

Augmented Curation of Unstructured Clinical Notes from a Massive EHR System Reveals Specific Phenotypic Signature of Impending COVID-19 Diagnosis

FNU Shweta^{1*}, Karthik Murugadoss^{2*}, Samir Awasthi², AJ Venkatakrisnan², Arjun Puranik², Martin Kang², Brian W. Pickering¹, John C. O'Horo¹, Philippe R. Bauer¹, Raymund R. Razonable¹, Paschalis Vergidis¹, Zelalem Temesgen¹, Stacey Rizza¹, Maryam Mahmood¹, Walter R. Wilson¹, Douglas Challener¹, Praveen Anand², Matt Liebers², Zainab Doctor², Eli Silvert², Hugo Solomon², Tyler Wagner², Gregory J. Gores¹, Amy W. Williams¹, John Halamka¹, Venky Soundararajan²⁺, Andrew D. Badley¹⁺

¹ Mayo Clinic, Rochester MN, USA

² Inference, Cambridge MA, USA

* Joint first authors

+ Address correspondence to: ADB (Badley.Andrew@mayo.edu), VS (venky@inference.net)

Understanding the temporal dynamics of COVID-19 patient phenotypes is necessary to derive fine-grained resolution of the pathophysiology. Here we use state-of-the-art deep neural networks over an institution-wide machine intelligence platform for the augmented curation of 8.2 million clinical notes from 14,967 patients subjected to COVID-19 PCR diagnostic testing. By contrasting the Electronic Health Record (EHR)-derived clinical phenotypes of COVID-19-positive (COVID_{pos}, n=272) versus COVID-19-negative (COVID_{neg}, n=14,695) patients over each day of the week preceding the PCR testing date, we identify diarrhea (2.8-fold), change in appetite (2-fold), anosmia/dysgeusia (28.6-fold), and respiratory failure (2.1-fold) as significantly amplified in COVID_{pos} over COVID_{neg} patients. The specific combination of cough and diarrhea has a 4-fold amplification in COVID_{pos} patients during the week prior to PCR testing, and along with anosmia/dysgeusia, constitutes the earliest EHR-derived signature of COVID-19 (4-7 days prior to typical PCR testing date). This study introduces an Augmented Intelligence platform for the real-time synthesis of institutional knowledge captured in EHRs. The platform holds tremendous potential for scaling up curation throughput, with minimal need for training underlying neural networks, thus promising EHR-powered early diagnosis for a broad spectrum of diseases.

Coronavirus disease 2019 (COVID-19) is a respiratory infection caused by the novel Severe Acute Respiratory Syndrome coronavirus-2 (SARS-CoV-2). As of April 15, 2020, according to WHO there have been more than 1.9 million confirmed cases worldwide and more than 123,000 deaths attributable to COVID-19 (<https://covid19.who.int/>) The clinical course and prognosis of patients with COVID-19 varies substantially, even among patients with similar age and comorbidities¹. Following exposure and initial infection with SARS-CoV-2, likely through the upper respiratory tract, patients can remain asymptomatic although active viral replication may be present for weeks before symptoms manifest^{1,2}. The asymptomatic nature of initial SARS-CoV-2 infection in the majority of patients may be exacerbating the rampant community transmission observed³. It remains unknown which patients become symptomatic, and in those who do, the timeline of symptoms remains poorly characterized and non-specific. Symptoms may include fever, fatigue, myalgias, loss of appetite, loss of smell (anosmia), and altered sense of taste, in addition to the respiratory symptoms of dry cough, dyspnea, sore throat, and rhinorrhea, and well as gastrointestinal symptoms of diarrhea, nausea, and abdominal discomfort⁴. A small proportion of COVID-19 patients progress to severe illness requiring hospitalization or intensive care management; among these individuals, mortality owing to Acute Respiratory Distress Syndrome (ARDS) is higher⁵. The estimated average time from symptom-onset to resolution can range from three days to more than three weeks, with a high degree of variability⁶. The COVID-19 public health crisis demands a data science-driven and temporal pathophysiology-informed precision medicine approach for its effective clinical management.

Here we introduce a platform for the augmented curation of the full-spectrum of patient phenotypes from 8,22,9092 clinical notes of the Mayo Clinic EHRs for 14,967 patients with confirmed positive/negative COVID-19 diagnosis by PCR testing (see **Methods**). The platform utilizes state-of-the-art transformer neural networks on the unstructured clinical notes to automate entity recognition (e.g. diseases, drugs, phenotypes), quantify the strength of contextual associations between entities, and characterize the nature of association into positive, negative, or other sentiments. We identify specific gastro-intestinal, respiratory, and sensory phenotypes, as well as some of their specific combinations, that appear to be indicative of impending COVID_{pos} diagnosis by PCR testing. This highlights the potential for neural networks-powered EHR curation to facilitate a significantly earlier diagnosis of COVID-19 than currently thought feasible.

Results

The clinical determination of the COVID-19 status for each patient was conducted using the SARS-CoV-2 PCR (RNA) test approved for human nasopharyngeal and oropharyngeal swab specimens under the U.S. FDA emergency use authorization (EUA)⁶. This PCR test resulted in 14,695 COVID_{neg} patient diagnoses and 272 COVID_{pos} patient diagnoses. In order to investigate the time course of COVID-19 progression in patients, we used BERT-based deep neural networks to extract symptoms and their putative synonyms from the clinical notes for a few weeks prior to, and a few weeks post, the date when the COVID-19 diagnosis test was taken (see **Methods; Table 1**). For the purpose of this analysis, all patients were temporally aligned, by setting the date of COVID-19 PCR testing to 'day 0', and the proportion of patients demonstrating each symptom derived from the EHR over each day of the week preceding and post PCR testing was tabulated (**Table 2**). As a negative control, we included a non-COVID-19 symptom 'dysuria'.

In the COVID_{pos} patients, diarrhea occurs in 43 of 272 patients (15.8%) in the week prior to the PCR testing date, whereas in the COVID_{neg} patients, only 822 of 14,695 patients (5.6%) have confirmed diarrhea in the week prior to the PCR test date. The amplified probability (**Table 1**; 2.8-fold; p-value = 8.4E-13) of diarrhea in the week preceding PCR testing for COVID_{pos} patients is quite noteworthy. Some of these undiagnosed COVID-19 patients that experience diarrhea may be unintentionally shedding SARS-CoV-2 fecally⁷. Incidentally, epidemiological surveillance by waste water monitoring conducted recently in the state of Massachusetts observed copious SARS-CoV-2 RNA⁸. The amplification of diarrhea in COVID_{pos} over COVID_{neg} patients in the week preceding PCR testing highlights the importance and necessity for washing hands often.

Change in appetite/intake is amplified in the week preceding PCR testing in COVID_{pos} over COVID_{neg} patients (**Table 1**, 2.0-fold amplification; p-value = 0.0026). Altered or diminished sense of taste or smell (dysgeusia or anosmia) is significantly amplified in COVID_{pos} over COVID_{neg} patients in the week preceding PCR testing (**Table 1**; 28.6-fold amplification; p-value = 5.1E-36). This result suggests that anosmia is likely a significant early indicator of COVID-19, including in otherwise asymptomatic patients.

Respiratory failure is modestly enriched in the week prior to PCR testing in COVID_{pos} over COVID_{neg} patients (2.1-fold amplification; p-value = 0.01; **Table 1**). Among other common phenotypes, diaphoresis manifests in 31 of 272 patients (11.4%) and fatigue in 37 of 272 patients (13.6%) during the week prior to COVID_{pos} PCR testing. In contrast, for the COVID_{neg} patients, diaphoresis occurs in 825 of 14,695 patients (5.6%) and fatigue in 1,279 of 14,695 (8.7%) during the week prior to the PCR test. This corresponds to a 2-fold amplification of diaphoresis (p-value = 4.7E-05) and 1.56-fold amplification of fatigue (p-value = 0.005). Headache occurs in 35 of 272 COVID_{pos} patients (12.9%) and in 1,023 of 14,695 COVID_{neg} patients (7.0%), reflecting a 1.9-fold amplification in COVID_{pos} patients in the week prior to the PCR test (p-value = 0.0002). Cough has a 1.3-fold amplification (p-value = 0.01) in COVID_{pos} over COVID_{neg} patients in the week preceding PCR testing (**Table 1**). Fever/chills occur in 67 of 272 COVID_{pos} patients (24.6%) and in 2726 of 14,695 COVID_{neg} patients (18.6%) in the week prior to the PCR test. This suggests that fever/chills

is somewhat nonspecific to COVID-19 patients. Finally, dysuria was included as a negative control for COVID-19, and consistent with this assumption, 1 out of 272 COVID_{pos} patients (0.4%) and 91 out of 14,695 COVID_{neg} patients (0.6%) had dysuria during the week preceding PCR testing (**Table 1**).

Next, we considered the 351 possible pairwise conjunctions of 27 phenotypes for COVID_{pos} versus COVID_{neg} patients in the week prior to the PCR testing date (**Table S4**). Given that an altered sense of smell or taste (anosmia/dysgeusia) occurs in very few of these conjunctions within the COVID_{pos} patients till date, and the fact that independently anosmia/dysgeusia is already a significant signature of impending COVID_{pos} diagnosis (based on the above results), here we only remark on the other 325 possible pairwise symptom combinations. The combination of cough and diarrhea is noted to be particularly significant in COVID_{pos} over COVID_{neg} patients during the week preceding PCR testing; i.e. cough and diarrhea co-occur in 36 of 272 COVID_{pos} patients (13.2%) and in 486 of 14,695 COVID_{neg} patients (3.3%) indicating a 4-fold amplification of this specific symptom combination as a signature of impending COVID-19 diagnosis (BH corrected p-value = 9.3E-19). Another enriched combination of symptoms in COVID_{pos} over COVID_{neg} patients in the week preceding PCR testing is diaphoresis and diarrhea that co-occur in 21 of 272 COVID_{pos} (7.7%) and 204 of 14,695 COVID_{neg} (1.4%) patients. This corresponds to a 5.6-fold enrichment (BH corrected p-value = 1.8E-17) in the COVID_{pos} patient group and suggests diaphoresis and diarrhea as another symptom combination preceding COVID_{pos} PCR test results.

We further investigated the temporal evolution of the proportion of patients with each symptom over the week prior to PCR testing. Cough and diarrhea were found to be the early indicators that significantly discriminate COVID_{pos} from COVID_{neg} patients. In particular, between 3 to 7 days prior to PCR testing, cough is amplified in the COVID_{pos} patient cohort over the COVID_{neg} patient cohort with an amplification of 5.5-fold (day -7, p-value = 7.0E-16), 5.3-fold (day -6, p-value = 5.4E-15), 4.7-fold (day -5, p-value = 1.4E-10), 3.9-fold (day -4, p-value = 3.6E-08), and 3.8-fold (day -3, p-value = 7.3E-08). The intriguing diminishing odds of cough as a symptom from 7 to 3 days preceding the PCR testing date, suggests this may be a notable temporal pattern. Likewise, diarrhea is amplified in the COVID_{pos} patient cohort over the COVID_{neg} patient cohort with an amplification of 5.7-fold (day -7, p-value = 6.1E-11), 5-fold (day -6, p-value = 5.5E-08), 3.2-fold (day -5, p-value = 3E-03), and 4.2-fold (day -4, p-value = 7E-06). Following the enriched odds of diarrhea and cough, we further find that change in appetite may be considered a subsequent symptom of impending COVID-19 diagnosis. This is because, change in appetite is amplified in the COVID_{pos} cohort over the COVID_{neg} cohort on day -4 (3.4-fold, p-value = 4.4E-03), day -3 (3.7-fold, p-value = 2.9E-03), and day -2 (4.2-fold, p-value = 6.5E-05). Finally, as the day of PCR testing ensues, fever/chills are enriched in the COVID_{pos} over the COVID_{neg} cohort, with 1.6-fold (day 0, p-value =), 2.4-fold (day 1, p-value =) and 2-fold (day 2, p-value = 1.7E-04) respectively. Similarly, cough is also amplified in the COVID_{pos} over the COVID_{neg} cohort on day 0 (1.9-fold, p-value = 2.1E-06), day 1 (2.8-fold, p-value = 2.3E-12) and day 2 (2-fold, p-value = 2.8E-03) post the PCR testing date. These observations characterize the temporal evolution of specific phenotypes that are enriched in COVID_{pos} patients preceding and post the PCR testing date.

While explicit identification of SARS-CoV-2 in patients prior to the PCR testing date was not conducted, such prospective validation of our augmented EHR curation approach is being initiated. Nevertheless, this high-resolution temporal overview of the EHR-derived clinical phenotypes as they relate to the SARS-CoV-2 PCR diagnostic testing date for 14,967 patients has revealed specific enriched signals of impending COVID-19 onset. These clinical insights can help modulate social distancing measures and appropriate clinical care for individuals exhibiting the specific gastro-intestinal (diarrhea, change in appetite/intake), sensory (anosmia, dysgeusia) and respiratory phenotypes identified herewith, including for patients awaiting conclusive COVID-19 diagnostic testing results (e.g. by SARS-CoV-2 RNA RT-PCR).

Discussion

In order to identify potential cells and tissue types that may be associated with the EHR-derived clinical phenotypes observed above for COVID-19 patients, we analyzed Single Cell RNA-seq data using the nferX platform (see **Methods**)⁹. Given recent studies implicating the necessity of both ACE2 and TMPRSS2 for the SARS-CoV-2 lifecycle¹⁰, we scouted for human cells that co-express both genes. This co-expression analysis revealed that specific cell types from the small intestine/colon, nasal cavity, respiratory system, pancreas, urinary tract, and gallbladder co-express both ACE2 and TMPRSS2 (**Figure 1, Figure S1**). Notably, multiple small intestine cell types co-express the two genes. These cell types include enterocytes, enteroendocrine cells, stem cells, goblet cells, and Paneth cells. In the pancreas, the cell types included ductal cells and acinar cells. The kidney cells co-expressing TMPRSS2 and ACE2 include proximal tubular cells, pelvic epithelial cells and type A intercalated cells. Co-expression of TMPRSS2 and ACE2 is also observed in the epithelial cells of the olfactory nasal cavity and the respiratory tract as well as in type II pneumocytes (albeit at comparatively lower level). While the identified tissues showing co-expression of ACE2 and TMPRSS2 in the gastro-intestinal, respiratory, and sensory systems correlate with the clinical phenotypes of early COVID-19 infection as described above, these insights are conceivably from normal/healthy tissues. This highlights the need for meticulous bio-banking of COVID-19 patient-derived biospecimen and their characterization via single cell RNA-seq and other molecular technologies.

Primary prevention is the most effective method to minimize spread of contagious infectious viruses such as SARS-CoV-2 (**Figure S2**). In addition to population-based strategies such as social distancing, there are significant ongoing efforts to develop a prophylactic solution (**Table S1**). As the immunodominant humoral immune response in patients is directed against the SARS-CoV2 spike protein, many vaccines under investigation target this viral protein. It remains to be determined whether anti-spike protein antibodies induced by natural infection or by vaccines induce neutralizing antibody responses. Chloroquine and its analogues have been shown to inhibit virus replication in-vitro²⁸. Whether Chloroquine or Hydroxychloroquine have meaningful effects of SARS-CoV2 replication in patients remains to be understood, and are the subject of clinical trials, both as post-exposure prophylaxis and as treatment (**Table S1**). Hydroxychloroquine was approved by FDA for emergency use in hospitalized COVID-19 patients who are not eligible for clinical trials on April 7, 2020 based on limited clinical data, but concerns have been raised about toxicity and risk of sudden death²⁹.

Our findings from the EHR analysis of COVID-19 progression can aid in a human pathophysiology enabled summary of the experimental therapies being investigated for COVID-19 (**Figure 2, Table S1**). Some of the earliest phases of intervention attempt to inhibit the entry/replication of SARS-CoV-2 by modulating critical host targets (e.g. renin angiotensin aldosterone system/RAAS inhibitors, ACE2 analogs, serine protease inhibitors) or directly inhibiting the function of viral proteins (e.g. viral RNA-dependent RNA polymerase inhibitors, protease inhibitors, convalescent plasma, synthetic immunoglobulins) (**Box 1, Table S2**). In patients with more advanced stages of disease progression, who suffer from respiratory abnormalities, therapeutics are being advanced to target the inflammatory response that can lead to Acute Respiratory Disease Syndrome (ARDS) and is associated with high mortality (**Box 1**). These include anti-GM-CSF agents, anti-IL-6 agents, JAK inhibitors, and complement inhibitors. Another emerging option for patients at this stage is convalescent plasma, which has shown some clinical benefits in cases of COVID-19 and related viral diseases (SARS-1, MERS) at various stages of severity (**Box 1**). Administration of convalescent plasma containing active specific antiviral antibodies may prevent or attenuate progression to severe disease. Expanded access to convalescent plasma for treatment of patients with COVID-19 has been approved by the FDA for emergency IND use and is available through a nationwide program led by Mayo Clinic (**Box 1**).

In those who become symptomatic, it is imperative that diagnostic testing is done, at dedicated testing sites if available, to confirm diagnosis (**Figure S2**). Meanwhile, patients are recommended to self-quarantine at home, use mask protection when social distancing cannot be obtained, and continue supportive measures. For patients with mild symptoms, such measures may be sufficient given the self-limited nature of viral syndromes. In the event of symptom exacerbation, often marked by worsening respiratory distress, medical evaluation is warranted, and possible hospitalization. The mainstay of treatment for COVID-19, remains supportive care, and as needed supplemental oxygen. Experimental therapies intended to block SARS-CoV2 viral entry and inhibit steps in the viral life cycle necessary for viral replication have been proposed at this early stage (**Figure 2**). The goal of these therapies is to reduce viral load, thus reducing the chance of overwhelming immune reaction by delaying progression of the disease.

Among the proposed treatment options for COVID-19, corticosteroid should be avoided outside a clinical trial, as suggested by the IDSA, until further clinical evidence can be established (www.idsociety.org/practice-guideline/covid-19-guideline-treatment-and-management). This is because there has been conflicting evidence and guidance on steroid use in COVID-19³⁰. While steroids can play a role in control of inflammation, a collection of clinical evidence from steroid use in other coronavirus outbreaks suggest that the use of corticosteroids might exacerbate COVID-19-associated lung injury³¹.

As patients progress to severe or critical diseases, the primary objective of COVID-19 management is to provide respiratory support and control immune overactivation (**Figure 3, Figure S3**). Patients whose condition deteriorates to critical status primarily decompensate from a respiratory standpoint, but may also develop multi-organ failure (respiratory failure, cardiac failure, renal failure, hypercoagulable state, thrombotic microangiopathy), as well as severe inflammatory responses similar to cytokine release syndrome and eventually reactive hemophagocytic lymphohistiocytosis syndrome. A major manifestation of respiratory decompensation and cytokine release syndrome is acute respiratory distress syndrome (ARDS). Critical care support such as mechanical support from noninvasive to invasive mechanical ventilation and in, some instances, extracorporeal support, vasopressors, renal replacement therapy, anticoagulation, and are paramount to survival of these critically ill patients per SCC guidelines (SCCM/ESICM 2020). On the other hand, drugs such as immunomodulatory agents often used to treat cytokine release syndrome, may allow for some degree of improvement or recovery either leading into or during severe and critical disease (**Figure 2**).

This study demonstrates how the highly unstructured institutional knowledge can be synthesized using deep learning and neural networks³². Expanding beyond one institution's COVID-19 diagnostic testing and clinical care to the EHR databases of other academic medical centers and health systems will provide a more comprehensive view of clinical phenotypes enriched in COVID_{pos} over COVID_{neg} patients in the days preceding confirmed diagnostic testing. This requires leveraging a privacy-preserving federated software architecture that enables each medical center to retain the span of control of their de-identified EHR databases, while enabling the machine learning models from partners to be deployed in their secure cloud infrastructure. Such seamless multi-institute collaborations over an Augmented Intelligence platform that puts patient privacy and HIPAA-compliance first, is being advanced actively over the Mayo Clinic's Clinical Data Analytics Platform Initiative (CDAP). The capabilities demonstrated in this study for rapidly synthesizing over 8.2 million unstructured clinical notes to develop an EHR-powered clinical diagnosis framework will be further strengthened through such a universal biomedical research platform.

A caveat of relying solely on EHR inference is that mild phenotypes that may not lead to a presentation for clinical care, such as anosmia, may go unreported in otherwise asymptomatic patients. As at-home serology-based tests for COVID-19 with high sensitivity and specificity are approved, capturing these symptoms will become increasingly important in order to facilitate the continued development and refinement of disease models. EHR-integrated digital health tools may help address this need.

As we continue to understand the diversity of COVID-19 patient outcomes through holistic inference of EHR systems, it is equally important to invest in uncovering the molecular mechanisms and gain cellular/tissue-scale pathology insights through large-scale patient-derived biobanking and multi-omics sequencing. As the anecdotal single cell RNA-seq (scRNA-seq) based co-expression analysis of ACE2 and TMPRSS2 on normal human samples conducted here highlights, the rich heterogeneity of cell types constituting various host tissues can be investigated in great detail by scRNA-seq. To correlate patterns of molecular expression from scRNA-seq with EHR-derived phenotypic signals of COVID-19 disease progression, a large-scale bio-banking system has to be created. Such a system will enable deep molecular insights into COVID-19 to be gleaned and triangulated with SARS-CoV-2 tropism and patient outcomes.

Ultimately, connecting the dots between the temporal dynamics of COVID_{pos} and COVID_{neg} clinical phenotypes across diverse patient populations to the multi-omics signals from patient-derived bio-specimen will help advance a more holistic understanding of COVID-19 pathophysiology. This will set the stage for a precision medicine approach to the diagnostic and therapeutic management of COVID-19 patients.

BOX 1

Experimental therapies targeting entry and replication of SARS-CoV-2

RAAS inhibitors and ACE2 analogs: One class of experimental therapies intended to inhibit viral entry and early disease in COVID-19 includes Renin Angiotensin Aldosterone System (RAAS) inhibitors and recombinant ACE2 (Table S2). ACE2 is the primary host receptor for SARS-CoV-2, while serine protease TMPRSS2 is implicated in the spike protein priming after viral binding¹⁰. Recombinant ACE2 has been proposed as an early COVID-19 therapy based on in-vitro data¹¹. At this time, the effect of RAAS inhibitors is uncertain in the context of COVID-19. Studies have investigated how ACE expression is modulated by coronavirus infection, and how that relates to lung injury¹¹. Trials are ongoing with Angiotensin Receptor Blockers (ARBs) for treatment of COVID-19 by diminishing downstream harmful effects of angiotensin receptor activation (Figure 2).

Serine Protease inhibitors: Given the TMPRSS2 involvement in viral entry (Figure 2), serine protease inhibitors such as Camostat are now under evaluation in trials and should also be considered in the early stages of SARS-CoV-2 infection.

Viral RNA-dependent RNA polymerase inhibitors: Of these, Remdesivir, a nucleoside analog, has attracted much attention for in-vivo inhibition of SARS-CoV-2, and a recent observational study of 53 patients who received Remdesivir under compassionate use found that 68% of patients demonstrated improvement in respiratory status after a 10 day regimen¹². Another nucleoside analog, Galidesivir, is also under evaluation in patients. Yet another viral replication inhibitor in clinical trials is Favipiravir (Figure 2). Favipiravir is a broad spectrum viral RNA dependent RNA polymerase inhibitor that is shown to have in-vivo activity against a wide range of RNA viruses. In one RCT of 240 patients, Favipiravir was found to improve the clinical recovery rate of COVID-19 relative to Umifenovir, a viral entry inhibitor¹³ (Table S1).

HIV Protease inhibitors: This class of medication is widely proposed and used off-label based on postulates that HIV and HCV proteases share structural similarities with those of SARS-CoV-2¹⁴. Of these, Lopinavir/Ritonavir (combination) has shown promise but was found to have a non-significant benefit in a Randomized Clinical Trial (RCT) of 199 patients in China¹⁵, while Darunavir has shown no significant activity against SARS-CoV-2 in-vitro (Table S1)¹⁶. Multiple randomized, controlled clinical trials are now underway in the USA to determine efficacy of these drugs in the treatment of COVID-19.

Other Antiviral Agents: Another emerging option for patients at this stage is convalescent plasma (Figure 2), which has shown clinical benefits in cases of COVID-19¹⁷ and related viral diseases (SARS-1, MERS) at various stages of severity^{18,19}. Administration of convalescent plasma containing specific antiviral antibodies may prevent or attenuate progression to severe disease. Expanded access to convalescent plasma for treatment of patients with COVID-19 is available through a program led by Mayo Clinic²⁰. Synthetic hyperimmune globulins are also under development and evaluation.

Agents being advanced that target the inflammatory response in COVID-19

Anti-GM-CSF agents -- A xenograft study found that granulocyte monocyte colony stimulating factor (GM-CSF) neutralization with Lenzilumab significantly reduced production of inflammatory cytokines²¹, offering evidence for efficacy of anti-GM-CSF agents in prevention of CART-induced cytokine release syndrome (CRS). Lenzilumab has been approved by the FDA for emergency IND use for CRS in COVID-19, while others such as Mavrilimumab and Gimsilumab aimed at controlling undesired inflammation from myeloid activation will be evaluated in clinical trials.

Anti-IL-6 agents: IL-6 is a pro-inflammatory cytokine, regarded as a driver of CRS²² (Figure 2A-C). A recent report suggests IL-6 as a biomarker for respiratory failure in COVID-19²². As such, anti-IL-6 agents including Tocilizumab, Sarilumab and Siltuximab are being evaluated in randomized trials (Table S1), and used off-label in severe COVID-19 patients. Tocilizumab was approved for the treatment of CRS in 2017. An observational study of 21 patients with severe COVID-19 pneumonia treated with Tocilizumab showed promising results^{23,24}.

Anti-JAK agents: A number of immunomodulatory agents not linked to CRS are also under trial for COVID-19 (Figure 2). Janus kinase (JAK) inhibitors such as Baricitinib, Fedratinib, and Ruxolitinib, indicated for Rheumatoid Arthritis and Myelofibrosis, have been tested in xenograft models for Chimeric Antigen Receptor (CAR) T-cell therapy induced CRS²⁵. Ruxolitinib is available under an expanded access program in USA for severely ill COVID-19 patients (Table S1) and trials are underway in other countries.

Anti-Complement agents: A recent study found that SARS-CoV-2 also binds to MASP2, a key driver of the complement activation pathway, leading to complement hyperactivation in COVID-19 patients²⁶. Inhibitors of the terminal complement pathway such as Eculizumab have been tried in individuals with improvements observed after administration in China.

Agents targeting ventilation/perfusion defects in COVID-19-induced ARDS

Vasodilators: A recent report based on 16 cases in Italy and Germany noted that, contrary to the established understanding in ARDS, COVID-19 patients in ARDS retain relatively high lung compliance²⁷ and demonstrate ventilation/perfusion defects likely arising from perfusion dysregulation and hypoxic vasoconstriction. Therefore, patients with COVID-19 in ARDS may benefit from vasodilators to address this pathophysiologic mechanism. A trial is underway in China for use of inhaled nitric oxide in patients with mechanical ventilation (Table S1).

Methods

Augmented curation of SARS-CoV2-positive patient charts

The nferX Augmented Curation technology was leveraged to rapidly curate the charts of SARS-CoV-2-positive patients. First, we read through the charts of 100 patients and identified and grouped symptoms into sets of synonymous words and phrases. For example, “SOB”, “shortness of breath”, and “dyspnea”, among others, were grouped into “shortness of breath”. We did the same for diseases and medications. For the SARS-CoV2-positive patients, we identified a total of 26 symptom categories (**Table S3**) with 145 synonyms or synonymous phrases. Together, these synonyms and synonymous phrases capture a multitude of ways that symptoms related to COVID-19 are described in the Mayo Clinic Electronic Health Record (EHR) databases.

Next, for charts that had not yet been manually curated, we used state-of-the-art BERT-based neural networks³² to classify symptoms as being present or not present based on the surrounding phraseology. The neural network used to perform this classification was trained using nearly 250 different phenotypes and 20000 sentences; it achieves over 96% recall for positive/negative sentiment classification. We went through individual sentences and either accepted the sentences or rejected and reclassified them. The neural networks were actively re-trained as curation progressed, leading to stepwise increases in curation efficiency and model accuracy. In step 1 of this process, we labeled 11433 sentences, 8737 of which were labeled as either ‘present’ or ‘not present.’ The model trained on this data set (80%-20% training/test split) achieved F1 scores of 0.93 and 0.84 for ‘present’ and ‘not present’ classifications, respectively. The model was then applied to an additional 3688 sentences in step 2, rapidly corrected by a human for classification errors and re-trained to generate a newer version of the model. Step 3 was an iteration of step 2 on an additional 3369 sentences. The model achieved F1 scores of 0.96/0.91 after step 2 and 0.96/0.96 after step 3 for the classification of ‘present’/‘not present.’ Due to the augmented nature of this approach, steps 2 and 3 required successively less input from the human annotator.

This model was applied to 80,148 clinical notes from the 272 COVID_{pos} patients and 8.2 million clinical notes from the 14,695 COVID_{neg} patients. First, the difference between the date on which a particular note was written and the PCR testing date of the patient corresponding to that note formed the relative date measure for that note. The PCR testing date was treated as ‘day 0’ with notes preceding it assigned ‘day-1’, ‘day-2’ and so on. BERT-based neural networks were applied on each note to provide a set of symptoms that were present at that point of time for the patient in question. This map was then inverted to determine for each symptom and relative date the set of unique patients experiencing that symptom.

For each synonymous group of symptoms, we computed the count and proportion of COVID_{pos} and COVID_{neg} patients that were deemed to have that symptom in at least one note between 1 and 7 days prior to their PCR test. We additionally computed the ratio of those proportions which indicates the extent of prevalence of the symptom in the COVID_{pos} cohort as compared to the COVID_{neg} cohort. A standard 2-proportion z hypothesis test was performed, and a p-value was reported for each symptom.

To capture the temporal evolution of symptoms in the COVID_{pos} and COVID_{neg} cohorts, the process described above was repeated considering counts and proportions for each day independently.

Pairwise analysis of phenotypes was performed by considering 351 phenotypic pairs from the original set of 27 individual phenotypes. For each pair, we calculated the number of patients in the COVID_{pos} and COVID_{neg} cohorts wherein both phenotypes occurred at least once in the week preceding PCR testing. With these patient proportions, a 2-proportion z test p-value was computed. Benjamini-Hochberg correction was applied to account for multiple hypothesis testing.

This research was conducted under IRB 20-003278, “Study of COVID-19 patient characteristics

with augmented curation of Electronic Health Records (EHR) to inform strategic and operational decisions”. All analysis of EHRs was performed in the privacy-preserving environment secured and controlled by the Mayo Clinic. nference and the Mayo Clinic subscribes to the basic ethical principles underlying the conduct of research involving human subjects as set forth in the Belmont Report and strictly ensures compliance with the Common Rule in the Code of Federal Regulations (45 CFR 46) on the Protection of Human Subjects.

Analysis of cell-types expressing ACE2 and TMPRSS2 using single cell RNAseq

Since the successful entry of virus in the cell requires priming by cellular host protease – TMPRSS2, we hypothesized that cells that express both TMPRSS2+ and ACE2+ cells could harbor SARS-CoV-2 during the course of infection. Thus, we probed for the expression of ACE2 and TMPRSS2 in all the single-cell studies from human tissues available on the nferX Single Cell platform (<https://academia.nferx.com/>). For all the tissues that we profiled, we ensured that there are a minimum of 100 cells in the cell population and that there is a minimum of 1% of the cells in the cell population co-expressing (non-zero expression) both TMPRSS2 and ACE2 expression.

Figure Legends

Figure 1. Clinical phenotypes of COVID-19 and their connection to single cell RNA-seq co-expression of ACE2-TMPRSS2. Severity of COVID-19 and associated clinical conditions are shown. Cell types co-expressing SARS-CoV-2 infectivity determinants ACE2 and TMPRSS2 determined by single cell RNA-seq are mapped onto the COVID-19 pathophysiology summary.

Figure 2. Pathophysiology of COVID-19, associated treatments, and the underlying molecular mechanisms. While there is no established treatment strategy for COVID-19, several classes of therapeutics have emerged for the medical management of the disease, on the basis of their known mechanisms of action and the pathophysiology of COVID-19.

Acknowledgments

We thank Murali Aravamudan, Ajit Rajasekharan, and Rakesh Barve for their thoughtful review and feedback on this manuscript. We also thank Andrew Danielsen, Jason Ross, Jeff Anderson, Ahmed Hadad, and Sankar Ardhanari for their support that enabled the rapid completion of this study.

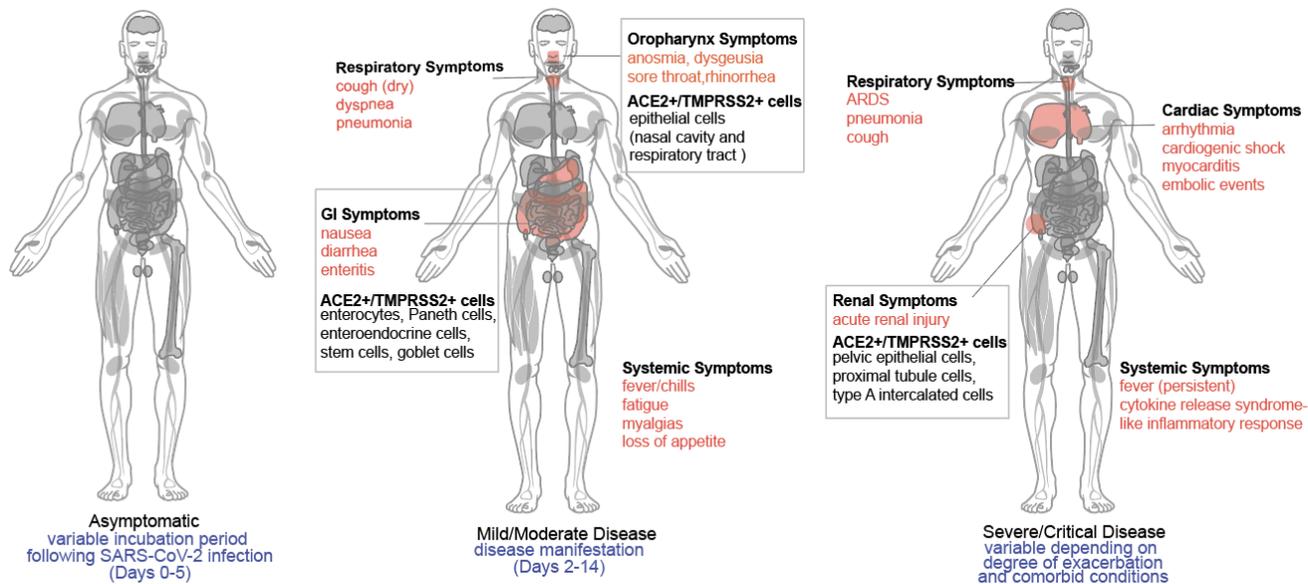


Figure 1

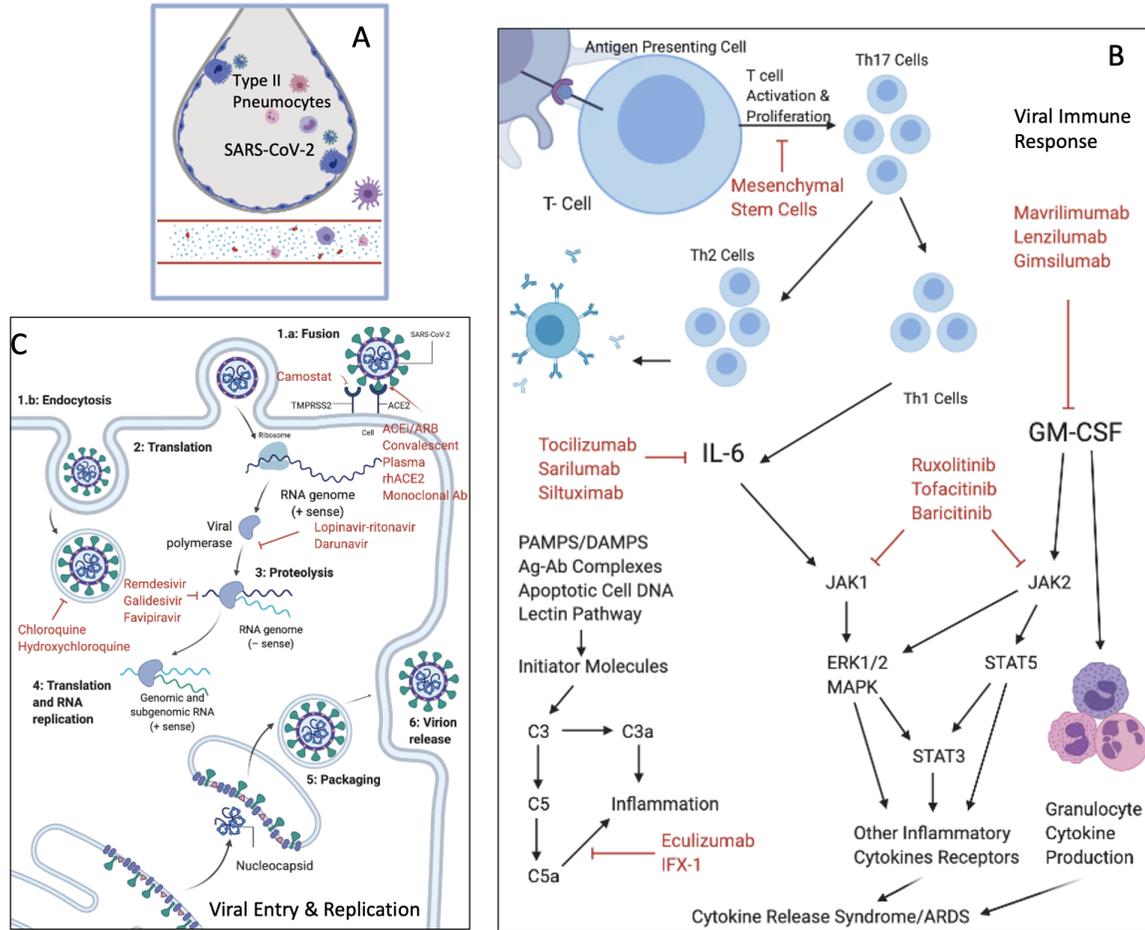


Figure 2

Table 1. Augmented curation of the unstructured clinical notes from the EHR reveals specific clinically confirmed phenotypes that are amplified in COVID_{pos} patients over COVID_{neg} patients in the week prior to the SARS-CoV-2 PCR testing date. The key COVID_{pos} amplified phenotypes in the week preceding PCR testing (i.e. day = -7 to day = -1) are highlighted in gray. The ratio of COVID_{pos} to COVID_{neg} proportions represents the fold change amplification of each phenotype in the COVID_{pos} patient set.

Phenotype	COVID _{pos} count (N=272)	COVID _{neg} count (N = 14695)	COVID _{pos} proportion (N=272)	COVID _{neg} proportion (N=14695)	(COVID _{pos} / COVID _{neg}) relative ratio	2-tailed p-value
Altered or diminished sense of taste or smell	9	17	0.03	0.00	28.60	5.07E-36
Diarrhea	43	822	0.16	0.06	2.83	8.45E-13
Respiratory Failure	12	309	0.04	0.02	2.10	0.010
Diaphoresis	31	825	0.11	0.06	2.03	4.71E-05
Change in appetite/intake	18	485	0.07	0.03	2.01	0.0026
Hemoptysis*	2	55	0.01	0.00	1.96	0.34
Dermatitis	10	285	0.04	0.02	1.90	0.04
Headache	35	1023	0.13	0.07	1.85	0.0002
Conjunctivitis*	2	59	0.01	0.00	1.83	0.39
Fatigue	37	1279	0.14	0.09	1.56	0.005
Otitis*	1	39	0.00	0.00	1.39	0.75
Cough	68	2766	0.25	0.19	1.33	0.010
Fever / chills	67	2726	0.25	0.19	1.33	0.011
Chest pain/pressure	38	1690	0.14	0.12	1.21	0.21
Myalgia/Arthralgia	43	1974	0.16	0.13	1.18	0.26
Generalized symptoms	51	2352	0.19	0.16	1.17	0.22
GI upset	23	1166	0.08	0.08	1.07	0.75
Productive cough	11	562	0.04	0.04	1.06	0.85
Congestion	22	1268	0.08	0.09	0.94	0.75
Neurological symptoms	14	820	0.05	0.06	0.92	0.76
Respiratory difficulty	31	1951	0.11	0.13	0.86	0.36
Rhinitis	12	865	0.04	0.06	0.75	0.30
Wheezing	10	789	0.04	0.05	0.68	0.22
Dysuria	1	91	0.00	0.01	0.59	0.60
Pharyngitis	9	968	0.03	0.07	0.50	0.030
Cardiac symptoms	1	134	0.00	0.01	0.40	0.35
Dry mouth	0	26	0.00	0.00	0.00	0.49

Table 2. Temporal analysis of the EHR clinical notes for the week preceding PCR testing (i.e. day -7 to day -1), the day of PCR testing (day 0), and the subsequent pair of days (day 1, day 2) in COVID_{pos} and COVID_{neg} patients. Temporal enrichment for each symptom is quantified using the ratio of COVID_{pos} patient proportion over the COVID_{neg} patient proportion for each day. The patient proportions in the rows labeled 'Positive (n = 272)' and 'Negative (n = 14695)' are represented as percentages.

Phenotype	COVID-19 (N = 14967)	Day = -7	Day = -6	Day = -5	Day = -4	Day = -3	Day = -2	Day = -1	Day = 0 (PCR test)	Day = 1	Day = 2
Altered or diminished sense of taste or smell	Positive (n = 272)	1.10	1.10	1.10	0.00	0.37	0.00	0.00	3.68	4.04	2.57
	Negative (n = 14695)	0.01	0.00	0.01	0.03	0.01	0.02	0.06	0.10	0.08	0.06
	Ratio (Positive/Negative)	162.08	-	81.04	0.00	27.01	0.00	0.00	36.02	49.52	42.02
	p-value	6E-28	3.9E-37	2E-22	0.79	4.4E-05	0.81	0.68	2.1E-46	2.2E-61	3.4E-36
Diarrhea	Positive (n = 272)	4.41	3.68	2.21	3.31	2.21	2.21	4.04	5.88	6.25	4.78
	Negative (n = 14695)	0.77	0.73	0.68	0.80	0.84	0.80	3.72	3.30	1.78	1.49
	Ratio (Positive/Negative)	5.74	5.05	3.24	4.16	2.64	2.77	1.09	1.78	3.51	3.21
	p-value	6.2E-11	5.5E-08	3E-03	7E-06	0.02	0.01	0.78	0.02	6.8E-08	1.5E-05
Change in appetite/intake	Positive (n = 272)	1.47	1.47	1.47	1.84	1.84	2.57	1.84	4.41	0.74	1.84
	Negative (n = 14695)	0.55	0.67	0.69	0.54	0.50	0.61	1.68	2.27	1.59	1.03
	Ratio (Positive/Negative)	2.67	2.21	2.14	3.42	3.70	4.20	1.09	1.94	0.46	1.79
	p-value	0.05	0.11	0.13	4.4E-03	2.9E-03	6.5E-05	0.84	6.7E-03	0.26	0.19
Respiratory failure	Positive (n = 272)	1.84	1.84	2.21	2.57	1.84	1.84	1.84	4.04	2.21	2.57
	Negative (n = 14695)	0.62	0.63	0.64	0.71	0.73	0.75	1.67	2.43	1.84	1.59
	Ratio (Positive/Negative)	2.97	2.94	3.45	3.60	2.52	2.46	1.10	1.66	1.20	1.62
	p-value	0.01	0.01	1.7E-03	4.2E-04	0.04	0.04	0.83	0.09	0.66	0.20
Headache	Positive (n = 272)	3.68	1.47	2.21	1.10	1.47	0.74	2.57	6.25	4.41	2.94
	Negative (n = 14695)	0.60	0.69	0.65	0.58	0.65	0.69	4.71	3.23	1.78	1.18
	Ratio (Positive/Negative)	6.14	2.12	3.41	1.91	2.27	1.06	0.55	1.93	2.47	2.48
	p-value	4.5E-10	0.13	1.9E-03	0.26	0.10	0.94	0.10	5.7E-03	1.4E-03	8.8E-03
Cough	Positive (n = 272)	6.99	6.99	5.51	5.51	5.51	2.94	6.62	19.49	15.07	6.99
	Negative (n = 14695)	1.26	1.32	1.16	1.42	1.46	1.64	15.22	10.52	5.30	3.57
	Ratio (Positive/Negative)	5.55	5.29	4.74	3.88	3.77	1.79	0.43	1.85	2.84	1.96
	p-value	6.99E-16	5.44E-15	1.38E-10	3.63E-08	7.3E-08	0.10	8.5E-05	2.1E-06	2.3E-12	2.8E-03
Fatigue	Positive (n = 272)	2.94	2.94	2.94	2.21	2.21	1.84	4.41	8.09	1.84	2.94
	Negative (n = 14695)	1.04	1.03	0.98	0.97	0.94	1.29	5.42	4.34	2.17	1.78
	Ratio (Positive/Negative)	2.82	2.84	3.00	2.28	2.35	1.42	0.81	1.86	0.85	1.66
	p-value	2.61E-03	2.45E-03	1.39E-03	0.04	0.03	0.43	0.47	2.9E-03	0.71	0.15
Fever / chills	Positive (n = 272)	6.25	7.35	3.68	5.15	4.41	3.31	6.99	18.75	15.81	9.93
	Negative (n = 14695)	2.08	2.04	1.91	2.07	2.15	2.48	14.20	11.80	6.70	4.91
	Ratio (Positive/Negative)	3.01	3.60	1.93	2.49	2.05	1.34	0.49	1.59	2.36	2.02
	p-value	2.6E-06	2E-09	0.04	4.9E-04	0.01	0.38	7E-04	4.6E-04	4E-09	1.7E-04
Dysuria	Positive (n = 272)	0.00	0.00	0.00	0.00	0.00	0.37	0.00	0.37	0.74	0.37
	Negative (n = 14695)	0.03	0.06	0.07	0.07	0.07	0.09	0.40	0.35	0.19	0.18
	Ratio (Positive/Negative)	0.00	0.00	0.00	0.00	0.00	4.16	0.00	1.04	3.86	2.08
	p-value	0.76	0.68	0.67	0.67	0.65	0.14	0.30	0.97	0.05	0.46

References

1. Guan, W.-J. *et al.* Clinical Characteristics of Coronavirus Disease 2019 in China. *N. Engl. J. Med.* (2020) doi:10.1056/NEJMoa2002032.
2. Verity, R. *et al.* Estimates of the severity of coronavirus disease 2019: a model-based analysis. *Lancet Infect. Dis.* (2020) doi:10.1016/S1473-3099(20)30243-7.
3. Hoehl, S. *et al.* Evidence of SARS-CoV-2 Infection in Returning Travelers from Wuhan, China. *N. Engl. J. Med.* **382**, 1278–1280 (2020).
4. Xiao, F. *et al.* Evidence for Gastrointestinal Infection of SARS-CoV-2. *Gastroenterology* (2020) doi:10.1053/j.gastro.2020.02.055.
5. Zhang, B. *et al.* Clinical characteristics of 82 death cases with COVID-19. *medRxiv* 2020.02.26.20028191 (2020).
6. COVID - Overview: Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2) RNA Detection, Varies. <https://www.mayocliniclabs.com/test-catalog/Overview/608825>.
7. Xu, Y. *et al.* Characteristics of pediatric SARS-CoV-2 infection and potential evidence for persistent fecal viral shedding. *Nat. Med.* **26**, 502–505 (2020).
8. Wu, F. *et al.* SARS-CoV-2 titers in wastewater are higher than expected from clinically confirmed cases. *medRxiv* 2020.04.05.20051540 (2020).
9. Venkatakrishnan, A. J. *et al.* Knowledge synthesis from 100 million biomedical documents augments the deep expression profiling of coronavirus receptors. *bioRxiv* 2020.03.24.005702 (2020) doi:10.1101/2020.03.24.005702.
10. Hoffmann, M. *et al.* SARS-CoV-2 Cell Entry Depends on ACE2 and TMPRSS2 and Is Blocked by a Clinically Proven Protease Inhibitor. *Cell* (2020) doi:10.1016/j.cell.2020.02.052.
11. del Pozo, F. P. *et al.* Inhibition of SARS-CoV-2 infections in engineered human tissues using clinical-grade soluble human ACE2. *Cell* (2020) doi:10.1016/j.cell.2020.04.004.
12. Grein, J. *et al.* Compassionate Use of Remdesivir for Patients with Severe Covid-19. *N. Engl. J. Med.* (2020) doi:10.1056/NEJMoa2007016.
13. Chen, C. *et al.* Favipiravir versus Arbidol for COVID-19: A Randomized Clinical Trial. *medRxiv* 2020.03.17.20037432 (2020).
14. Chen, H. *et al.* First Clinical Study Using HCV Protease Inhibitor Danoprevir to Treat Naive and Experienced COVID-19 Patients. *medRxiv* 2020.03.22.20034041 (2020).
15. Cao, B. *et al.* A Trial of Lopinavir-Ritonavir in Adults Hospitalized with Severe Covid-19. *N. Engl. J. Med.* (2020) doi:10.1056/NEJMoa2001282.
16. De Meyer, S. *et al.* Lack of Antiviral Activity of Darunavir against SARS-CoV-2. *medRxiv* 2020.04.03.20052548 (2020).
17. Shen, C. *et al.* Treatment of 5 Critically Ill Patients With COVID-19 With Convalescent Plasma. *JAMA* (2020) doi:10.1001/jama.2020.4783.
18. Bloch, E. M. *et al.* Deployment of convalescent plasma for the prevention and treatment of COVID-19. *J. Clin. Invest.* (2020) doi:10.1172/JCI138745.
19. Duan, K. *et al.* The feasibility of convalescent plasma therapy in severe COVID-19 patients: a pilot study. *medRxiv* 2020.03.16.20036145 (2020).
20. Convalescent Plasma COVID-19 (Coronavirus) Treatment – Mayo Clinic. <https://www.uscovidplasma.org/>.
21. Sterner, R. M. *et al.* GM-CSF inhibition reduces cytokine release syndrome and neuroinflammation but enhances CAR-T cell function in xenografts. *Blood* **133**, 697–709 (2019).
22. Shimabukuro-Vornhagen, A. *et al.* Cytokine release syndrome. *J Immunother Cancer* **6**, 56 (2018).
23. Herold, T. *et al.* Level of IL-6 predicts respiratory failure in hospitalized symptomatic COVID-19 patients. *medRxiv* 2020.04.01.20047381 (2020).
24. Mingfeng, X. X. H. *et al.* Effective Treatment of Severe COVID-19 Patients with Tocilizumab. *ChinaXiv.org* <http://www.chinaxiv.org/abs/202003.00026>.
25. Kenderian, S. S. *et al.* Ruxolitinib Prevents Cytokine Release Syndrome after CART Cell Therapy without Impairing the Anti-Tumor Effect in a Xenograft Model. *Blood* **128**, 652–652 (2016).
26. Gao, T. *et al.* Highly pathogenic coronavirus N protein aggravates lung injury by MASP-2-mediated complement over-activation. *medRxiv* 2020.03.29.20041962 (2020).
27. Gattinoni, L. *et al.* Covid-19 Does Not Lead to a 'Typical' Acute Respiratory Distress Syndrome. *Am. J. Respir.*

Crit. Care Med. (2020) doi:10.1164/rccm.202003-0817LE.

28. Vincent, M. J. *et al.* Chloroquine is a potent inhibitor of SARS coronavirus infection and spread. *Virology* **2**, 69 (2005).
29. Borba M *et al.* Chloroquine diphosphate in two different dosages as adjunctive therapy of hospitalized patients with severe respiratory syndrome in the context of coronavirus (SARS-CoV-2) infection: Preliminary safety results of a randomized, double-blinded, phase IIb clinical trial (CloroCovid-19 Study) (2020).
30. Wang, Y. *et al.* Early, low-dose and short-term application of corticosteroid treatment in patients with severe COVID-19 pneumonia: single-center experience from Wuhan, China. *medRxiv* 2020.03.06.20032342 (2020).
31. Russell, C. D., Millar, J. E. & Baillie, J. K. Clinical evidence does not support corticosteroid treatment for 2019-nCoV lung injury. *Lancet* **395**, 473–475 (2020).
32. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. (2018).

SUPPLEMENTARY MATERIAL

Figure S1. Cell-types connected to pathophysiology of COVID-19 as inferred from high expression of ACE2 and TMPRSS2 in human scRNA seq datasets. A scatter plot depicting the expression of ACE2 and TMPRSS2 inferred from the single-cell RNA-seq profiling of human tissues using nferX single cell platform. The x-axis represents the mean $\ln(\text{cp10k}+1)$ expression of ACE2 in all the cells and the y-axis represents the mean $\ln(\text{cp10k}+1)$ expression of TMPRSS2 in the corresponding cell-types from respective tissues. The colors on the scatter plot depicts the tissue origins. The size of the points on the scatter plot represents the percentage of single cells in the cell-type that co-express ACE2 and TMPRSS2 (non-zero expression).

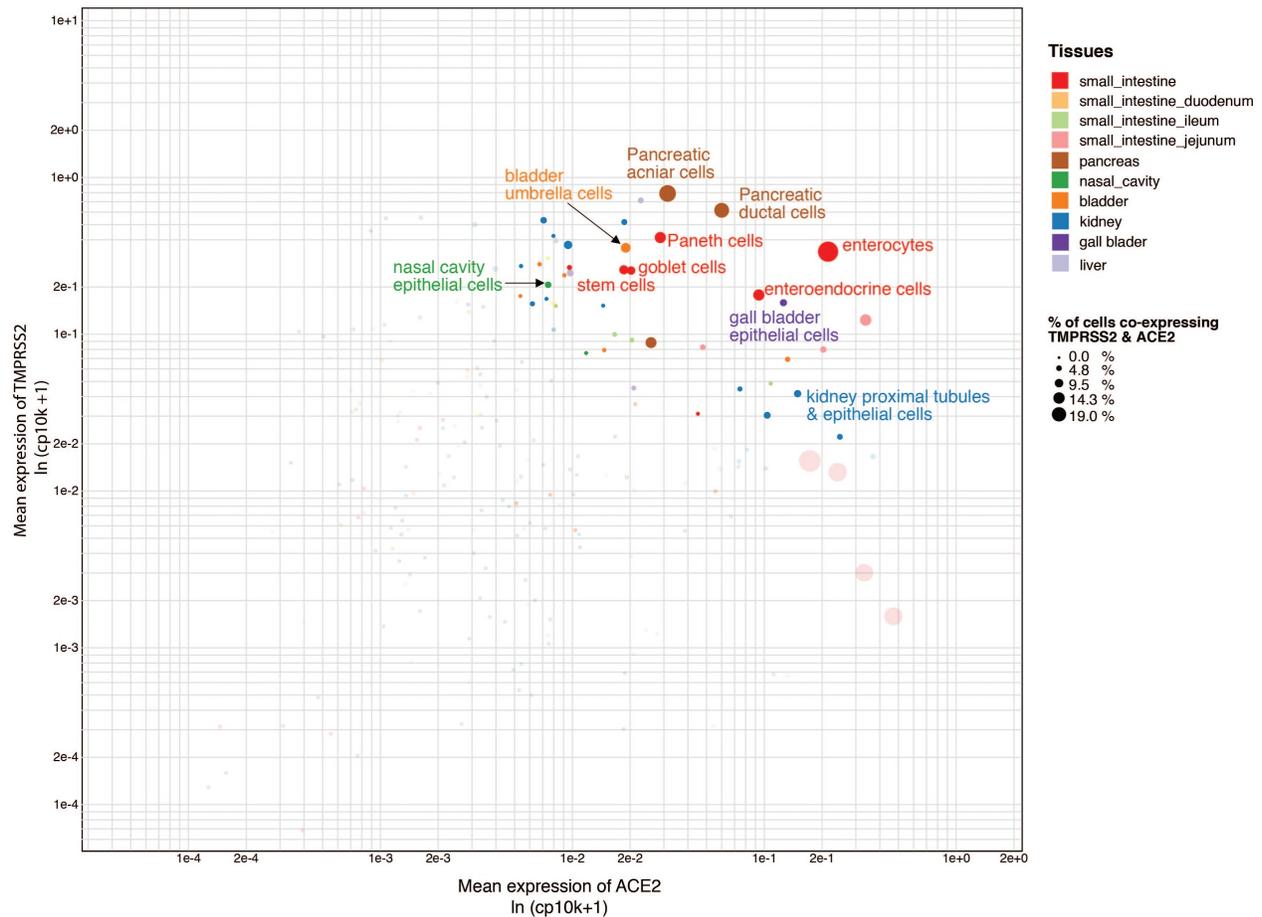


Figure S2. Disease progression of COVID-19 can be divided into multiple stages, and appropriate therapeutics can be chosen based on the specific pathophysiological mechanisms. Using nferX Knowledge Synthesis, the most associated molecular markers at each step of disease progression are also identified (see *Supplementary Methods* for details on nferX knowledge synthesis). In order to capture biomedical literature based associations, the nferX platform defines two scores: a “local score” and a “global score”, as described previously (Park, J. et al. Recapitulation and Retrospective Prediction of Biomedical Associations Using Temporally-enabled Word Embeddings. doi:10.1101/627513).

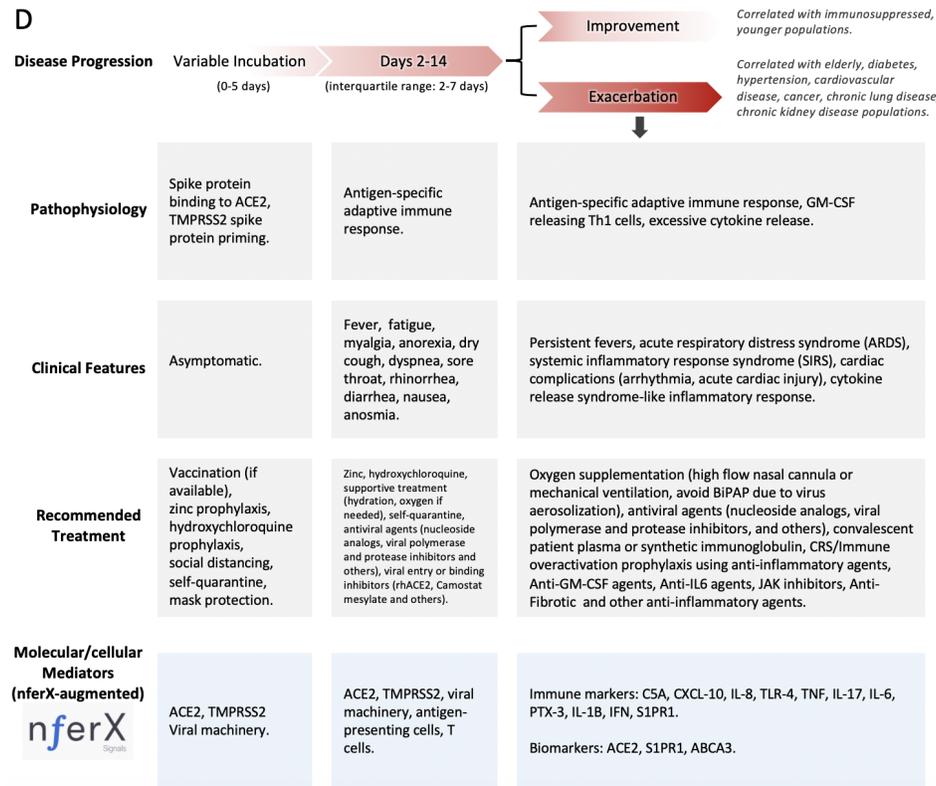
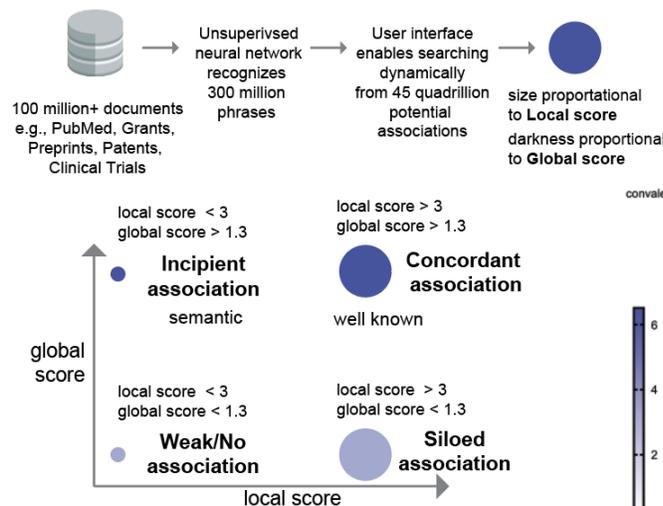


Figure S3. nferX-derived associations of COVID-19 treatment options to clinical phenotypes. (A) Schematic of the derivation of nferX local and global scores quantifying associations between concepts from across the literature. **(B)** Heatmap of nferX Local Scores capturing associations discussed in the literature between select COVID-19 treatment drugs and COVID-19 related phenotypes. In order to capture biomedical literature based associations, the nferX platform defines two scores: a “local score” and a “global score”, as described previously (Park, J. et al. Recapitulation and Retrospective Prediction of Biomedical Associations Using Temporally-enabled Word Embeddings. doi:10.1101/627513).

A Literature-synthesized association between two phrases
(e.g. 'drug name' and 'disease phenotype')



**B Association of nferX local scores with symptoms
and clinical phenotypes associated with COVID-19**

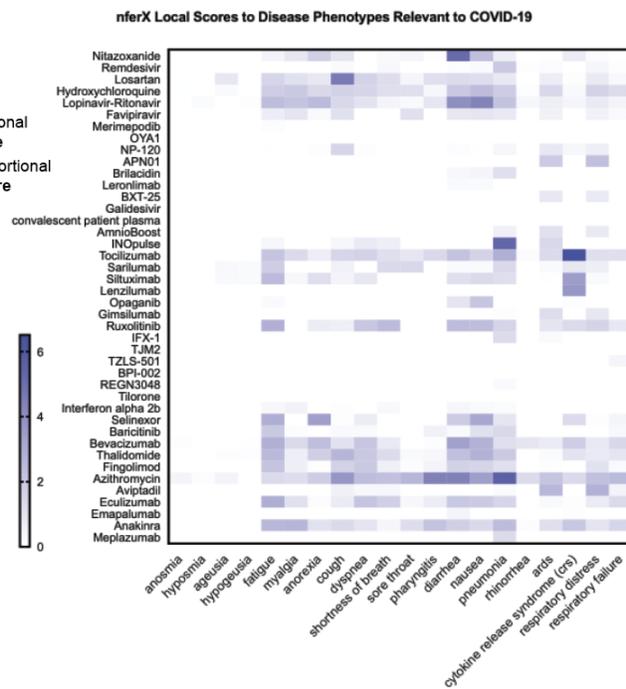


Table S1. Collection of treatments undergoing clinical trials for COVID-19.

Drug Class	Drug	Proposed Mechanism of Action	NCT Number
Renin-Angiotensin-Aldosterone System (RAAS) Inhibitors and recombinant ACE2	ACE inhibitors	Inhibits viral entry and downstream activation of angiotensin mediated inflammation. Effects of ACE inhibition under evaluation	NCT04330300, NCT04322786 NCT04330300
	Angiotensin Receptor Blockers	Inhibits angiotensin receptor. Given that SARS-CoV-2 uses ACE2 in viral entry, ARBs are under evaluation for potential efficacy and/or risk in COVID..	NCT04335123, NCT04328012,and more
	Recombinant human ACE2 (rhACE2)	SARS-CoV infection reduces ACE2 expression and that administration of recombinant ACE2 following infection can lead to reversal of lung injury	NCT04287686, NCT04335136
Antiviral Agents	Remdesivir	Nucleoside analog reported to improve respiratory status in 70% of patients (of 53 patients total) in an observation study. RCTs underway.	NCT04292899, NCT04292730 and more
	Favipiravir	Viral RNA polymerase inhibitor.RCT showed favipiravir improved the clinical recovery rate of COVID relative to umifenovir, a viral entry inhibitor	NCT04336904, NCT04333589 and more
	Lopinavir/ritonavir	Protease inhibitor used in HIV infection. Found to have a non-significant benefit in an RCT of 199 patients in China. Larger clinical studies in progress.	NCT04295551, NCT04331470 and more
	Umifenovir	Broad-spectrum antiviral that inhibits membrane fusion between virus and target hose. Initially developed for influenza.	NCT04260594, NCT04273763
	Ribavirin	Nucleoside analog with broad spectrum activity against RNA and DNA virus.	NCT04276688
	Emtricitabine/tenofovir	Nucleoside reverse-transcriptase inhibitor used for HIV infection being evaluated for COVID prevention in healthcare workers	NCT04334928
	Oseltamivir	sialidase used in influenza. Proposed for role of sialic acid in viral entry of coronavirus	NCT02735707, NCT04261270
	DAS181	inhaled sialidase, initially developed for parainfluenza infection, proposed given the role of sialic acid in the viral entry of coronavirus	NCT04324489, NCT04298060 and more
	Darunavir	Protease inhibitor used in HIV infection, proposed for potential antiviral efficacy.	NCT04252274
	ASC09	a novel protease inhibitor similar in structure to darunavir. Initially developed for HIV infection.	NCT04261907
	Meplazumab	Shown to bind to spike protein of SARS-CoV-2 and found to reduce severity of COVID in RCT of 17 patients in China ¹	NCT04275245
	Interferon Beta-1A	Shown to have in-vitro inhibitory activity against SARS-CoV2 ²	NCT04315948, NCT02735707 and more
	Interferon Beta-1B		
	Levamisole + Budesonide+ Formoterol inhaler	Found to bind to Papaine Like Protease, which is implicated in SARS-CoV-2 virulence	NCT04331470
Camostat Mesilate	Serine protease proposed for in-vitro inhibition of MERS-CoV entry ³	NCT04321096, NCT04338906	
	Convalescent plasma	Antibodies in recovered patients. Synthetic antibodies under development	NCT04264858
Anti-IL6 Agents	Tocilizumab	IL-6 inhibitor proposed to reduce risk of cytokine release syndrome in COVID	NCT04317092, NCT04320615, and more
	Siltuximab	IL-6 inhibitor proposed to reduce risk of cytokine release syndrome in COVID	NCT04330638, NCT04329650

	Sarilumab	IL-6R inhibitor proposed to reduce risk of cytokine release syndrome in COVID	NCT04315298, NCT04324073, and more
Complement Inhibitors	Eculizumab	proposed for anti-C5a activity to inhibit complement overactivation in COVID	NCT04288713
	IFX-1		NCT04333420
JAK Inhibitors	Baricitinib	JAK inhibitor proposed for anti-inflammatory activity to attenuate immune overactivation in severe COVID	NCT04340232, NCT04320277 and more
	Tofacitinib		NCT04332042
	Ruxolitinib		NCT04337359, NCT04331665 and more
Anti-GM-CSF Agents	Mavrilimumab	anti-GM-CSF agent proposed for anti-inflammatory effects in the event of COVID-induced cytokine release syndrome	NCT04337216
	Lenzilumab	anti-GM-CSF agent. Approved under emergency IND use.	
	Gimsilumab	anti-GM-CSF agent about to begin clinical trials.	
Other Anti-Inflammatory or Anti-Fibrotic Agents	Piclidenoson	initially developed for rheumatoid arthritis. Binds to A3AR, which may attenuate CRS. Also shown to have antiviral activity against RNA virus ⁴	NCT04333472
	Anakinra	IL-1a antagonist initially developed for rheumatoid arthritis.	NCT04330638, NCT04339712
	Emapalumab	Anti-IFN gamma agent proposed for anti-inflammatory effect in combination with anakinra	NCT04324021
	Thalidomide	anti-TNF agent used in multiple myeloma and graft-versus-host diseases proposed for immunomodulatory effects	NCT04273581
	Colchicine	Proposed for inhibition of inflammation caused by NLRP3 inflammasome, through inhibition of tubulin polymerization, and potential effects on cellular adhesion molecules and inflammatory chemokines	NCT04326790, NCT04328480 and more
	Fingolimod	Immunomodulatory agent used in multiple sclerosis.	NCT04280588
	Aviptadil	Synthetic VIP peptide proposed for immunomodulatory effects ⁵	NCT04311697
	BLD-2660	Selective CAPN inhibitor initially developed for pulmonary fibrosis. Studies found that CAPN inhibition led to reduce replication of SARS-CoV-1 ⁶	NCT04334460
	Bevacizumab	anti-VEGF antibody proposed to attenuate increases in vascular permeability in the COVID-induced vascular inflammation	NCT04305106, NCT04275414
	Bromhexine	Increases secretion of mucus components in the respiratory tract and alleviates respiratory inflammation. ⁷	NCT04340349, NCT04273763
	CD24Fc	Initially developed for GVHD Shown to reduce inflammation by binding to DAMP and Siglec G/10 to modulate immune response ⁸	NCT04317040
	Carrimycin	Under evaluation in China for efficacy against upper respiratory tract diseases	NCT04286503
	Nintedanib	tyrosine kinase inhibitor initially developed for idiopathic pulmonary fibrosis, likely approved for anti-inflammatory and anti-fibrotic activity	NCT04338802
	PUL-042	TLR 2/6/9 agonist proposed to prevent COVID	NCT04313023, NCT04312997
Defibrotide	Oligonucleotide initially approved in hepatic veno-occlusive diseases. Likely proposed for effect in modulation of endothelial injury.	NCT04335201	

	Ibrutinib	BTK inhibitor used for B-cell malignancies, but also has activity against TEC family kinase. In-vivo evidence for lung protection in viral infection ⁹	
Vasodilators or other agents targeting effects of hypoxic vasoconstriction	Sildenafil	Vasodilator likely proposed to attenuate perfusion dysregulation and hypoxic vasoconstriction in COVID ARDS	NCT04304313
	Inhaled nitric oxide		NCT04306393
	Sargramostim	inhaled sargramostim, which induces hematopoiesis, proposed to attenuates acute hypoxic respiratory failure in COVID	NCT04326920
Vaccines	mRNA-1273	Clinical trial sponsored by NIAID underway at Kaiser Permanente Washington Health Research Institute (KPWHRI) in Seattle.	NCT04283461
	INO-4800	Inovio announced an accelerated timeline for the development of the vaccine on 03 March. Trial underway.	NCT04336410
	BCG Vaccine	Suggested to be protective against COVID with limited data ¹⁰	NCT04328441, NCT04327206
Mesenchymal Stem Cell (MSC) Therapy	NestCell ^Å ®, Human Umbilical Derived MSC I, Dental Pulp MSC, Wharton's Jelly-MSCs	MSCs shown to reduce nonproductive inflammation and affect tissue regeneration and is being evaluated in patients with ARDS	NCT04315987 NCT04336254 NCT04313322 And more
	MSC + Ruxolitinib	MSC+ ruxolitinib being evaluated in severe COVID-19	ChiCTR2000029580
	Exosomes allogeneic adipose MSC (MSCs-Exo)	Pilot pilot clinical trial performed to explore the safety and efficiency of aerosol inhalation of the MSCs-Exo in patients with severe COVID-19	NCT04276987

Table S1 - References

1. Bian, Huijie, et al. "Meplazumab treats COVID-19 pneumonia: an open-labelled, concurrent controlled add-on clinical trial." medRxiv (2020).
2. Hensley, L., Fritz, E., Jahrling, P., Karp, C., Huggins, J. and Geisbert, T., 2020. Interferon-B 1A And SARS Coronavirus Replication.
3. Yamamoto, Mizuki, et al. "Identification of nafamostat as a potent inhibitor of Middle East respiratory syndrome coronavirus S protein-mediated membrane fusion using the split-protein-based cell-cell fusion assay." Antimicrobial agents and chemotherapy 60.11 (2016): 6532-6539.
4. Cohen, Shira, and Pnina Fishman. "Targeting the A3 adenosine receptor to treat cytokine release syndrome in cancer immunotherapy." Drug design, development and therapy 13 (2019): 491.
5. Chorny, Alejo, et al. "Vasoactive intestinal peptide induces regulatory dendritic cells with therapeutic effects on autoimmune disorders." Proceedings of the National Academy of Sciences 102.38 (2005): 13562-13567.
6. Barnard, Dale L., et al. "Inhibition of severe acute respiratory syndrome-associated coronavirus (SARSCoV) by calpain inhibitors and β -D-N4-hydroxycytidine." Antiviral Chemistry and Chemotherapy 15.1 (2004): 15-22.
7. Zanasi, Alessandro, Massimiliano Mazzolini, and Ahmad Kantar. "A reappraisal of the mucoactive activity and clinical efficacy of bromhexine." Multidisciplinary respiratory medicine 12.1 (2017): 7.
8. Chen, Guo-Yun, et al. "CD24 and Siglec-10 selectively repress tissue damage-induced immune responses." Science 323.5922 (2009): 1722-1725.
9. Florence, Jon M., et al. "Inhibiting Bruton's tyrosine kinase rescues mice from lethal influenza-induced acute lung injury." American Journal of Physiology-Lung Cellular and Molecular Physiology 315.1 (2018): L52-L58.
10. Dayal, Devi, and Saniya Gupta. "Connecting BCG Vaccination and COVID-19: Additional Data." medRxiv (2020).

Table S2. Disease Progression of COVID-19 and Associated Treatment

Diseases Progression	Symptoms and other Clinical Indicators	Recommended Treatment
Prior to or at time of Exposure	Asymptomatic	<ul style="list-style-type: none"> • Vaccination, when available • Zinc Prophylaxis • Hydroxychloroquine Prophylaxis • Social Distancing • Self-Quarantine • Mask Protection
Early Infection, without hospitalization	Fever Anosmia/Ageusia Mild respiratory symptoms (cough, dyspnea, sore throat, rhinorrhea) GI symptoms (diarrhea, nausea, abdominal discomfort) Focal consolidation on CXR; Focal ground glass opacity on CT	<ul style="list-style-type: none"> • Zinc • Hydroxychloroquine • Supportive Treatment (Hydration, Oxygen if needed) • Self-Quarantine • Antiviral agents (nucleoside analogs, viral polymerase and protease inhibitors, and others) • Viral entry or binding inhibitors (rhACE2, Camostat mesylate and others):
Moderate diseases, requiring supplemental oxygen	Persistent Fever Moderate respiratory and GI symptoms Dyspnea with increased oxygen requirement Oxygen desaturation (<93%) Diffuse Ground Glass Opacity on CT Chest	<ul style="list-style-type: none"> • Oxygen supplementation (high flow nasal cannula, avoid BiPAP due to virus aerosolization) • Antiviral agents (nucleoside analogs, viral polymerase and protease inhibitors, and others) • Convalescent patient plasma or Synthetic immunoglobulin • CRS/Immune overactivation prophylaxis using anti-inflammatory agents • Anti-GM-CSF Agents • Anti-IL6 Agents • JAK Inhibitors
Severe diseases requiring mechanical ventilation, with immune overactivation	ARDS Cytokine Release Syndrome-like inflammatory response Multi-organ involvement Elevation of IL-6, Ferritin & other inflammatory markers Diffuse bilateral coalescent opacities	Anti-Inflammatory agents, including <ul style="list-style-type: none"> • Anti-IL6 Agents • Complement Inhibitors • Anti-GM-CSF Agents • JAK Inhibitors • Anti-Fibrotic Agents

Table S3. Symptoms and their synonyms used for the EHR analysis.

Symptom/Finding	Synonyms/related entities identified in EHR
Fever / chills	fever, chills, fevers, temp, tactile fever, felt warm, subjective fever
Altered or diminished sense of taste or smell	Anosmia/Dysgeusia, decreased taste, altered sense of taste and smell, change in his sense of taste and smell, no sense of smell or taste, anosmia, change in smell and taste, dysgeusia, lost her sense of smell, ageusia, everything smells and tastes terrible, bitter taste in her mouth, change in taste, change in smell, decrease in smell, decrease in taste and smell, loss of smell and taste, altered smell, decreased sense of taste
Diarrhea	diarrhea, loose stools, diarrhea/vomiting, loose stool, watery BM, watery diarrhea, soft stools
GI upset	nausea, posttussive emesis, emesis, vomiting, diarrhea/vomiting, Nausea/vomiting/abdominal pain, stomach cramping
Wheezing	wheezing
Respiratory difficulty	lower respiratory symptoms, increased oxygen demands, labored breathing, shortness of breath, SOB, Dyspnea, dyspnea on exertion, tachypnea, tachypneic
Cough	cough, dry cough, nonproductive cough, +cough, non-productive cough, cough(NP/P)
Hemoptysis	hemoptysis, blood-tinged sputum
Productive cough	cough productive, productive cough, cough(NP/P)
Chest pain/pressure	chest congestion, chest heaviness, chest pain, pleuritic chest pain, chest tightness, chest discomfort
Congestion	congestion/rhinorrhea, sinus congestion, sinus pressure, congestion, nasal congestion, head congestion, stuffy nose
Rhinitis	tickling in nose, sneezing, Rhinitis, sniffles, congestion/rhinorrhea, runny nose, rhinorrhea
Myalgia/Arthralgia	aches and pains, body aches, body ache, myalgias, sore neck, myalgia, muscle aches, muscle discomfort, joints became sore, arthralgia, arthralgias, ache, back pain, achy joints
Generalized symptoms	cold, generalized weakness, malaise, weakness, weak, influenza like symptoms, feeling run down, activity change
Fatigue	fatigue, lethargy, energy level is poor, fatigued, sleeping more than usual
Diaphoresis	sweaty, sweats, diaphoresis, night sweats, sweating
Pharyngitis	sore throat, tingly throat, scratchy throat, throat discomfort
Headache	headache, HA, headaches, HA's, sinus headache
Dry mouth	dry mouth
Change in appetite/intake	decreased appetite, anorexia, appetite is poor, weight loss, not eating and drinking, no appetite, appetite change, Little appetite
Conjunctivitis	watery eyes, watery eyes with redness, pink eye, red eye, red eyes
Neuro	confusion, agitation, delirium, dizziness
Cardiac	palpitations, Lightheadedness
Otitis	earache
Dermatitis	rash
Angioedema	face was red, swollen and itchy

Table S4. Pairwise analysis of phenotypes in the COVID_{pos} and COVID_{neg} cohorts.

Phenotype 1	Phenotype 2	COVID _{pos} count (N=272)	COVID _{neg} count (N = 14695)	COVID _{pos} proportion (N=272)	COVID _{neg} proportion (N=14695)	(COVID _{pos} / COVID _{neg}) relative ratio	raw 2-tailed p-value	BH-corrected p-value
Altered or diminished sense of taste or smell	Headache	5	3	0.02	0.0002	90.04	8.32E-38	2.92E-35
Altered or diminished sense of taste or smell	Fever / chills	8	12	0.03	0.0008	36.02	1.82E-37	3.19E-35
Altered or diminished sense of taste or smell	Cough	8	13	0.03	0.0009	33.25	1.32E-35	1.55E-33
Altered or diminished sense of taste or smell	Congestion	5	7	0.02	0.0005	38.59	4.73E-25	4.15E-23
Cough	Diarrhea	36	486	0.13	0.03	4.00	9.31E-19	6.53E-17
Altered or diminished sense of taste or smell	Fatigue	5	11	0.02	0.0007	24.56	1.16E-18	6.80E-17
Diaphoresis	Diarrhea	21	204	0.08	0.01	5.56	1.83E-17	9.17E-16
Altered or diminished sense of taste or smell	Myalgia/Arthralgia	5	13	0.02	0.0009	20.78	1.58E-16	6.93E-15
Altered or diminished sense of taste or smell	Chest pain/pressure	4	8	0.01	0.0005	27.01	2.92E-16	1.14E-14
Diarrhea	Headache	23	254	0.08	0.02	4.89	3.43E-16	1.20E-14
Altered or diminished sense of taste or smell	Wheezing	3	4	0.01	0.0003	40.52	4.27E-16	1.36E-14
Altered or diminished sense of taste or smell	GI upset	4	9	0.01	0.0006	24.01	5.36E-15	1.57E-13
Diarrhea	Fever / chills	36	608	0.13	0.04	3.20	2.36E-13	6.36E-12
Altered or diminished sense of taste or smell	Diarrhea	4	12	0.01	0.0008	18.01	3.76E-12	9.44E-11
Altered or diminished sense of taste or smell	Respiratory difficulty	4	14	0.01	0.001	15.44	8.88E-11	2.08E-09
Altered or diminished sense of taste or smell	Pharyngitis	2	3	0.01	0.0002	36.02	1.63E-10	3.57E-09
Altered or diminished sense of taste or smell	Change in appetite/intake	3	8	0.01	0.0005	20.26	2.57E-10	5.31E-09
Chest pain/pressure	Diarrhea	21	324	0.08	0.02	3.50	1.89E-09	3.69E-08
Diarrhea	Generalized symptoms	28	513	0.10	0.03	2.95	2.58E-09	4.77E-08
Diarrhea	Fatigue	22	362	0.08	0.02	3.28	6.11E-09	1.07E-07
Altered or diminished sense of taste or smell	Rhinitis	2	4	0.01	0.0003	27.01	7.45E-09	1.25E-07
Altered or diminished sense of taste or smell	Respiratory failure	2	4	0.01	0.0003	27.01	7.45E-09	1.19E-07
Diaphoresis	Headache	16	234	0.06	0.02	3.69	4.49E-08	6.86E-07
Chest pain/pressure	Headache	23	450	0.08	0.03	2.76	4.69E-07	6.87E-06
Change in appetite/intake	Cough	16	270	0.06	0.02	3.20	1.38E-06	1.93E-05
Diaphoresis	Generalized symptoms	25	533	0.09	0.04	2.53	1.59E-06	2.15E-05
Cough	Fatigue	35	872	0.13	0.06	2.17	2.04E-06	2.66E-05
Chest pain/pressure	Diaphoresis	22	460	0.08	0.03	2.58	4.45E-06	5.57E-05
Headache	Hemoptysis	2	7	0.01	0.0005	15.44	4.56E-06	5.52E-05
Conjunctivitis	Diaphoresis	2	7	0.01	0.0005	15.44	4.56E-06	5.34E-05