

ARTICLE

Open Access

Initial whole-genome sequencing and analysis of the host genetic contribution to COVID-19 severity and susceptibility

Fang Wang¹, Shujia Huang^{2,3}, Rongsui Gao¹, Yuwen Zhou^{2,4}, Changxiang Lai¹, Zhichao Li^{2,4}, Wenjie Xian¹, Xiaobo Qian^{2,4}, Zhiyu Li¹, Yushan Huang^{2,4}, Qiyuan Tang¹, Panhong Liu^{2,4}, Ruikun Chen¹, Rong Liu², Xuan Li¹, Xin Tong^{1,2}, Xuan Zhou¹, Yong Bai^{1,2}, Gang Duan¹, Tao Zhang^{1,2}, Xun Xu^{1,2,5}, Jian Wang^{2,6}, Huanming Yang^{2,6}, Siyang Liu², Qing He¹, Xin Jin^{2,3} and Lei Liu¹

Abstract

The COVID-19 pandemic has accounted for millions of infections and hundreds of thousand deaths worldwide in a short-time period. The patients demonstrate a great diversity in clinical and laboratory manifestations and disease severity. Nonetheless, little is known about the host genetic contribution to the observed interindividual phenotypic variability. Here, we report the first host genetic study in the Chinese population by deeply sequencing and analyzing 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, we conducted both single-variant and gene-based association tests among five severity groups including asymptomatic, mild, moderate, severe, and critical ill patients after the correction of potential confounding factors. Pedigree analysis suggested a potential monogenic effect of loss of function variants in *GOLGA3* and *DPP7* for critically ill and asymptomatic disease demonstration. Genome-wide association study suggests the most significant gene locus associated with severity were located in *TMEM189-UBE2V1* that involved in the IL-1 signaling pathway. The p.Val197Met missense variant that affects the stability of the TMPRSS2 protein displays a decreasing allele frequency among the severe patients compared to the mild and the general population. We 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 genomic study of Chinese patients provides genetic insights into the phenotypic difference among the COVID-19 patient groups and highlighted genes and variants that may help guide targeted efforts in containing the outbreak. Limitations and advantages of the study were also reviewed to guide future international efforts on elucidating the genetic architecture of host–pathogen interaction for COVID-19 and other infectious and complex diseases.

Introduction

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. Since the late December of 2019, the 2019 novel coronavirus diseases (COVID-19) has spread rapidly throughout the world, resulting in more than five million confirmed cases and hundreds of thousands deaths in less than 6 months^{2,3}. The disease was caused by the infection of a novel enveloped RNA betacoronavirus that has been

Correspondence: Siyang Liu (liusiyang@genomics.cn) or Qing He (heqingjoe@163.com) or Xin Jin (jinxin@genomics.cn) or Lei Liu (liuleiszsdsrmyy@163.com)

¹The Third People's Hospital of Shenzhen, National Clinical Research Center for Infectious Disease, The Second Affiliated Hospital of Southern University of Science and Technology, Shenzhen, Guangdong 518112, China

²BGI-Shenzhen, Shenzhen, Guangdong 518083, China

Full list of author information is available at the end of the article

These authors contributed equally: Fang Wang, Shujia Huang, Rongsui Gao, Yuwen Zhou, Changxiang Lai, Zhichao Li

© The Author(s) 2020



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

named severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), which is the seventh coronavirus species that causes respiratory disease in humans^{4,5}. The virus causes serious respiratory illnesses such as pneumonia, lung failure, and even death⁶. Until now, there is no specific therapeutics and vaccine available for its control. Continuing epidemiological and molecular biological study to better understand, treat and prevent COVID-19 are urgently needed.

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 interindividual 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., non-pneumonia and mild pneumonia), while 15%–20% would progress to severe or critical stage with a high probability of respiratory failure^{8–10}. Patients with severe disease had more prominent laboratory abnormalities including lymphocytopenia and leukopenia than those with non-severe disease^{11,12}. In addition, not all people exposed to SARS-CoV-2 were infected according to the epidemiological observation of the patients' close contacts^{13,14}. Notably, previous studies have indicated that genetic background plays an essential role in determining the host responses to infections by HIV^{15–17}, HBV¹⁸, HCV¹⁹, influenza^{20–23}, SARS-CoV^{24,25}, numerous common viruses²⁶, etc. Those studies highlighted the HLA alleles and several genes involved in the interferon production and viral replication pathway and indicates that genetic factors may also play an important role to explain the interindividual clinical variability among patients infected by SARS-CoV-2.

Till now, the global genetic community has been actively investigating in the genetic contribution to COVID-19. A recent twin study in UK suggests a 30%–50% genetic heritability for self-reported symptoms of COVID-19 and the predictive disease onset²⁷, indicating a very strong genetic background predisposing the COVID-19 patients' clinical manifestation and susceptibility. An earlier studies comparing the distribution of ABO blood group from 1775 patients infected with SARS-CoV-2 with 3694 normal people from Wuhan city and 23,386 people from Shenzhen city suggested that blood group A had a significantly higher risk for COVID-19 (OR = 1.20, $p = 0.02$) while blood group O had the lower risk²⁸. Using allele frequency and expression quantitative loci (eQTL) information of general healthy population from 1000 genome project and others, a few studies investigate the mutation frequency spectrum in different populations in candidate genes such as *ACE2* and *TRMPSS2*^{29–31}. Genome-wide association test on array

data from the UK Biobank participants with a positive and negative polymerase chain reaction (PCR)-tests also reveals a few suggestive genes²⁶. The COVID-19 host genetics initiative was established to encourage generation, sharing and meta-analysis of the genome-wide association summary statistics data around the world³². International collaborative efforts are necessary to elucidate the role of host genetic factors defining the severity and susceptibility of the SARS-CoV-2 virus pandemic.

Herein, we report the first genetic study of COVID-19 disease severity in China by deeply analyzing the association between the genetic variants present in the patients' genome and their disease progression. We have recruited 332 hospitalized patients from a designated infectious disease hospital in Shenzhen City³³. The patients display varying clinical and laboratory features and were categorized as asymptomatic, mild, moderate, severe, and critical cases according to the criteria made by the Chinese Center for Disease Control and Prevention⁶. To maximize the statistical power given the relatively small hospitalized sample size and for accurate detection of extremely rare variants, we conducted deep whole genome sequencing (average 46×) 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 variability³⁴. 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 (GWAS) and by evaluating the difference of allele frequency of the protein truncating variants and HLA alleles among the patient groups. In addition, we performed joint-calling of the genetic variants of the unrelated COVID-19 patients ($n = 284$) and the publicly available Chinese genomes from the 1000 genome project³⁵ ($n = 301$, ~7×) and 665 selected Chinese genomes from the Chinese Reference Panel Population (manuscript in preparation, ~30×) to explore potential genetic factors that may contribute genetic susceptibility of SARS-CoV-2 infection.

Results

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

The 332 recruited patients with laboratory-confirmation of SARS-COV-2 infection were quarantined and treated in the Shenzhen Third Hospital. We extracted and analyzed the clinical symptoms, laboratory assessment and 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.

A total of 25 (7.5%), 12 (3.6%), 225 (67.8%), 53 (16.0%), and 17 (5.1%) patients were defined as asymptomatic, mild,

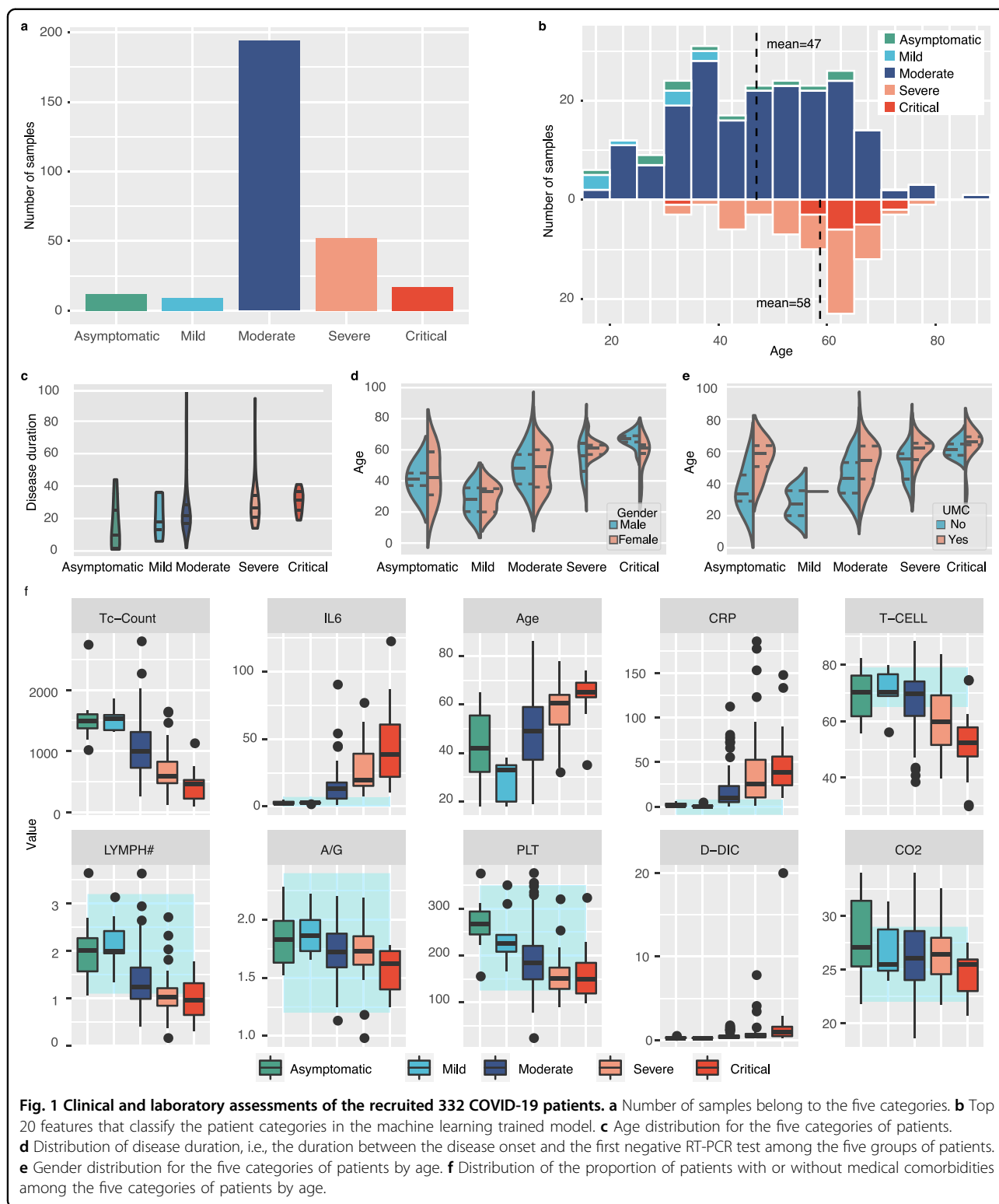


Fig. 1 Clinical and laboratory assessments of the recruited 332 COVID-19 patients. **a** Number of samples belong to the five categories. **b** Top 20 features that classify the patient categories in the machine learning trained model. **c** Age distribution for the five categories of patients. **d** Distribution of disease duration, i.e., the duration between the disease onset and the first negative RT-PCR test among the five groups of patients. **e** Gender distribution for the five categories of patients by age. **f** Distribution of the proportion of patients with or without medical comorbidities among the five categories of patients by age.

moderate, severe, and critically ill, respectively, according to the most severe stage they encountered during the disease course following the Chinese CDC criteria⁶ (Fig. 1a). The criteria for the severity based on the clinical manifestation

were detailed in the section “Materials and methods”. A broader definition of the mild group includes the asymptomatic, mild and moderate patients, and the severe group includes the severe and critically ill patients.

The patients displayed several clinical presentations typical to COVID-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%) (Supplementary Fig. S1). Less than 10% of the patients had also experienced diarrhea, chest and abdominal pain, shortness of breath and anosmia. More than 50% of the patients had at least one medical comorbidities (e.g., hypertension). Consistent with previous report, the broadly defined severe patients tend to be older (severe: average 45 years old vs. mild average 58 years old, t test $p = 0.03$, Fig. 1b), suffer from a longer course of disease between the onset and the first negative reverse-transcriptase (RT)-PCR test outcome (Fig. 1c) and shorter exposure time (Supplementary Fig. S2). In addition, the severe patient group consist of more males than females (severe 66.7% vs. mild 41.3%, χ^2 test $p = 4.3e-4$, Fig. 1d) and tend to undergo medical comorbidities more frequently (severe 58.8% vs. mild 45.1%, χ^2 test $p = 0.07$, Fig. 1e) than the mild patients.

During hospitalization, a series of 64 laboratory assessments including a complete blood count and blood chemical analysis, assessment of liver function, assessment of renal functions, test of humoral immunity, test of coagulation, measure of electrolyte, and measure of blood gas electrolyte (Supplementary Fig. S3) and a time-series evaluation of T lymphocyte subgroups (Supplementary Fig. S4) were performed for each of the patients to monitor their disease status and progression. Using a tree-based machine learning prediction model³⁶, we computed the local interaction effects of the 64 laboratory assessment features as well as three demographic features including age, gender, and w/o medical comorbidities for classification of the patient severity category (Supplementary Fig. S5). The top ten features of greatest importance that contribute to a severer disease outcome include decreased lymphocyte counts (Tc-Count, T-CELL, LYMPH#) and platelet counts, evaluated interleukin 6, C-reactive protein and D-dimer, increased age and decreased A/G and CO2 (Fig. 1f), consistent with previous reports³⁷. We applied the top 20 features of importance to assign a severity score for each patient to reflect their disease status (Supplementary Fig. S6).

Deep whole-genome sequencing and genetic variants identified

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 $46.26\times$ vs. severe $46.71\times$) (Fig. 2a). We conducted variation detection and genotyping using the GATK joint genotyping framework to avoid potential batch effect derived from individual variant calling. Bioinformatics

analysis and the data quality control process were described in details in the section “Material and methods”.

Among the 332 patients, we identified a total of 22.2 million variants including 17.9 million biallelic single nucleotide polymorphism, 1.75 million biallelic small insertions and deletions, and 2.49 million multiallelic variants (Fig. 2b, Supplementary Table S1). The average transition/transversion (ts/tv) ratio is 2.12 and the proportion of heterozygous vs. homozygous variants among all the samples is 1.29, consistent with the statistical expectation³⁸ and indicates good quality of the variant calls (Supplementary Fig. S7). Among the 284 unrelated patients, we have identified 398 K variants that result in an alteration of the protein coding sequence, 5147 of which were predicted loss of function variants that included 1973 frameshift, 1528 stop-gained, 954 splice donor, and 692 splice acceptor variants (Fig. 2c). Totally, 261 of those variants were uniquely presented in the COVID-19 patients (5.07%) and have not been previously reported in the 1000 genome and the gnomAD studies^{35,39,40}. On average, each patient possessed 343 predicted loss of function variants in their genome (Supplementary Fig. S8). We evaluated whether loss of function variation might enrich in severe patients in the following analysis.

Potential monogenic genetic effects using the pedigree and population strategies

Our first question was whether there might be monogenic cause for the young but critically ill patients, or on the other side, the old but asymptomatic patients. We tried to tackle this question using both the family and the population strategy. There were in total 35 pedigree families involved in the study. Their disease severity was positively correlated with their ages where older people tend to experience severer disease progression (Supplementary Table S2). Nonetheless, there were two families consisting of patients that did not follow the trend (Supplementary Tables S2 and S3). Family KING8 was composed of four members (Fig. 2d). Patient 2780 was a 35-year-old woman without previous comorbidity when she was infected by SARS-CoV-2 and progressed in critically severe COVID-19. Her 6.75-year-old daughter (Patient 2822) was an asymptomatic patient while her 61-year-old mother (Patient 4902) and her 34-year-old husband were moderate patients infected by the virus. We investigated the loss of function variants that were uniquely present in patient 2780 but not her mother and daughter and that were rare (<0.005) in the general population (1000 genome and the gnomAD). In total, there were five rare loss of function variants meeting the criteria (Supplementary Table S4). Among the five, four were uniquely presented as heterozygous genotypes in patient II while one SNV (rs143359233), resulting in a splice acceptor alternation in gene *GOLGA3*, were present

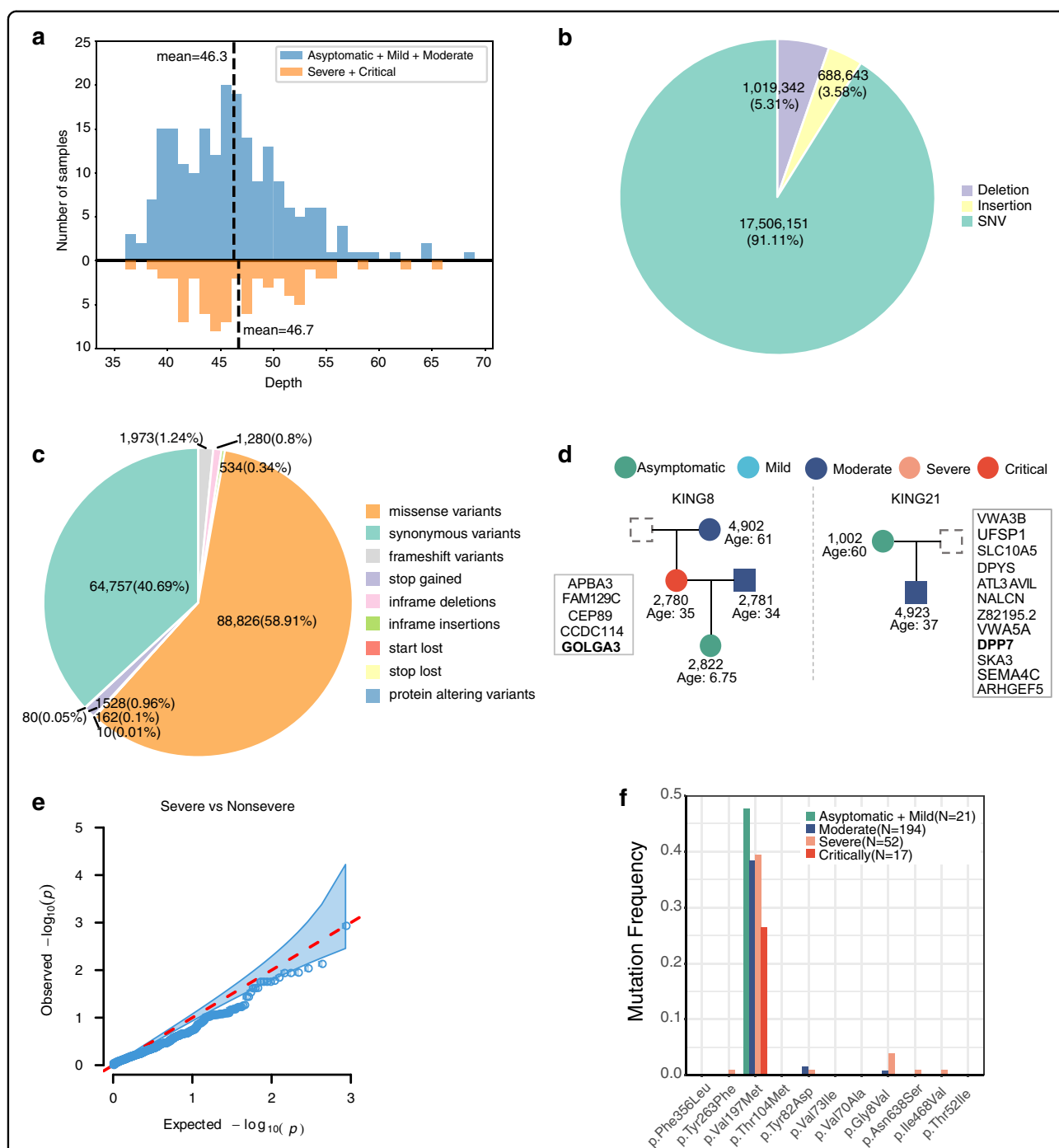


Fig. 2 Deep whole-genome sequencing and genetic variation among the patients. **a** Sequencing depth distribution. **b** Proportions and numbers of types (SNP, Indel) of genetic variants identified from the patients. **c** Proportions and numbers of functional consequences of the genetic variants among the patients. **d** Pedigree gene discovery. **GOLGA3** and **DPP7** were marked in bold for presence of rare recurrent loss of function variants. **e** Mutation burden association test for loss of function between the severe and non-severe patients. **f** Allele frequency distribution for all the missense and loss of function variants present in *ACE2* and *TRMPSS2* genes.

twice among the critically ill group of patients (the other critically ill patient is a 65-year-old male patient). In another noteworthy family KING21, a 60-year-old female patient (Patient 1002) remained to be asymptomatic

during the infection while her 37-year-old son displayed moderate symptoms (Patient 4923). There were eleven loss of function variants that were presented in patient 1002 but not in her moderately infected son and that were

absent in the severe patients (Supplementary Table S5). Particularly, a 1-bp insertion in gene *DPP7* was present twice in another 40-year-old female asymptomatic patient besides Patient 1002. This insertion has been known to destroy *DPP7*'s transcription in whole blood and several other tissues and organs according to GTEx Portal (Supplementary Fig. S9). Notably, the *DPP7* is a dipeptidyl peptidase that plays a role in innate immune system⁴¹ and another dipeptidyl peptidase 4 (*DPP4*) is the host receptor for the binding of the MERS-CoV envelope spike glycoprotein⁴². Nonetheless, due to lack of sampling of more family members, such as the father of Patient 2780 and the parents or siblings of Patient 1002, we were not able to resolve the causality of the candidate genetic variants to disease severity. Replication with more family patients with extreme phenotypes sequenced recently and in the near future around the globe will help elucidate the impact of those loss of function variation among the COVID-19 patients.

As an alternative strategy powered more by unrelated samples, we further investigated difference of genetic burden of loss of function variants between the severe and the mild groups of patients. The severe and the critical patients tend to have slightly more loss of function insertions and deletions than the asymptomatic, mild and the moderate groups ($p = 0.004$) in a logistic regression taking the number of loss of function variants as variable and the patients' age, gender, the 20 principle components and effective sequencing depth as covariates (Supplementary Fig. S10). When performing a mutation burden test for each of the 16,801 genes that have more than one variant among the 284 unrelated patients, we did not identify genes that were enriched in loss of function variants in the severe and critical patients (Fig. 2e). On the other hand, we found two heterozygous loss of function variants located in *MST1R* and *RASA2* that were only present in the asymptomatic and mild patients (Supplementary Fig. S11). The *MST1R* encodes the macrophage stimulating one receptor expressed on the ciliated epithelia of the mucociliary transport apparatus of the lung and follows an autosomal dominant inheritance mode for susceptibility to nasopharyngeal carcinoma⁴³. Nonetheless, because those loss of function variants were only present in one patient among the mild groups, the statistical significance observed here was not robust. In conclusion, using this strategy, we did not identify strong evidence for loss of function mutation burden difference between the mild and the severe groups.

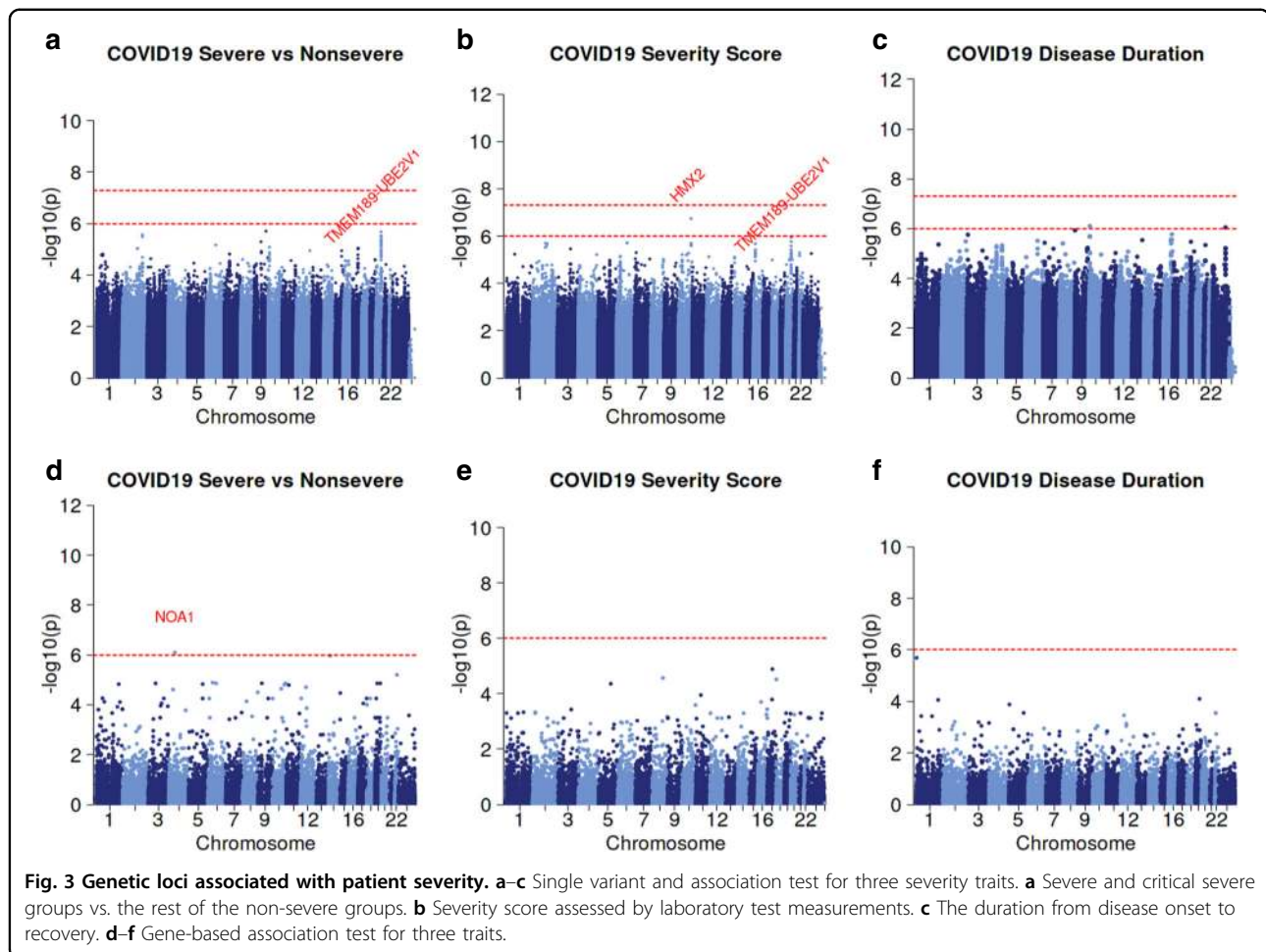
Particularly, we have inspected the missense and loss of function variants present in the SARS-CoV-2 S protein host cellular receptor gene *ACE2* and the S protein primer gene *TMPRSS2* that were known to have played a critical role in controlling the viral entry into the host cell, as well as a few other genes that were predicted to play a role in

the host–pathogen interaction network like *SLC6A19*, *ADAM17*, *RPS6*, *HNRNPA1*, *SUMO1*, *NACA*, and *BTF3*⁴⁴. The majority of the functional variants have minor allele frequency (MAF) less than 1% except for the p.Val197Met missense variant in *TMPRSS2* (Fig. 2f). Although not statistically significant, the p.Val197Met variant (rs12329760) displayed a higher allele frequency in the asymptomatic and mild group compared to the rest of the group (asymptomatic: 0.46, mild: 0.50, moderate: 0.38, severe: 0.39, critical severe: 0.26). 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 (0.14–0.29) and the Europeans (0.17–0.23) (Supplementary Fig. S12). Computational protein modeling suggested that the p.Val197Met *TMPRSS2* isoform decrease the stability of the *TMPRSS2* protein, promote the binding to S-protein and inhibit its binding with *ACE2*⁴⁴. The decreasing allele frequency in the severe patient groups supports that the p.Val197Met is related to the disease outcomes of COVID-19. The other genes did not display significant allele frequency difference among the patient groups (Supplementary Fig. S13).

Genome-wide association of common and rare variants with COVID-19 severity

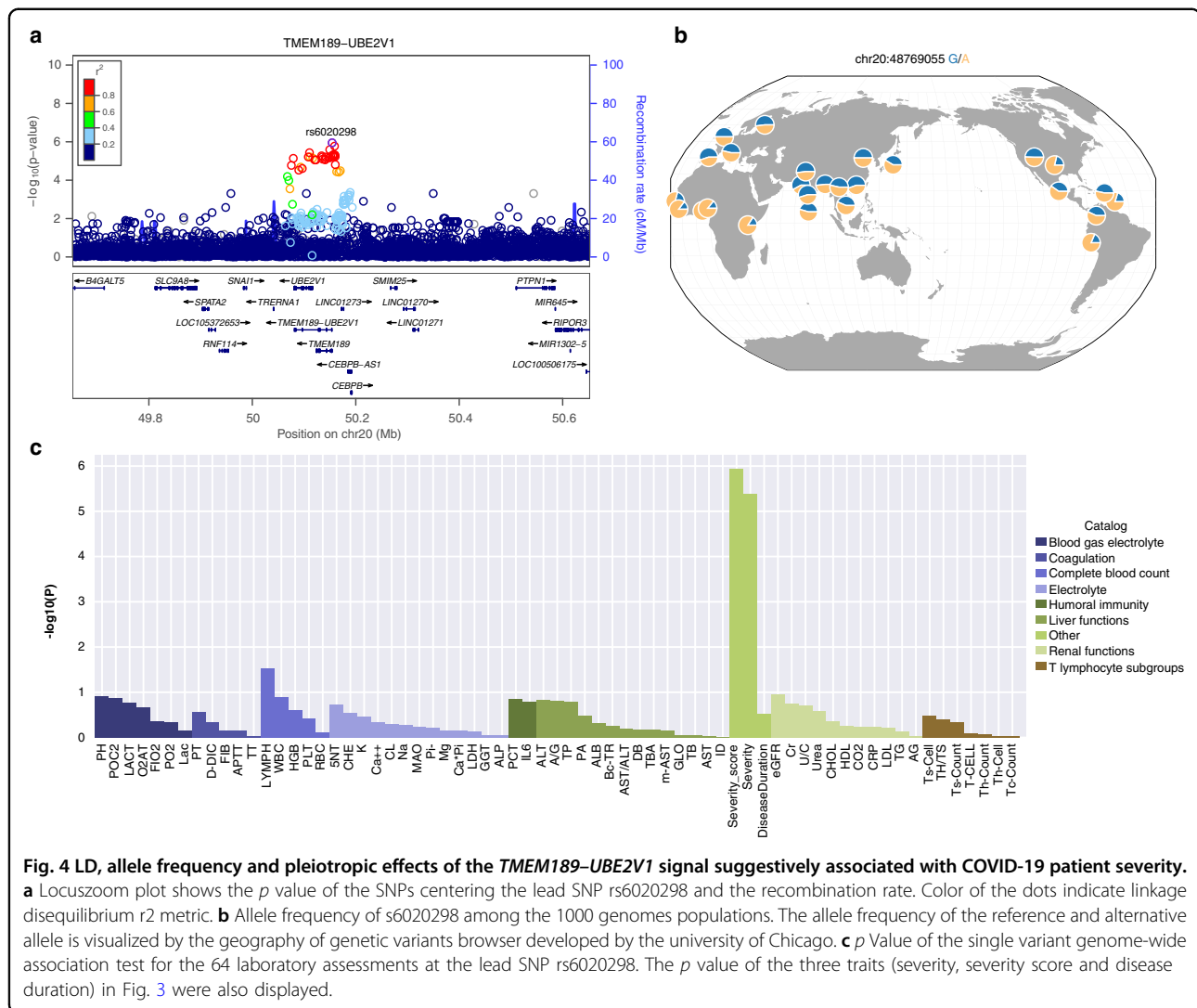
To further investigate genetic effects for the patient severity, we performed genome-wide single variant association test and sequence kernel association test (SKAT) analysis of three traits implicating patient severity. We defined the first trait as a dichotomous classification of the broadly defined “severe group” that consists of the severe and critical ill patients ($n = 70$) and the “mild group” ($n = 262$) that consists of the asymptomatic, mild, and moderate patients. We defined the second trait as a quantitative measurement of the severity level trained from the demographic features such as age, gender, and the 64 laboratory assessments ($n = 332$) (Supplementary Figs. S5 and S6). We used the disease duration from the electronic health records as the third trait which corresponds to the duration of time between the complained disease onset and the first laboratory confirmed PCR test negative outcome ($n = 233$) (Fig. 1d). Power analysis indicates that given 80% statistical power, we will be able to identify associations between genotypes and phenotypes for variants with MAF greater than 0.2 and with a relative genetic risk contribution greater than 2 given the current sample size for dichotomous trait and similarly for the quantitative trait (Supplementary Fig. S14). Principal component analysis of the patients suggests little genetic differentiation (Supplementary Figs. S15 and S16).

We tested all the QC-passed 19.6 million biallelic variants for association 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 resulting p values was very close to the null expectation ($\lambda = 0.996\text{--}1.1$, Supplementary Fig. S17) indicating that stratification was adequately controlled. The most significant single-nucleotide polymorphisms (SNP) rs6020298 is located in the intron of a read-through transcript *TMEM189-UBE2V1* in the 20q13.13 region (Fig. 3a, b). The rs6020298 (hg38 chr20:50152518, A allele frequency severe vs. non-severe: 0.59 vs. 0.45) marks a suggestive significant association signal for both the severe and mild binary trait (logistic regression $p = 4.1\text{e-}6$, OR = 1.2) and the quantitative measurement of the severity score (linear regression $p = 1.1\text{e-}6$, beta = 0.35). SNPs in linkage disequilibrium with rs6020298 ($r^2 > 0.8$) also affect the gene *UBE2V1* and *TMEM189* (Fig. 4a). The *UBE2V1* gene encodes the ubiquitin-conjugating enzyme E2 variant 1. Both the *UBE2V1* and *TMEM189-UBE2V1* have been involved in the interleukin-1 (IL-1) signaling pathway⁴¹ and suggested to work together with *TRIM5* to promote innate immune signaling⁴⁵. IL-1 is elevated in COVID-19 patients especially the severe and critical patients who suffer from the cytokine storm and severe inflammation⁴⁶.

Clinical trial using IL-1 blockade on critical patients results in an improvement in respiratory function in 72% of the patients⁴⁷. The lead SNP rs6020284 has a MAF close to 0.5 among the worldwide populations except for the African population (AF = 0.13) (Fig. 4b). It is also an eQTL for *LINC01273*, *TMEM189* among several tissues including the lung where the risk A allele increases the *TMEM189* and *LINC01273* expression in several tissues (Supplementary Fig. S18). This may indicate that an inborn evaluated *TMEM189* expression in the patients may promote IL-1 signaling and predisposes the patients toward a poorer outcome against the COVID-19 infection. However, given the limited sample size in this study and that the intermediate pathways between *TMEM189* and IL-1 production is still unclear, more replication and functional validation efforts should be made to re-evaluate this association signal. Notably, the *TMEM189-UBE2V1* locus has been associated with monocyte percentage of leukocytes and granulocyte percentage of myeloid white cells⁴⁸. We did not observe nominal association ($p < 0.05$) at the lead SNP rs6020298 with all the 64 laboratory assessments among the patients (Fig. 4c). Therefore, the observed signal



is not supposed to be confounded by individual variability on blood cell types. There is no strong genetic association with the disease durations (Fig. 3c). We evaluated the rs6020298 signal in the European population using the host genetic initiative data. There is no significant association for this locus among the European population (Supplementary Fig. S19). We also performed a GWAS on the 81,193 copy number variations identified from the 332 individuals. We did not identify copy number variations that were significantly associated with the disease severity (Supplementary Fig. S20).

As for the time when this study was conducted, two loci-3p21.31 and 9q34.2 were reported to have been strongly associated with the risk of progressing in critical severe COVID-19 in the Spanish and Italian population⁴⁹. We compared the association signal and the allele frequency of these two loci reported by the Spanish and Italian study with that in the Chinese population (Table 1). For the

3p21.31 locus, although the risk allele of the lead SNP (rs11385942) is a common variant present in the European, African, and the South Asian populations, it is almost absent in the Chinese and the general East Asian population (allele frequency ~ 0) (Table 1). Therefore, we didn't identify any variants that were significantly associated with the patient severity in this region (Supplementary Fig. S21). On the other hand, the allele frequency of the lead SNP (rs657152) at 9q34.2 is similar around the globe ($\sim 0.42-0.63$). We found an increasing of risk for critically severe COVID-19 patient with the risk allele that determines the blood type A in the 9q34.2 loci compared to the other blood types for critically severe patients although it is not statistically significant. The comparison analysis here suggests that prediction of genetic risk should consider the genetic diversity from different populations.

We further performed optimal SKAT gene-based association test on the functional variants including a total of

Table 1 Comparison of the allele frequency and genome-wide association signals for two associated loci in European population.

Compared information rs11385942 3p21.31 rs657152 9q34.2, ABO			
<i>Allele frequency</i>			
CHROM	chr3	chr9	
POS (hg38)	45,834,967	133,263,862	
Risk allele	GA	C	
Other allele	G	A	
All patients (N = 284)	0	0.453	
Asymptomatic (N = 12)	0	0.413	
Moderate (N = 194)	0	0.44	
Severe (N = 52)	0	0.469	
Critical (N = 17)	0	0.618	
ChinaMAP	0.00396	0.424	
1000G_EAS	0.005	0.628	
1000G_EUR	0.0805	0.601	
1000G_SAS	0.296	0.596	
1000G_AFR	0.053	0.539	
gnomAD_EAS	0.00061	0.633	
<i>Association</i>			
Italian	OR (95% CI)	1.53–2.48	1.22–1.59
	p value	7.02E–08	5.31E–05
Spain	OR (95% CI)	1.76–4.42	1.17–1.60
	p value	1.17E–05	2.81E–03
Chinese	OR (95% CI)	NA	0.38–1.42
	p Value	NA	0.693

99,166 missense and loss of function variants that were predicted to have high or moderate impacts by variant effect predictor among the patients. The *NOAI* gene tend to higher mutation burden in the severe group ($P = 8.1e-07$) (Fig. 3d). This gene encodes the GTPase that functions in the mitochondrion and has been associated with platelet count and leukocyte count⁴¹. We did not identify other genes that are genome-wide significantly associated with the severity score or the disease duration (Fig. 3e, f).

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

Manifestation of numerous infectious diseases are closely related to the genetic variants across the major histocompatibility complex genes, i.e., the human leukocyte antigen (HLA) genes, which play an essential role in presenting the antigen determinant epitopes from the pathogens to the T cell or B cell to activate the host

immune response^{50,51}. In the 2003 SARS outbreak, caused by the SARS coronavirus (SARS-CoV) related to SARS-CoV-2, the HLA-B*46:01 was reported to be associated with infection severity in East Asian patients²⁴. Herein, we investigated the genetic effect from HLA genes on the COVID-19 patient severity. We re-aligned all the reads mapped to the eight HLA haplotypes in the human reference genome (GRCh38) and all the unaligned reads and typed the three class I HLA genes (A–C) and four class II HLA genes (DPB1, DQA1, DQB1, and DRB1) using the xHLA⁵² and the SOAP-HLA approach⁵³. 4-digit haplotyping resolution was achieved for 99% of the patients for all the genes except for DQA1 where three patients were only typed to the 2-digit resolution. We observed zero mendelian error rate for the typing results using the family members involved in the study. We investigated whether some HLA alleles may significantly differ between the broadly defined severe (severe and critical, $n = 69$) and mild (asymptomatic, mild and moderate, $n = 215$) groups of unrelated patients using a logistic regression 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-A, 51 HLA-B, 28 HLA-C, 20 DPB1, 21 DQA1, 16 DQB1 and 32 DRB1 alleles were displayed in Fig. 5 and Supplementary Table S6. Among the class I HLA genes, C*14:02 (severe 8.7% vs. mild 4.6%, OR = 4.7, $P = 3e-3$), B*51:01 (severe 10.1% vs. mild 5.8%, OR = 3.3, $p = 7e-3$), A*11:01 (severe 29.7% vs. 26.2%, OR = 2.3, $p = 8.5e-3$) were the top three most significant alleles between the two groups that predispose the patients entering the severe stage (Table 2). The HLA-A*11:01, B*51:01, and C*14:02 is in strong linkage equilibrium with each other and thus represents one haplotype. This haplotype has an average allele frequency 2.4%–3.6% among the Chinese populations according to the HLA Allele Frequency Net Database⁵⁴. In our study, we find that this haplotype is more prevalent in the severe patients compared with the mild patients.

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 covariates²⁴, our analysis does not 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 severe patients (10.1%) than among the mild patients (12.8%). Class II HLA genes is less significantly associated with the disease severity compared to the Class I genes (Table 2). DRB1*14:04 (severe 2% vs. mild 0.5%, $p = 0.01$), 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 DRB1*12:01

Table 2 Nominal association of HLA allele and severity by logistic regression.

	Severe	Non-severe	OR	SE	P
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

Severe group indicates severe and critical patients.

Non-severe group includes asymptomatic, mild, and moderate patients.

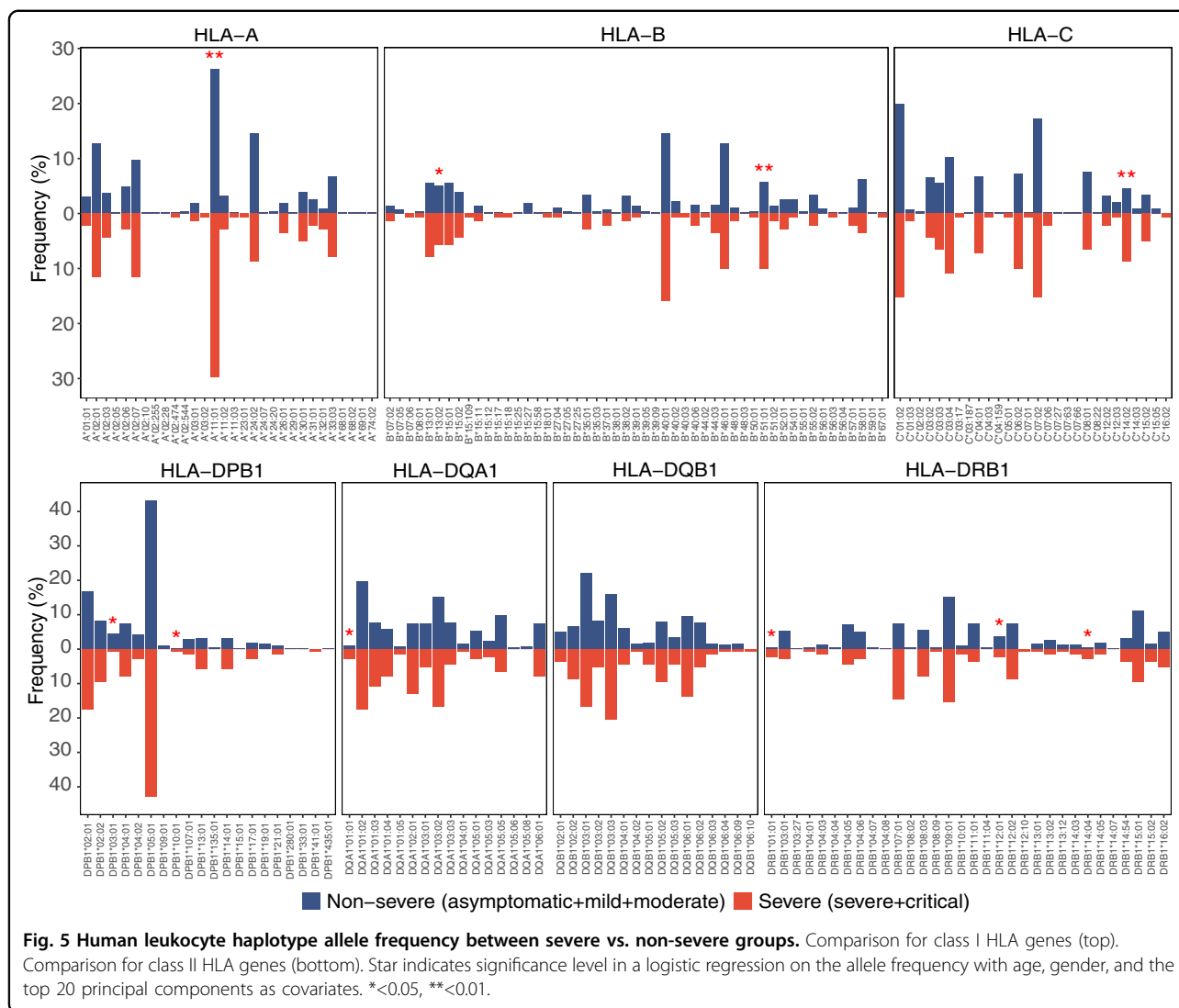
(severe 2.2% vs. mild 3.7%) might display a protective effect.

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 the genetic contribution to patient severity. Mapping genes related to infection susceptibility is more difficult. An ideal design commands a comparison between people who are exposed or not exposed to the pathogen. This is hard to meet because early detection and isolation of infected patients are the primary containment strategies against an outbreak⁵⁶. Therefore, we choose another approach to investigate genetic susceptibility by comparing the 284 unrelated hospitalized patients (the Case) with two general populations including 301 Chinese individuals in 1000 genome project³⁵ (the Control I) and 665 individuals recruited from the Chinese Reference Panel program (CNPR, paper in preparation, the Control II). Control I and Control II differ in terms of the similarity of the adopted sequencing protocol compared to the Case. All the technical components are almost the same between the Case and Control 2 except for sequencing depth (case 46× vs. control 2 30×). On the other hand, various factors are different between the Case and Control 1, including types of sample (case fresh blood vs. control 1 cell line), sequencing technology (case MGI's nanoball sequencing vs. control 1 Illumina sequencing), sequencing read cycles (case 100 bp pair-end vs. control 1 150 bp pair-end) and the sequencing depth (case average 46× vs. average 7×). Study like this can reveal genetic difference between the infected population and the general population if any and if not, instruct on what cautions should be taken when comparing the disease cohorts vs. the general in the whole genome sequencing context. The comparison of the infected patients and the general

population was also conducted by the COVID-19 host genetic initiative.

We analyzed the data carefully by jointly genotype the samples from their individual gvcf files using the GATK best practices³⁸ instead of simply merging the population vcf files of the case and the control. Principle component analysis indicates that population structure is the dominant confounding factor and sequencing induced batch effects were difficult to identify in the PCs (Supplementary Figs. S22 and S23). Similarly, we conducted both single variant and gene-based association tests for the two case-control data sets using the top 20 PCs, gender and age (age was not available for 1 KGP samples and was used for the CNPR alone) as covariates Fig. 6. Surprisingly, in the single association test for the high and moderate impact variants, many variants in the HLA region displayed significant associations between the COVID-19 patients and the 1 KGP Chinese (Fig. 6a) even though the inflation was seemingly adequately controlled (Supplementary Fig. S24). In the gene-based association test, we observed significantly different mutation burdens in the immunoglobulin loci (Fig. 6b). However, this was not replicated when we compared the COVID-19 patients with the 665 CNRP individuals (Fig. 6c, d). Therefore, we inferred that the association signals between the 1 KGP and the COVID-19 patients were probably due to sequencing batch effects. As the fresh blood of an infected individual contains numerous somatic mutated B-cells, patients tend to accumulate more mutations in the immunoglobulin genes⁵⁷. As many studies try to directly compare the allele frequency between the general population and the COVID-19 patients^{29,31}, our discoveries remind us of the necessity for re-evaluation of the significant hits given distinct experimental protocol for case and control.

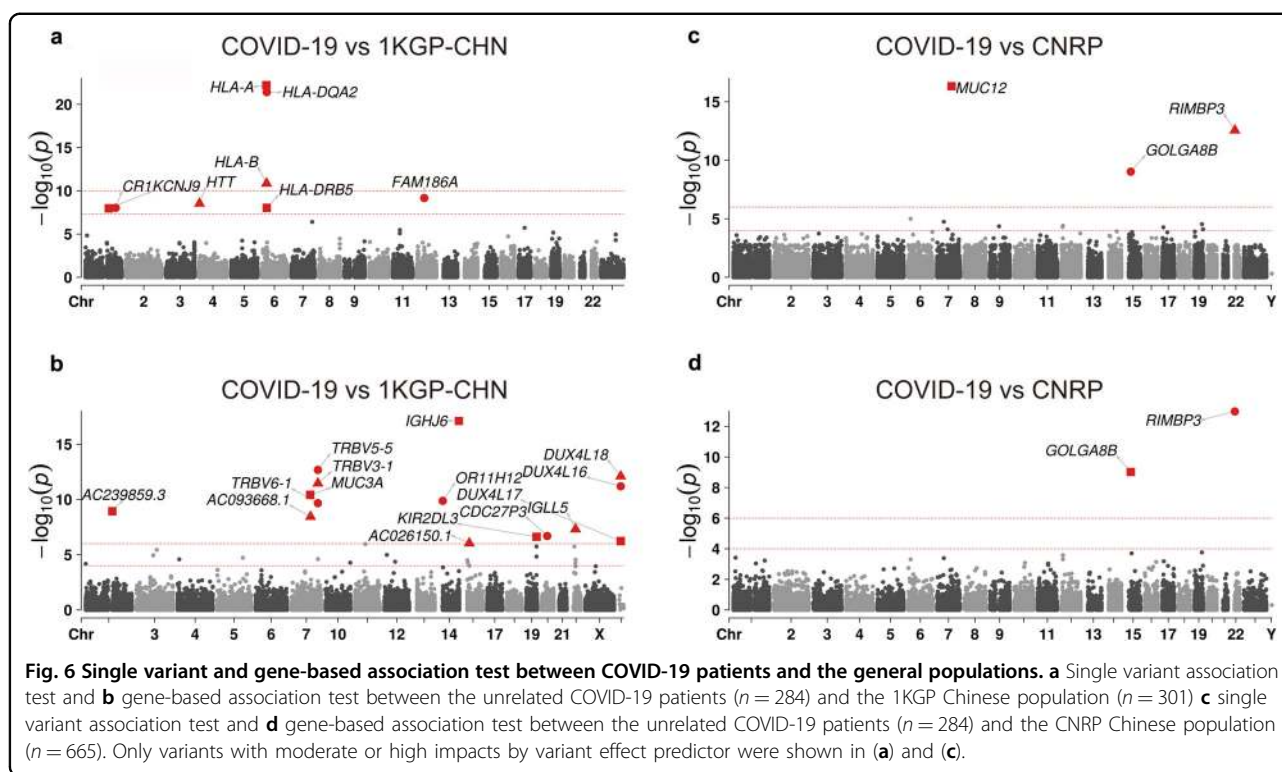


In the single variant association test between the COVID-19 patients and the CNPR who were sequenced using the same experimental protocol and were laboratory PCR tested negative, we identified genome-wide significant associated signals tagged by a novel missense variant (Patient T allele frequency = 0.34, CNPR T_AF = 0.14, OR = 18, $p = 4.7e-17$) in *MUC2*; a missense variant rs200584390 (Patient G allele frequency = 0.31, CNPR G_AF = 0.09, 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* (Fig. 6c). Gene-based association test also indicates that *RIMBP3* and *GOLGA8B* were different between the patients and the CNPR (Fig. 6d). Those discoveries require further replication and interpretation when more sequencing data for patients and for general populations become available worldwide³².

Discussion

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

We revealed that the disease progression after the SARS-CoV-2 infection can be determined by both the



monogenic and complex genetic basis. In the investigation of potential monogenic effects, using a pedigree gene-mapping strategy, we identified a recurrent loss of function mutation in gene *GOLGA3* among the critically ill patients and a recurrent loss of function 1-bp insertion in gene *DPP7* among the asymptomatic patients. Both genes were related to the host immune response to the viral infection. We did not identify further genes that has a large monogenic effect using the population gene-based association test strategy.

For a list of candidate genes proposed by previous study on the molecular virology of SARS-CoV-2, we identified that the missense variant rs12329760 in *TMPRSS2* was less frequent among the critical patients compared to the rest of the patients and the general population. This variant results in an alteration of the valine to the methionine at the 197th amino acids (p.Val197Met) that has been predicted to decrease the *TMPRSS2* protein stability and ACE2 binding⁴⁴. On the other hand, our study using Chinese samples did not support the assumption²⁹ that host genetic factors in the essential SARS-CoV 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 *TMEM189–BE2V1* and *TEMEM189–UBE2V1* that are known to function in the IL-1 signaling pathway^{41,45}. The lead 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 critical patients. While COVID-19 severe patients demonstrate elevated IL-1 compared to the mild patients and the general population⁴⁶, our study suggests potential correlation between genetic variability in this gene and the disease severity.

Notably, the HLA-A*11:01, B*51:01, and C*14:02 alleles were significantly more prevalent in the severe and critical severe patients compared to the mild and the moderate patients after careful control of population structure and demographic characters such as age and gender. The three alleles were in linkage disequilibrium with each other and has been previously reported to have a 2%–3% population allele frequency in Dai and Jinpo minorities in China⁵⁴ and the B*51:01 has been previously linked to the Behcet’s disease⁶¹, a kind of rheumatic disease. We were not able to access the role of HLA-B*46:01, although it has been predicted as the worst presenting HLA alleles to the SARS-CoV-2 proteome⁵⁵ and linked to the SARS 2003 outbreak²⁴.

Surprisingly, GWAS using the COVID-19 patients as the case and the 1000 genome Chinese population as the control suggested an enrichment of significantly associated signals in the HLA region and mutation burden in the immunoglobulin genes. Nonetheless, this was not replicated when we compared the patients to another independent Chinese population. A lot of efforts in the

genetic field have been made and there may be more in the future to investigate genetic susceptibility of the SARS-COV-2 infection by directly comparing two or more general populations with the COVID-19 patients^{31,32}. Therefore, cautions should be taken to properly control the batch effects. Replication is essential and perhaps a joint or meta-analysis strategy can rule out the real from the false statistical signals.

Some limitations of the study should be noted. Power analysis indicates that sample size of around 300 is barely sufficient to identify genome-wide significant genetic variants with MAF greater than 0.2 and odds ratio greater than 1.8 given type I error rate 0.05. We do not have power to detect causal variants beyond this risk and allele frequency scenario. In addition, although the study of hospitalized patients in a designated hospital includes all severe patients, the design has a limited presentation of the asymptomatic patients (7.5%) which ratio has been estimated to be 30.8% (95% confidence interval: 7.7%–53.8%)⁶². Given that RT-PCR test and the seroprevalence immunoglobulin M and G antibody tests targeting the SARS-CoV-2 has been 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 and around the world to understand the host genetic background underlying the varying clinical outcome of the patients, this work represents the first genetic study on the Chinese hospitalized patients where high quality sequencing data were generated and systematic analysis on the genomic and clinical data were conducted. Our results highlight several genetic factors involved in the immune responses including genes involved in the viral entry in the host cells, genes 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 samples, the genetic assay approach, the bioinformatics and the statistical genetic analysis for COVID-19 as well as other infection and complex disease. The publicly available summary statistics will encourage international collaborative efforts to understand the host–pathogen interaction and to contain the COVID-19 outbreak.

Material and methods

Patient recruitment and definition of phenotypes

A total of 332 patients were recruited from Jan 11th 2020 to Apr 2020 in Shenzhen Third People's Hospital, the only referral hospital in Shenzhen City, China³³. All were confirmed with SARS-COV-2 infection using real-time RT-PCR assay of nasal and pharyngeal swab specimens. The demographic, epidemiological, clinical, and

laboratory assessments were extracted from the electronic medical records of the patients. This study was approved by the ethics commissions of the Shenzhen Third People's Hospital Ethics Committee with a waiver of informed consent. According to the 5th edition of the national treatment guideline of COVID19 in China and the Chinese CDC criteria⁶, the patients were diagnosed as asymptomatic, mild, moderate, severe, and critically ill according to the most severe stage they experienced during the disease course. The asymptomatic, mild and the moderate groups of patients do not experience pneumonia. When meeting any one of the following criteria: (1) RR > 30, (2) oxygen level < 93%, (3) PaO₂/FiO₂ < 300 mmHg, and (4) disease progression greater than 50% area in CT scan, a patient is categorized as severe patients. Patients experienced one of the following: (1) respiratory failure and requires mechanical ventilation, (2) shock, and (3) complicated by failure of other organs and requires intensive care monitoring were classified as critically ill.

Assignment of severity score to each patient

A machine learning XGBoost-based model was developed to predict ordinal 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 did not have any recordings. The remaining 52 laboratory test items with missing values were further imputed by missForest algorithm⁶⁴. The missForest is a nonparametric method to impute missing values using random forest model in an iterative fashion. Then the originally ordered severity levels of asymptomatic, mild, moderate, severe, and critical were assigned integer values of 1–5, respectively. The numeric representations retained the ordinal levels of severity. We applied the reduction framework mentioned in Li et al.⁶⁵, where the ordinal regression was reduced to binary classification. The reduction framework of extended binary classification was then integrated within XGBoost model. Moreover, we selected the most predictive laboratory test items using SHAP (SHapley Additive exPlanations) algorithm⁶⁶. The SHAP is a game theoretic approach to explain the output of a given machine learning model using Shapley values from game theory and their related extensions. We finally trained the XGBoost-based ordinal regression model using the selected laboratory test items. As a result, the prediction outcome produced by the final model was typically a real number reflecting severity level that was used in the downstream analysis. We used 100 base estimators for missForest, maximum iteration of 10, and the criterion was mean squared error. For the XGBoost-based ordinal regression model, we used 500 base estimators and learn rate of 0.5. In general, the

hyper-parameters of models in this study were chosen by combining grid search of fivefold cross validation and manual tuning.

DNA extraction, library construction, and deep whole-genome sequencing

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

Genome alignment and variant detection

We used Sentieon Genomics software (version: sentieon-genomics-201911) to perform genome alignment and variant detection⁶⁷. Analysis pipeline were built according to the recommendation in the Broad institute best practices described in <https://gatk.broadinstitute.org/hc/en-us/sections/360007226651-Best-Practices-Workflows>. Sequencing reads were mapped to hg38 reference genome using BWA algorithm. For each sample, after remove duplicates, Indel realignment and base quality score recalibration, SNP and short Indel variants were detect using the Sentieon Haplotyper algorithm with option `--emit_mode gvcf` to generate an individual GVCF file. Then the GVCF files for all samples were subjected to Sentieon GVCFTyper algorithm to perform joint variant calling. Copy number variations were identified using manta v1.6⁶⁸, an assembly-based caller using the default parameters following the protocol in a recent HGDP project⁶⁹.

Variant quality score recalibration and filtration

Variant quality score recalibration were perform using Genome Analysis Toolkit (GATK version 4.1.2). Known variant files were downloaded from the GATK bundle. For indel recalibration, we used Mills_and_1000G_gold_standard indels as the positive training and true set. For SNP recalibration, we used hapmap_3.3, 1000G_omni2.5, and 1000G_phase1.snps as positive training sets, hapmap_3.3 as true set, and dbSNP_v146 as the known set. The metrics DP, QD, MQRankSum, ReadPosRankSum, FS, SOR were used in the recalibration process. The truth-sensitivity-filter-level were set to 99.0 for both the SNPs and the Indels. Finally, variants with quality score ≥ 100 were selected for further analysis.

Variant effect prediction

Annotation of the genes mentioned in the paper and the annotation of the existence of the variants in database

such as dbSNP, GnomeAD, 1 KGP was carried out using Variant Effect Predictor⁷⁰ using the default parameters.

Familial relationship and principal component analysis

PLINK (v1.9)⁷¹ and KING (v2.1.5)⁷² was applied to detect the kinship relatedness between each pair of the individuals. 48 patients from 16 families were detected as related to each other. For several allele frequency-based approach, we exclude the related patients and thus the sample size was restricted to 284. PCA was performed using a subset of autosomal bi-allelic SNPs on the unrelated patients using PLINK (v1.9). The PC-AiR module (PCA in related samples) in the Genesis R package was used to conduct PCA analysis for the 332 patients including the related family members. Several restrictions were employed to select the final 614,963 SNPs for PCA analysis, including $\text{MAF} \geq 1\%$ (common and low-frequency variants), genotyping rate $\geq 90\%$, Hardy–Weinberg–Equilibrium $P > 0.000001$, and removing one SNP from each pair with $r^2 \geq 0.5$ (in windows of 50 SNPs with steps of 5 SNPs).

Genotype–phenotype association tests

We have applied both the rvtest⁷³ and the SAIGE⁷⁴ approaches to carry out logistic regression, linear regression, burden test, the SKAT and the optimal SKAT-O algorithm for the genotype–phenotype association tests using the default parameters. For all the association tests, we have used the gender, the age and the top 20 principal components from the principal component analysis as the covariates. Exception is for the GWAS between the 1 KGP and the COVID-19 patients as age is not available for the 1 KGP data set. Independent loci were defined as significant variants clustered in a 1 Mbp window. The lead SNP was defined as the SNP in the 1 Mbp window that has most significant, i.e., smallest p value. The genomic inflation factor, GC lambda, attenuation ratio, LD score regression intercept and the SNP heritability were estimated using the LD score regression approach⁷⁵. The qqman R package was applied to generate the manhattan and qqplot. We defined genome-wide significance for single variant association test as $5e-8$, suggestive significance as $1e-5$ and for gene-based association test as $1e-6$.

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 algorithm²² 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.

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 discussion of the results and advice. We would also like to acknowledge Ms. Juehan Liu's contribution on literature search and data visualization on this study.

Author details

¹The Third People's Hospital of Shenzhen, National Clinical Research Center for Infectious Disease, The Second Affiliated Hospital of Southern University of Science and Technology, Shenzhen, Guangdong 518112, China. ²BGI-Shenzhen, Shenzhen, Guangdong 518083, China. ³School of Medicine, South China University of Technology, Guangzhou, Guangdong 510006, China. ⁴BGI Education Center, University of Chinese Academy of Sciences, Shenzhen, Guangdong 518083, China. ⁵Guangdong Provincial Key Laboratory of Genome Read and Write, BGI-Shenzhen, Shenzhen, Guangdong 518120, China. ⁶James D. Watson Institute of Genome Science, Hangzhou, Zhejiang 310008, China

Author contributions

Conceptualization: L.L., X.J., Q. H., S. Liu, and S.H.; Methodology: S.H., S. Liu, Y.Z., and X.T.; Formal analysis: S.H., Y.Z., X.Q., Zhi. L., P.L., Y.H., R.L., X.T., Y.B., and S. Liu; Resources: F.W., R.G., C.L., W.X., Zhi.L., Q.T. R.C. X.L., X.Z., and G.D.; Data curation: S. Liu, F.W., R.G., and C.L.; Writing—original draft: S. Liu; Writing—review and editing: S. Liu and all; Supervision: X.X., J.W., and H.Y.; Project administration: X.J., S.H., S. Liu, and F.C.; Funding acquisition: F.W., X.J., and S. Liu.

Data availability

The data that support the findings of this study, including the allele frequency for the five groups of patients at all the 19.6 million biallelic genetic variants and the genome-wide association test summary statistics have been deposited in CNSA (China National GeneBank Sequence Archive, <https://db.cngb.org/cnsa/>) in Shenzhen, China with accession number CNP0001107. The release of the data was approved by the Ministry of Science and Technology of China (Project ID: 2020BAT0262).

Conflict of interest

The authors declare that they have no conflict of interest.

Publisher's note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Supplementary Information accompanies the paper at (<https://doi.org/10.1038/s41421-020-00231-4>).

Received: 12 June 2020 Accepted: 3 October 2020

Published online: 10 November 2020

References

- Johnson, N. P. A. S. & Mueller, J. Updating the accounts: global mortality of the 1918–1920 'Spanish' influenza pandemic. *Bull. Hist. Med.* **76**, 105–115 (2002).
- Zhu, N. et al. A novel coronavirus from patients with pneumonia in China, 2019. *N. Engl. J. Med.* **382**, 727–733 (2020).
- John Hopkins University and Medicine. COVID-19 Map—Johns Hopkins Coronavirus Resource Center. (John Hopkins Coronavirus Resource Center, 2020).
- Jiang, S., Du, L. & Shi, Z. An emerging coronavirus causing pneumonia outbreak in Wuhan, China: calling for developing therapeutic and prophylactic strategies. *Emerg Microbes Infect.* **9**, 275–277 (2020).
- Shi, Z. & Hu, Z. A review of studies on animal reservoirs of the SARS coronavirus. *Virus Res.* **133**, 74–87 (2008).
- Wu, Z. & McGoogan, J. M. Characteristics of and Important lessons from the coronavirus disease 2019 (COVID-19) outbreak in China: summary of a report of 72 314 cases from the Chinese center for disease control and prevention. *JAMA* **323**, 1239–1242 (2020).
- Kenney, A. D. et al. Human genetic determinants of viral diseases. *Annu. Rev. Genet.* **51**, 241–263 (2017).
- Guan, W. et al. Clinical characteristics of coronavirus disease 2019 in China. *N. Engl. J. Med.* **382**, 1708–1720 (2020).
- Fu, L. et al. Clinical characteristics of coronavirus disease 2019 (COVID-19) in China: a systematic review and meta-analysis. *J. Infect.* **80**, 656–665 (2020).
- Liu, Z., Bing, X. & Zhi, X. Z. The epidemiological characteristics of an outbreak of 2019 novel coronavirus diseases (COVID-19) in China. *Zhonghua Liu Xing Bing Xue Za Zhi* **41**, 145–151 (2020).
- Qin, C. et al. Dysregulation of immune response in patients with COVID-19 in Wuhan, China. *Clin. Infect. Dis.* **71**, 762–768 (2020).
- Yang, X. et al. Clinical course and outcomes of critically ill patients with SARS-CoV-2 pneumonia in Wuhan, China: a single-centered, retrospective, observational study. *Lancet Respir. Med.* **8**, 475–481 (2020).
- Nishiura, H. Backcalculating the Incidence of Infection with COVID-19 on the Diamond Princess. *J. Clin. Med.* **9**, 657 (2020).
- Hu, Z. et al. Clinical characteristics of 24 asymptomatic infections with COVID-19 screened among close contacts in Nanjing, China. *Sci. China Life Sci.* **63**, 706–711 (2020).
- Fellay, J. et al. Common genetic variation and the control of HIV-1 in humans. *PLoS Genet.* **5**, e1000791 (2009).
- Fellay, J. et al. A whole-genome association study of major determinants for host control of HIV-1. *Science* **317**, 944–947 (2007).
- Pereyra, F. et al. The major genetic determinants of HIV-1 control affect HLA class I peptide presentation. *Science* **330**, 1551–1557 (2010).
- Kamatani, Y. et al. A genome-wide association study identifies variants in the HLA-DP locus associated with chronic hepatitis B in Asians. *Nat. Genet.* **41**, 591–595 (2009).
- Ge, D. et al. Genetic variation in IL28B predicts hepatitis C treatment-induced viral clearance. *Nature* **461**, 399–401 (2009).
- Schulert, G. S. et al. Whole-exome sequencing reveals mutations in genes linked to hemophagocytic lymphohistiocytosis and macrophage activation syndrome in fatal cases of H1N1 influenza. *J. Infect. Dis.* **213**, 1180–1188 (2016).
- Wang, Z. et al. Early hypercytokinemia is associated with interferon-induced transmembrane protein-3 dysfunction and predictive of fatal H7N9 infection. *Proc. Natl. Acad. Sci. USA* **111**, 769–774 (2014).
- Everitt, A. R. et al. IFITM3 restricts the morbidity and mortality associated with influenza. *Nature* **484**, 519–523 (2012).
- Yang, X. et al. Interferon-inducible transmembrane protein 3 genetic variant rs12252 and influenza susceptibility and severity: a meta-analysis. *PLoS One* **10**, e0124985 (2015).
- Lin, M. et al. Association of HLA class I with severe acute respiratory syndrome coronavirus infection. *BMC Med. Genet.* <https://doi.org/10.1186/1471-2350-4-9> (2003).
- Ching, J. C. et al. Significance of the myxovirus resistance A (MxA) Gene –123C>A single-nucleotide polymorphism in suppressed interferon β induction of severe acute respiratory syndrome coronavirus infection. *J. Infect. Dis.* **201**, 1899–1908 (2010).
- Kachuri, L. et al. The landscape of host genetic factors involved in infection to common viruses and SARS-CoV-2. *medRxiv* <https://doi.org/10.1101/2020.05.01.20088054> (2020).
- Williams, F. M. et al. Self-reported symptoms of covid-19 including symptoms most predictive of SARS-CoV-2 infection, are heritable. *medRxiv* <https://doi.org/10.1101/2020.04.22.20072124> (2020).
- Zhao, J. et al. Relationship between the ABO Blood Group and the COVID-19 Susceptibility. *medRxiv*. <https://doi.org/10.1101/2020.03.11.20031096> (2020).
- Cao, Y. et al. Comparative genetic analysis of the novel coronavirus (2019-nCoV/SARS-CoV-2) receptor ACE2 in different populations. *Cell Discov.* **6**, 4–7 (2020).
- Bhattacharyya, C. et al. Global Spread of SARS-CoV-2 Subtype with Spike Protein Mutation D614G is Shaped by Human Genomic Variations that Regulate Expression of TMPRSS2 and MX1 Genes. *bioRxiv* <https://doi.org/10.1101/2020.05.04.075911> (2020).

31. Renieri, A. et al. ACE2 variants underlie interindividual variability and susceptibility to COVID-19 in Italian population. *medRxiv* <https://doi.org/10.1101/2020.04.03.20047977> (2020).
32. The COVID-19 Host Genetics Initiative, a global initiative to elucidate the role of host genetic factors in susceptibility and severity of the SARS-CoV-2 virus pandemic. *Eur. J. Hum. Genet.* **28**, 715–718 (2020).
33. Cai, Q. et al. COVID-19 in a designated infectious diseases hospital outside Hubei Province, China. *Allergy Eur. J. Allergy Clin. Immunol.* **75**, 1742–1752 (2020).
34. Darbehshiti, F. & Rezaei, N. Genetic predisposition models to COVID-19 infection. *Med. Hypotheses* <https://doi.org/10.1016/j.mehy.2020.109818> (2020).
35. Auton, A. et al. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
36. Lundberg, S. M. et al. From local explanations to global understanding with explainable AI for trees. *Nat. Mach. Intell.* **2**, 56–67 (2020).
37. Jiang, X. et al. Towards an artificial intelligence framework for data-driven prediction of coronavirus clinical severity. *Comput. Mater. Contin.* **63**, 537–551 (2020).
38. DePristo, M. A. et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
39. Karczewski, K. J. et al. Variation across 141,456 human exomes and genomes reveals the spectrum of loss-of-function intolerance across human protein-coding genes. *bioRxiv* <https://doi.org/10.1101/531210> (2019).
40. Collins, R. L. et al. An open resource of structural variation for medical and population genetics. *bioRxiv* <https://doi.org/10.1101/578674> (2019).
41. Stelzer, G. et al. The GeneCards suite: From gene data mining to disease genome sequence analyses. *Curr. Protoc. Bioinform.* <https://doi.org/10.1002/cpbi5> (2016).
42. Wang, N. et al. Structure of MERS-CoV spike receptor-binding domain complexed with human receptor DPP4. *Cell Res.* **23**, 986–993 (2013).
43. Dai, W. et al. Whole-exome sequencing identifies MST1R as a genetic susceptibility gene in nasopharyngeal carcinoma. *Proc. Natl. Acad. Sci. USA* **113**, 3317–3322 (2016).
44. Sharma, S. et al. ACE2 Homo-dimerization, Human genomic variants and interaction of host proteins explain high population specific differences in outcomes of COVID19. *bioRxiv* <https://doi.org/10.1101/2020.04.24.050534> (2020).
45. Pertel, T. et al. TRIM5 is an innate immune sensor for the retrovirus capsid lattice. *Nature* **472**, 361–365 (2011).
46. Shi, Y. et al. COVID-19 infection: the perspectives on immune responses. *Cell Death Diff.* **27**, 1451–1454 (2020).
47. Cavalli, G. et al. Interleukin-1 blockade with high-dose anakinra in patients with COVID-19, acute respiratory distress syndrome, and hyperinflammation: a retrospective cohort study. *Lancet Rheumatol.* **2**, e325–e331 (2020).
48. Buniello, A. et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–D1012 (2019).
49. Ellinghaus, D. et al. Genome-wide association study of severe covid-19 with respiratory failure. *N. Engl. J. Med.* **383**, 1522–1534 (2020).
50. Hammer, C. et al. Amino acid variation in HLA class II proteins is a major determinant of humoral response to common viruses. *Am. J. Hum. Genet.* **97**, 738–743 (2015).
51. Tian, C. et al. Genome-wide association and HLA region fine-mapping studies identify susceptibility loci for multiple common infections. *Nat. Commun.* <https://doi.org/10.1038/s41467-017-00257-5> (2017).
52. Xie, C. et al. Fast and accurate HLA typing from short-read next-generation sequence data with xHLA. *Proc. Natl. Acad. Sci. USA* **114**, 8059–8064 (2017).
53. Cao, H. et al. An integrated tool to study MHC region: accurate SNV detection and HLA genes typing in human MHC Region using targeted high-throughput sequencing. *PLoS ONE* <https://doi.org/10.1371/journal.pone.0069388> (2013).
54. Gonzalez-Galarza, F. F. et al. Allele frequency net database. *Methods Mol. Biol.* https://doi.org/10.1007/978-1-4939-8546-3_4 (2018).
55. Nguyen, A. et al. Human leukocyte antigen susceptibility map for SARS-CoV-2. *J. Virol.* <https://doi.org/10.1128/jvi.00510-20> (2020).
56. Lai, S. et al. Effect of non-pharmaceutical interventions to contain COVID-19 in China. *Nature* **585**, 410–413 (2020).
57. Ju, B. et al. Potent human neutralizing antibodies elicited by SARS-CoV-2 infection. *bioRxiv* <https://doi.org/10.1101/2020.03.21.990770> (2020).
58. Rashkin, S., Jun, G., Chen, S. & Abecasis, G. R. Optimal sequencing strategies for identifying disease-associated singletons. *PLoS Genet.* <https://doi.org/10.1371/journal.pgen.1006811> (2017).
59. Wainschein, P. et al. Recovery of trait heritability from whole genome sequence data. *Yearb. Paediatr. Endocrinol.* <https://doi.org/10.1530/ey.16.14.15> (2019).
60. Carithers, L. J. & Moore, H. M. The Genotype-Tissue Expression (GTEx) Project. *Biopreserv. Biobank* **13**, 307–308 (2015).
61. Giza, M., Koftori, D., Chen, L. & Bowness, G. R. Is Behçet's disease a 'class 1-opathy'? The role of HLA-B*51 in the pathogenesis of Behçet's disease. *Clin. Exp. Immunol.* **191**, 11–18 (2018).
62. Nishiura, H. et al. Estimation of the asymptomatic ratio of novel coronavirus infections (COVID-19). *Int. J. Infect. Dis.* **94**, 154–155 (2020).
63. Chen, T. & Guestrin, C. XGBoost: A scalable tree boosting system. in *Proc. ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* <https://doi.org/10.1145/2939672.2939785> (2016).
64. Stekhoven, D. J. & Bühlmann, P. MissForest-Non-parametric missing value imputation for mixed-type data. *Bioinformatics* **28**, 112–118 (2012).
65. Li, L. & Lin, H. T. Ordinal regression by extended binary classification. *Adv. Neural Inf. Process. Syst.* <https://doi.org/10.7551/mitpress/7503.003.0113> (2007).
66. Lundberg, S. M. & Lee, S. I. A unified approach to interpreting model predictions. *Adv. Neural Inf. Process. Syst.* (2017).
67. Freed, D. N., Aldana, R., Weber, J. A. & Edwards, J. S. The Sentieon Genomics Tools—a fast and accurate solution to variant calling from next-generation sequence data. *bioRxiv* (2017).
68. Chen, X. et al. Manta: Rapid detection of structural variants and indels for germline and cancer sequencing applications. *Bioinformatics* **32**, 1220–1222 (2016).
69. Almarri, M. A. et al. Population structure, stratification, and introgression of human structural variation. *Cell* **182**, 189–199 (2020).
70. McLaren, W. et al. The ensembl variant effect predictor. *Genome Biol.* <https://doi.org/10.1186/s13059-016-0974-4> (2016).
71. Chang, C. C. et al. Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 7 (2015).
72. Manichaikul, A. et al. Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867–2873 (2010).
73. Zhan, X., Hu, Y., Li, B., Abecasis, G. R. & Liu, D. J. RVTESTS: An efficient and comprehensive tool for rare variant association analysis using sequence data. *Bioinformatics* **32**, 1423–1426 (2016).
74. Zhou, W. et al. Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nat. Genet.* **50**, 1335–1341 (2018).
75. Bulik-Sullivan, B. K. et al. LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).