depth whole genome sequencing and 426 analysis for the patients to obtain the greatest possible power given a small 427 sample size available so far. The study design enables the detection of very 428 rare and private functional variants for the patients⁵⁸ and ensures that the 429 potential causal variants are directly assayed to compensate the loss of power 430 due to poor linkage disequilibrium between the assayed and the causal 431 variants⁵⁹. 432

We revealed that the disease progression after the SARS-CoV-2 433 infection was a complex event and not explained by a monogenic model. The 434 435 severe and critical patients did not carry causal monogenic variants related to the disease severity in their genome. We identified that the missense variant 436 rs12329760 in TMPRSS2 was less frequent among the critical patients 437 compared to the rest of the patients and the general population. This variant 438 results in an alteration of the valine to the methionine at the 197th amino acids 439 (p.Val197Met) has been predicted to decrease the TMPRSS2 protein stability 440

and ACE2 binding⁴⁴. On the other hand, our study using Chinese samples did
not support the assumption³⁰ that host genetic factors in the essential SARSCoV receptor *ACE2* and some other genes involved in the host pathogen
interaction network might play a role in determining the patient's severity or
susceptibility.

In the genome-wide association analysis, a gene locus around 446 447 TMEM189-UBE2V1 displayed suggestively association with COVID-19 severity. 448 This gene locus contains genes such as UBE2V1 and TEMEM189-UBE2V1 that are known to function in the interleukin-1 signaling pathway^{45,46}. The lead 449 450 SNP rs60220284 is an eQTL where the risk allele A increases the gene expression of genes within the locus⁶⁰ and is more prevalent in the severe and 451 critical patients. While COVID-19 severe patients demonstrate elevated IL-1 452 compared to the mild patients and the general population⁴⁷, our study suggests 453 454 potential correlation between genetic variability in this gene and the disease 455 severity.

Notably, the HLA-A*11:01, B*51:01 and C*14:02 alleles were significantly 456 457 more prevalent in the severe and critical severe patients compared to the mild and the moderate patients after careful control of population structure and 458 459 demographic characters such as age and gender. The three alleles were in linkage disequilibrium with each other and has been previously reported to have 460 a 2-3% population allele frequency in Dai and Jinpo minorities in China⁵⁴ and 461 the B*51:01 has been previously linked to the Behcet's disease⁶¹, a kind of 462 rheumatic disease. We were not able to access the role of HLA-B*46:01, 463 although it has been predicted as the worst presenting HLA alleles to the SARS-464 CoV-2 proteome⁵⁵ and linked to the SARS 2003 outbreak²⁵. 465

Surprisingly, genome-wide association study using the COVID-19 466 patients as the case and the 1000 genome Chinese population as the control 467 suggested an enrichment of significantly associated signals in the HLA region 468 and mutation burden in the immunoglobulin genes. Nonetheless, this was not 469 replicated when we compared the patients to another independent Chinese 470 population. A lot of efforts in the genetic field have been made and there may 471 be more in the future to investigate genetic susceptibility of the SARS-COV-2 472 infection by directly comparing two or more general populations with the 473 COVID-19 patients^{32,33}. Therefore, cautions should be taken to properly control 474 the batch effects. Replication is essential and perhaps a joint-analysis effort can 475 rule out the real signals from the false delusion. 476

Some limitations of the study should be noted. Power analysis indicates 477 that sample size of around 300 is barely sufficient to identify genome-wide 478 significant genetic variants with minor allele frequency greater than 0.2 and 479 odds ratio greater than 1.8 given type I error rate 0.05. We don't have power to 480 detect causal variants beyond this risk and allele frequency scenario. In addition, 481 although the study of hospitalized patients in a designated hospital includes all 482 severe patients, the design has a limited presentation of the asymptomatic 483 patients (7.5%) which ratio has been estimated to be 30.8% (95% confidence 484 interval 7.7-53.8%)⁶². Given that RT-PCR test and the seroprevalence 485 immunoglobulin M and G antibody tests targeting the SARS-CoV-2 has been 486

widely adopted in China and around the globe, it will be important to identify
and study the extreme asymptomatic patients to understand the host factors
contributing to a capable control of the viral infection.

As we and the others are continuing to recruit patients and data in China 490 491 and around the world to understand the host genetic background underlying the varying clinical outcome of the patients, this work represents the first genetic 492 493 study on the Chinese hospitalized patients where high quality sequencing data 494 were generated and systematic analysis on the genomic and clinical data were conducted. Our results highlight several genetic factors involved in the immune 495 496 responses including genes involved in the viral entry in the host cells, genes 497 related to immune responses and the HLA alleles. This work is also an important and initial start to guide study design regarding the selection of 498 499 samples, the genetic assay approach, the bioinformatics and the statistical 500 genetic analysis for COVID-19 as well as other infection and complex disease. 501 The publicly available summary statistics will encourage international 502 collaborative efforts to understand the host-pathogen interaction and to contain 503 the COVID-19 outbreak.

504 505

506 Acknowledgements

The study was supported by National Natural Science Foundation of China
(31900487), Guangdong Provincial Key Laboratory of Genome Read and
Write (No. 2017B030301011) and China National GeneBank (CNGB). We
would like to acknowledge Fan Zhang from Illumina, Zilong Li and Kang Fang,
Defu Xiao from BGI, Xinjun Zhang from University of California, Los Angeles,
Emilia Huerta-Sanchez from Brown University and Rasmus Nielsen from
University of UC Berkeley for helpful discussions of the results and advice.

515 Author contributions

Conceptualization, L.L, X.J., Q. H., S. Liu, S.H.; Methodology, J.S. S. Liu,
Y.Z., X.T.; Formal Analysis, S.H., Y.Z., X.Q., Zhi.L., P.L., Y.H., R.L., X.T., Y.B.,
S. Liu; Resources, F.W., R.G., C.L,W.X., Zhi.L, Q.T. R.C. X.L, X.Z., G.D.;
Data Curation, S. Liu, F.W., R.G., C.L; Writing – Original Draft, S. Liu; Writing
– Review & Editing, All; Supervision, X.X., J.W., H.Y.; Project Administration,
X.J., S.H., S. Liu, and F.C.; Funding Acquisition, F.W., X.J. and S. Liu.

522 523

524 Data availability

The data that support the findings of this study, including the allele frequency for the five groups of patients at all the 20 million detected genetic variants and the genome-wide association test summary statistics have been deposited in CNSA (China National Genebank Sequence Archive) in Shenzhen, China with accession number CNP0001107 (https://db.cngb.org/cnsa/).

530

532 Material and Methods

533 Patient recruitment and definition of phenotypes

A total of 332 patients were recruited from Jan 11th 2020 to Apr 2020 in 534 Shenzhen Third People's Hospital, the only referral hospital in Shenzhen City, 535 China³⁴. All were confirmed with SARS-COV-2 infection using real-time 536 reverse-transcriptase- polymerase-chain-reaction (RT-PCR) assay of nasal 537 and pharyngeal swab specimens. The demographic, epidemiological, clinical 538 and laboratory assessments were extracted from the electronic medical records 539 of the patients. This study was approved by the ethics commissions of the 540 Shenzhen Third People's Hospital Ethics Committee with a waiver of informed 541 consent. According to the 5th edition of the national treatment guideline of 542 COVID19 in China and the Chinese CDC criteria⁶, the patients were diagnosed 543 as asymptomatic, mild, moderate, severe and critically severe according to the 544 most severe stage they experienced during the disease course. The 545 asymptomatic, mild and the moderate groups of patients do not experience 546 pneumonia. When meeting any one of the following criteria, 1) RR>30 2) 547 Oxygen level < 93% 3) PaO2/FiO2 < 300 mmHg 4) disease progression greater 548 than 50% area in CT scan, a patient is categorized as severe patients. Patients 549 experienced one of the following 1) respiratory failure and requires mechanical 550 551 ventilation 2) shock 3) complicated by failure of other organs and requires intensive care monitoring were classified as critically severe. 552

553

554 Assignment of severity score to each patient

A machine learning XGBoost-based model was developed to predict ordinal 555 556 severity scores using patients' phenotype data of 64 laboratory test results⁶³. We first filtered out the laboratory test items of which at least 50% of patients 557 did not have any recordings. The remaining 52 laboratory test items with 558 missing values were further imputed by missForest algorithm⁶⁴. The missForest 559 is a nonparametric method to impute missing values using random forest model 560 561 in an iterative fashion. Then the originally ordered severity levels of asymptomatic, mild, moderate, severe and critical were assigned integer values 562 of 1, 2, 3, 4 and 5, respectively. The numeric representations retained the 563 ordinal levels of severity. We applied the reduction framework mentioned in Li 564 et al⁶⁵, where the ordinal regression was reduced to binary classification. The 565 reduction framework of extended binary classification was then integrated 566 within XGBoost model. Moreover, we selected the most predictive laboratory 567 test items using SHAP (SHapley Additive exPlanations) algorithm⁶⁶. The SHAP 568 is a game theoretic approach to explain the output of a given machine learning 569 model using Shapley values from game theory and their related extensions. We 570 finally trained the XGBoost-based ordinal regression model using the selected 571 laboratory test items. As a result, the prediction outcome produced by the final 572 model was typically a real number reflecting severity level that was used in the 573 downstream analysis. We used 100 base estimators for missForest, maximum 574 575 iteration of 10, and the criterion was mean squared error. For the XGBoostbased ordinal regression model, we used 500 base estimators and learn rate of 576

577 0.5. In general, the hyper-parameters of models in this study were chosen by 578 combining grid search of 5-fold cross validation and manual tuning.

579

580 **DNA extraction, library construction and deep whole genome** 581 **sequencing**

582 Genomic DNA was extracted from frozen blood samples of the 332 patients 583 using Magnetic Beads Blood Genomic DNA Extraction Kit (MGI, Shenzhen, 584 China). At least 0.5µg was obtained for each individual and used to create WGS 585 library, which insert sizes 300-500bp for paired-end libraries according to the 586 BGI library preparation pipeline. Sequencing was conducted on the DNBSEQ 587 platform (MGI, Shenzhen, China) to generate 100bp paired-end reads.

588

589 Genome alignment and variant detection

590 We used Sentieon Genomics software (version: sentieon-genomics-201911) to perform genome alignment and variant detection⁶⁷. Analysis pipeline were built 591 according to the recommendation in the Broad institute best practices described 592 593 in https://gatk.broadinstitute.org/hc/en-us/sections/360007226651-Best-594 Practices-Workflows. Sequencing reads were mapped to hg38 reference 595 genome using BWA algorithm. For each sample, after remove duplicates, Indel 596 realignment and base quality score recalibration (BQSR), SNP and short Indel 597 variants were detect using the Sentieon Haplotyper algorithm with option --598 emit mode gvcf to generate an individual GVCF file. Then the GVCF files for 599 all samples were subjected to Sentieon GVCFtyper algorithm to perform joint 600 variant calling.

601

602 Variant Quality Score Recalibration and Filtration

Variant Quality Score Recalibration were perform using Genome Analysis 603 Toolkit (GATK version 4.1.2). Known variant files were downloaded from the 604 605 GATK bundle. For indel recalibration. used we Mills and 1000G gold standard indels as the positive training and true set. 606 SNP recalibration, we used hapmap 3.3, 1000G omni2.5, 607 For and 1000G_phase1.snps as positive training sets, hapmap_3.3 as true set, and 608 dbSNP v146 as the known set. The metrics DP, QD, MQRankSum, 609 ReadPosRankSum, FS, SOR were used in the recalibration process. The truth-610 sensitivity-filter-level were set to 99.0 for both the SNPs and the Indels. Finally, 611 variants with quality score \geq 100 were selected for further analysis. 612

613

614 Familial relationship and population structure analysis

PLINK (v1.9)68 and KING (v2.1.5)69 was applied to detect the kinship 615 relatedness between each pair of the individuals. 48 patients from 16 families 616 were detected as related to each other. For several allele frequency-based 617 approach, we exclude the related patients and thus the sample size was 618 restricted to 284. PCA was performed using a subset of autosomal bi-allelic 619 SNPs on the unrelated patients using PLINK (v1.9). The PC-AiR module 620 (Principal components analysis in related samples) in the Genesis R package 621 was used to conduct PCA analysis for the 332 patients including the related 622

family members. Several restrictions were employed to select the final 614,963 SNPs for PCA analysis, including minor allele frequency (MAF) \ge 1% (common and low-frequency variants), genotyping rate \ge 90%, Hardy-Weinberg-Equilibrium (HWE) P > 0.000001, and removing one SNP from each pair with r2 \ge 0.5 (in windows of 50 SNPs with steps of 5 SNPs).

628

629 Genotype-phenotype association analysis

We have applied both the rvtest⁷⁰ and the SAIGE⁷¹ approaches to carry out 630 logistic regression, linear regression, burden test, the sequence kernel analysis 631 test (SKAT) and the optimal SKAT-O algorithm for the genotype-phenotype 632 association tests using the default parameters. For all the association tests, we 633 have used the gender, the age and the top 20 principal components from the 634 principal component analysis as the covariates. Exception is for the GWAS 635 between the 1KGP and the COVID-19 patients as age is not available for the 636 1KGP data set. Independent loci were defined as significant variants clustered 637 in a 1Mbp window. The lead SNP was defined as the SNP in the 1Mbp window 638 that has most significant, i.e., smallest p value. The genomic inflation factor, GC 639 lambda, attenuation ratio, LD score regression intercept and the SNP heritability 640 were estimated using the LD score regression approach ⁷². The qqman R 641 642 package was applied to generate the manhattan and gaplot. We defined genome-wide significance for single variant association test as 5e-8, suggestive 643 significance as 1e-5 and for gene-based association test as 1e-6. 644

645

646 HLA typing

When performing HLA typing, we first extracted reads which aligned to HLA region of GRCh38 and unmapped reads from individual bam files. Then using xHLA algorithm²³ typing HLA class I(A B C gene) and II(DRB1 DQB1 DPB1) genes. DQA1 gene was typed using SOAP-HLA algorithm⁵³ for xHLA does not include this gene. We performed the association analysis between HLA types and the binary severe and mild groups using PLINK (version 1.90) using a logistic regression model, adjusted for age, gender and top 20 PCs.

- 654
- 655 656

657 **Reference**

- Johnson, N. P. A. S. & Mueller, J. Updating the accounts: global mortality of
 the 1918-1920 'Spanish' influenza pandemic. *Bull. Hist. Med.* (2002).
- 660 doi:10.1353/bhm.2002.0022
- 2. Zhu, N. *et al.* A novel coronavirus from patients with pneumonia in China,
 2019. *N. Engl. J. Med.* (2020). doi:10.1056/NEJMoa2001017
- John Hopkins University and Medicine. COVID-19 Map Johns Hopkins
 Coronavirus Resource Center. John Hopkins Coronavirus Resource Center
 (2020).
- 4. Jiang, S., Du, L. & Shi, Z. An emerging coronavirus causing pneumonia
 outbreak in Wuhan, China: calling for developing therapeutic and prophylactic
 strategies. *Emerging Microbes and Infections* (2020).
 doi:10.1080/22221751.2020.1723441

670	5.	Shi, Z. & Hu, Z. A review of studies on animal reservoirs of the SARS
671		coronavirus. Virus Res. (2008). doi:10.1016/j.virusres.2007.03.012
672	6.	Wu, Z. & McGoogan, J. M. Characteristics of and Important Lessons From the
673		Coronavirus Disease 2019 (COVID-19) Outbreak in China. JAMA (2020).
674		doi:10.1001/jama.2020.2648
675	7.	Kenney, A. D. et al. Human Genetic Determinants of Viral Diseases. Annu.
676		<i>Rev. Genet.</i> (2017). doi:10.1146/annurev-genet-120116-023425
677	8.	Guan, W. et al. Clinical characteristics of coronavirus disease 2019 in China.
678		N. Engl. J. Med. (2020). doi:10.1056/NEJMoa2002032
679	9	Full, et al. Clinical characteristics of coronavirus disease 2019 (COVID-19) in
680	0.	China: A systematic review and meta-analysis ./ Infect (2020)
681		doi:10.1016/i jinf 2020.03.041
682	10	Guan W L et al. Clinical Characteristics of Coronavirus Disease 2019 in
683	10.	China N Engl I Med (2020) doi:10.1056/NE IMoa2002032
687	11	The enidemiological characteristics of an outbreak of 2019 novel coronavirus
604 605	11.	discassos (COVID 10) in China. Zhonghua Liu Ying Bing Yuo Za Zhi (2020)
605		diseases $(00010-19)$ in China. Zhonghua Liu Xing Ding Xue Za Zhi (2020). doi:10.3760/cmp i ison 0.254.6450.2020.02.003
000 607	10	C_{10} C at al Discregulation of immune response in patients with COVID 10 in
007	12.	Wuban China Clin Infact Dia (2020) dai:10.1002/aid/aiaa.249
000	12	Vang X at al. Clinical course and outcomes of critically ill patients with SAPS
600	13.	Cold 2 programming in Muhan Ching: a single contered retrospective
690 CO1		cov-2 prieumonia in wurian, China. a single-centereu, retrospective,
691		005ervalional sludy. Lancel Respir. Med. (2020). doi: 10.1016/52215-
692	11	2000(20)30079-3 Nichiura H. Rackaelaulating the Incidence of Infaction with COV/ID 10 on the
693	14.	Dismond Dringson L Clin Mad (2020) doi:10.2200/jem0020657
694 COF	15	Diamond Princess. J. Clin. Med. (2020). doi: 10.3390/jCm9030657
695	15.	Hu, Z. <i>et al.</i> Clinical characteristics of 24 asymptomatic infections with COVID-
696		19 screened among close contacts in Nanjing, China. Sci. China Life Sci.
697	40	(2020). doi: 10.1007/S11427-020-1661-4
698	16.	Feliay, J. et al. Common genetic variation and the control of HIV-1 in numans.
699	47	PLOS Genet. (2009). doi:10.1371/journal.pgen.1000/91
700	17.	Feliay, J. <i>et al.</i> A whole-genome association study of major determinants for
/01	4.0	nost control of HIV-1. Science (80). (2007). doi:10.1126/science.1143/67
702	18.	Pereyra, F. et al. The major genetic determinants of HIV-1 control affect HLA
703		class I peptide presentation. Science (80). (2010).
704	4.0	doi:10.1126/science.11952/1
705	19.	Kamatani, Y. et al. A genome-wide association study identifies variants in the
706		HLA-DP locus associated with chronic hepatitis B in Asians. Nat. Genet.
707		(2009). doi:10.1038/ng.348
708	20.	Ge, D. <i>et al.</i> Genetic variation in IL28B predicts hepatitis C treatment-induced
709		viral clearance. Nature (2009). doi:10.1038/nature08309
710	21.	Schulert, G. S. et al. Whole-exome sequencing reveals mutations in genes
711		linked to hemophagocytic lymphohistiocytosis and macrophage activation
712		syndrome in fatal cases of H1N1 influenza. <i>J. Infect. Dis.</i> (2016).
713		doi:10.1093/infdis/jiv550
714	22.	Wang, Z. et al. Early hypercytokinemia is associated with interferon-induced
715		transmembrane protein-3 dysfunction and predictive of fatal H7N9 infection.
716		Proc. Natl. Acad. Sci. U. S. A. (2014). doi:10.1073/pnas.1321748111
717	23.	Everitt, A. R. et al. IFITM3 restricts the morbidity and mortality associated with
718		influenza. <i>Nature</i> (2012). doi:10.1038/nature10921
719	24.	Yang, X. et al. Interferon-inducible transmembrane protein 3 genetic variant

720 721		rs12252 and influenza susceptibility and severity: A meta-analysis. <i>PLoS One</i> (2015). doi:10.1371/journal.pone.0124985
722 723	25.	Lin, M. <i>et al.</i> Association of HLA class I with severe acute respiratory syndrome coronavirus infection. <i>BMC Med. Genet.</i> (2003). doi:10.1186/1471-
724		2350-4-9
725	26.	Ching, J. C. <i>et al.</i> Significance of the Myxovirus Resistance A (MxA) Gene
726		-123C>A Single-Nucleotide Polymorphism in Suppressed Interferon B
727		Induction of Severe Acute Respiratory Syndrome Coronavirus Infection, J.
728		Infect. Dis. 201 , 1899–1908 (2010).
729	27.	Kachuri, L. <i>et al.</i> The landscape of host genetic factors involved in infection to
730		common viruses and SARS-CoV-2. <i>medRxiv</i> (2020).
731		doi:10.1101/2020.05.01.20088054
732	28.	Williams, F. M. et al. Self-reported symptoms of covid-19 including symptoms
733		most predictive of SARS-CoV-2 infection, are heritable, <i>medRxiv</i> (2020).
734		doi:10.1101/2020.04.22.20072124
735	29.	Zhao, J. et al. Relationship between the ABO Blood Group and the COVID-19
736		Susceptibility. <i>medRxiv</i> (2020). doi:10.1101/2020.03.11.20031096
737	30.	Cao. Y. et al. Comparative genetic analysis of the novel coronavirus (2019-
738		nCoV/SARS-CoV-2) receptor ACE2 in different populations. Cell Discov. 6. 4–
739		7 (2020).
740	31.	Bhattacharyya, C. et al. Global Spread of SARS-CoV-2 Subtype with Spike
741		Protein Mutation D614G is Shaped by Human Genomic Variations that
742		Regulate Expression of TMPRSS2 and MX1 Genes. <i>bioRxiv</i> (2020).
743		doi:10.1101/2020.05.04.075911
744	32.	Renieri, A. et al. ACE2 variants underlie interindividual variability and
745		susceptibility to COVID-19 in Italian population. medRxiv (2020).
746		doi:10.1101/2020.04.03.20047977
747	33.	The COVID-19 Host Genetics Initiative, a global initiative to elucidate the role
748		of host genetic factors in susceptibility and severity of the SARS-CoV-2 virus
749		pandemic. Eur. J. Hum. Genet. (2020). doi:10.1038/s41431-020-0636-6
750	34.	Cai, Q. et al. COVID-19 in a designated infectious diseases hospital outside
751		Hubei Province, China. Allergy Eur. J. Allergy Clin. Immunol. (2020).
752		doi:10.1111/all.14309
753	35.	Darbeheshti, F. & Rezaei, N. Genetic predisposition models to COVID-19
754		infection. Medical Hypotheses (2020). doi:10.1016/j.mehy.2020.109818
755	36.	Auton, A. et al. A global reference for human genetic variation. Nature 526,
756		68–74 (2015).
757	37.	Lundberg, S. M. et al. From local explanations to global understanding with
758		explainable AI for trees. Nat. Mach. Intell. (2020). doi:10.1038/s42256-019-
759		0138-9
760	38.	Jiang, X. et al. Towards an Artificial Intelligence Framework for Data-Driven
761		Prediction of Coronavirus Clinical Severity. Comput. Mater. Contin. (2020).
762		doi:10.32604/cmc.2020.010691
763	39.	Depristo, M. A. et al. A framework for variation discovery and genotyping using
764		next-generation DNA sequencing data. Nat. Genet. (2011). doi:10.1038/ng.806
765	40.	McLaren, W. et al. The Ensembl Variant Effect Predictor. Genome Biol. (2016).
766		doi:10.1186/s13059-016-0974-4
767	41.	Karczewski, K. J. et al. Variation across 141,456 human exomes and genomes
768		reveals the spectrum of loss-of-function intolerance across human protein-
769		coding genes. <i>bioRxiv</i> (2019). doi:10.1101/531210

770	42.	Collins, R. L. et al. An open resource of structural variation for medical and
771		population genetics. <i>bioRxiv</i> (2019). doi:10.1101/578674
772	43.	Dai, W. et al. Whole-exome sequencing identifies MST1R as a genetic
773		susceptibility gene in nasopharyngeal carcinoma. Proc. Natl. Acad. Sci. U. S.
774		<i>A.</i> (2016). doi:10.1073/pnas.1523436113
775	44.	Sharma, S. et al. ACE2 Homo-dimerization, Human Genomic variants and
776		Interaction of Host Proteins Explain High Population Specific Differences in
777		Outcomes of COVID19. <i>bioRxiv</i> (2020). doi:10.1101/2020.04.24.050534
778	45.	Stelzer, G. et al. The GeneCards suite: From gene data mining to disease
779		genome sequence analyses. Curr. Protoc. Bioinforma. (2016).
780		doi:10.1002/cpbi.5
781	46.	Pertel, T. et al. TRIM5 is an innate immune sensor for the retrovirus capsid
782		lattice. Nature (2011). doi:10.1038/nature09976
783	47.	Shi, Y. et al. COVID-19 infection: the perspectives on immune responses. Cell
784		Death and Differentiation (2020), doi:10.1038/s41418-020-0530-3
785	48.	Cavalli, G. <i>et al.</i> Interleukin-1 blockade with high-dose anakinra in patients with
786		COVID-19, acute respiratory distress syndrome, and hyperinflammation; a
787		retrospective cohort study. <i>Lancet Rheumatol.</i> (2020), doi:10.1016/s2665-
788		9913(20)30127-2
789	49	Buniello, A, et al. The NHGRI-FBI GWAS Catalog of published genome-wide
790		association studies targeted arrays and summary statistics 2019 Nucleic
791		Acids Res (2019) doi:10.1093/nar/gkv1120
792	50	Hammer C. et al. Amino acid variation in HI A class II proteins is a major
793	00.	determinant of humoral response to common viruses Am J Hum Genet
79/		(2015) doi:10.1016/i.aiba.2015.09.008
795	51	Tian C. et al. Genome-wide association and HI A region fine-manning studies
796	01.	identify susceptibility loci for multiple common infections. Nat. Commun
797		(2017) doi:10.1038/s41467-017-00257-5
798	52	Xie C. et al. Fast and accurate HI A typing from short-read next-generation
799	02.	sequence data with xHLA Proc Natl Acad Sci U.S. A (2017)
800		doi:10.1073/pnas.1707945114
200 201	53	Cao, H. et al. An Integrated Tool to Study MHC Region: Accurate SNV
801 801	55.	Detection and HLA Conce Typing in Human MHC Region Using Targeted
00Z 002		High Throughout Sequencing <i>PL</i> os One (2012)
005		doi:10.1271/journal.nono.0060288
004 905	54	Conzeloz Colorza E E ot al Allele frequency net detabase in Methods in
805 806	54.	Mologular Piology (2018) doi:10.1007/078.1.4020.8546.2.4
000 907	55	Novecular Biology (2018). doi: 10.1007/976-1-4959-6546-5_4
807	55.	Nguyen, A. et al. numan leukocyte antigen susceptibility map for SARS-COV-
808	FG	2. J. VII'01. (2020). 001. 10. 1120/JVI.00510-20
809	50.	Lai, S. <i>et al.</i> Effect of non-pharmaceutical interventions to contain COVID-19 in China. Nature (2020). dai:40.4020/s44500.020.2003.v
810	67	China. Nature (2020). doi: 10.1038/s41586-020-2293-x
811	57.	Ju, B. et al. Potent numan neutralizing antibodies elicited by SARS-CoV-2
812	50	Infection. <i>bioRxiv</i> 2020.03.21.990770 (2020). doi:10.1101/2020.03.21.990770
813	58.	Rashkin, S., Jun, G., Chen, S. & Abecasis, G. R. Optimal sequencing
814		strategies for identifying disease-associated singletons. <i>PLoS Genet.</i> (2017).
815	F 0	doi:10.1371/journal.pgen.1006811
816	59.	vvainschtein, P. et al. Recovery of trait heritability from whole genome
817	00	sequence data. Yearb. Paediatr. Endocrinol. (2019). doi:10.1530/ey.16.14.15
818	60.	Carithers, L. J. & Moore, H. M. The Genotype-Tissue Expression (GTEx)
819		Project. Biopreservation and Biobanking (2015).

820		doi:10.1089/bio.2015.29031.hmm
821	61.	Giza, M., Koftori, D., Chen, L. & Bowness, P. Is Behcet's disease a 'class 1-
822	• • • •	opathy'? The role of HLA-B*51 in the pathogenesis of Behcet's disease.
823		Clinical and Experimental Immunology (2018), doi:10.1111/cei.13049
824	62.	Nishiura, H. <i>et al.</i> Estimation of the asymptomatic ratio of novel coronavirus
825	•=-	infections (COVID-19). International Journal of Infectious Diseases (2020).
826		doi:10.1016/i.iiid.2020.03.020
827	63	Chen T & Guestrin C XGBoost: A scalable tree boosting system in
828	00.	Proceedings of the ACM SIGKDD International Conference on Knowledge
829		Discovery and Data Mining (2016), doi:10.1145/2939672.2939785
830	64	Stekhoven D. I. & Bühlmann P. Missforest-Non-narametric missing value
831	04.	imputation for mixed-type data. <i>Bioinformatics</i> (2012)
832		doi:10.1093/bioinformatics/btr597
833	65	Li L & Lin H T Ordinal regression by extended binary classification in
831	00.	Advances in Neural Information Processing Systems (2007)
034 025		doi:10.7551/mitpress/7503.003.0113
83C 022	66	Lundhera S M & Lee S I A unified approach to interpreting model
030	00.	prodictions in Advances in Noural Information Processing Systems (2017)
027	67	Freed D. N. Aldena D. Weber, J. A. & Edwards, J. S. The Sentioon
020	07.	Conomics Tools A fast and accurate solution to variant calling from next
039		constant calling from the data bio Dviv (2017)
040	60	Chang C. C. at al Second generation DI INK: riging to the challenge of larger
841 842	00.	and risher deteasts. Circassiones 4, 7 (2015)
84Z	60	Manichaikul A at al Debust relationship inforence in genome wide
045	09.	manicifiation, A. et al. Robust relationship interence in genome-wide
844 045	70	Zhan X, Hu X, Li P, Abagagia C, P, & Liu D, L DV/TESTS: An afficient
845 84C	70.	Zildii, A., Hu, T., Li, D., Abecasis, G. R. & Liu, D. J. RVTESTS. All efficient
840		and comprehensive tool for fare variant association analysis using sequence
847	74	data. <i>Bioinionnatics</i> (2016). doi:10.1095/bioinionnatics/biw079
848	71.	Zhou, W. <i>et al.</i> Efficiently controlling for case-control imbalance and sample
849		dei:10.1028/c41588.018.0184.v
850	70	001.10.1038/S41588-018-0184-y
851	12.	buik-Sumvan, B. K. et al. LD Score regression distinguishes comounding from
852		polygenicity in genome-wide association studies. <i>Nat. Genet.</i> (2015).
853		dol: 10. 1038/ng.3211
854		
855		
856		
857		
858		
859		
860		
861		
862		
863		
864		





Figure 1. Clinical and laboratory assessments of the recruited 332 COVID-19 867 868 patients. (A) number of samples belong to the five categories (B) top 20 features that 869 classify the patient categories in the machine learning trained model (C) age distribution for the five categories of patients (D) distribution of disease duration, i.e. 870 the duration between the disease onset and the first negative RT-PCR test among the 871 five groups of patients (E) gender distribution for the five categories of patients by age 872 (F) distribution of the proportion of patients with or without medical comorbidities 873 among the five categories of patients by age. 874 875





Figure 2. Deep whole genome sequencing and genetic variation among the 880 patients. A) Sequencing depth distribution B) Proportions and numbers of types (SNP, 881 Indel) of genetic variants identified from the patients C) Proportions and numbers of 882 functional consequences of the genetic variants among the patients D) comparison of 883 loss of function variation burden for SNP, small insertions and deletions between the 884 severe and the non-severe patients E) Single variation association test for loss of 885 886 function mutation burden between the severe and non-severe patients F) allele frequency distribution for all the missense and loss of function variants present in 887 ACE2 and TRMPSS2 genes. 888

- 889
- 890



association test for three severity traits. (A) Severe and critical severe groups versus
 the rest of the non-severe groups. (B) Severity score assessed by laboratory test
 measurements. (C) the duration from disease onset to recovery (D)-(F) Gene-based
 association test for three traits.

917



919

920 Figure 4. LD, allele frequency and pleiotropic effects of the TMEM189-UBE2V1 921 signal suggestively associated with COVID-19 patient severity. A) Locuszoom plot shows the p-value of the SNPs centering the lead SNP rs6020298 and the 922 recombination rate. Color of the dots indicate linkage disequilibrium r2 metric. B) Allele 923 924 frequency of s6020298 among the 1000 genomes populations. The allele frequency of the reference and alternative allele is visualized by the geography of genetic 925 variants browser developed by the university of Chicago. C) P-value of the single 926 variant genome-wide association test for the sixty-four laboratory assessments at the 927 lead SNP rs6020298. The P-value of the three traits (Severity, Severity score and 928 929 Disease Duration) in Figure 3 were also displayed.

- 930
- 931 932
- 933
- 934
- 935
- - -



Figure 5. Human leukocyte haplotype allele frequency between severe vs non severe groups. Comparison for class I HLA genes (top). Comparison for class II HLA
 genes (bottom). Star indicates significance level in a logistic regression on the allele
 frequency with age, gender and the top twenty principal components as covariates.
 *<0.05, **<0.01



959 960

Figure 6. Single variant and gene-based association test between COVID-19 patients and the general populations. (A) single variant association test and (B) gene-based association test between the unrelated COVID-19 patients (N=284) and the 1KGP Chinese population (N=301) (C) single variant association test and (D) gene-based association test between the unrelated COVID-19 patients (N=284) and the CNRP Chinese population (N=665). Only variants with moderate or high impacts by variant effect predictor were shown in (A) and (C).

968

969

970

971

972

974	
975	

Table 1. Nominal association of HLA allele and severity by logistic rgression

	Severe	Non-severe	OR	SE	Р
C*14:02	0.086	0.047	4.75	0.52	0.003028
B*51:01	0.101	0.058	3.38	0.45	0.007017
A*11:01	0.297	0.263	2.33	0.32	0.008512
DRB1*14:04	0.029	0.005	15.1	1.06	0.01027
DRB1*01:01	0.022	0.005	13.7	1.13	0.02034
DPB1*03:01	0.008	0.044	0.09	1.15	0.03669
DQA1*01:01	0.029	0.009	6.05	0.87	0.03947
DRB1*12:01	0.022	0.037	0.18	0.87	0.04478
B*13:02	0.058	0.051	0.27	0.66	0.04935

976 977

978 Severe group indicates severe and critical patients

Non-severe group includes asymptomatic, mild and moderate patients