# Gene Expression Risk Scores for COVID-19 Illness Severity

Derick R Peterson[1], Andrea M Baran[1], Soumyaroop Bhattacharya[2], Angela R Branche[3], Daniel P Croft[4], Anthony M Corbett[1], Edward E Walsh[3,5], *Ann R Falsey[3,5] and *Thomas J Mariani[2]

[1]Department of Biostatistics and Computational Biology, [2]Division of Neonatology and Pediatric Molecular and Personalized Medicine Program, Department of Pediatrics, [3]Division of Infectious Diseases, Department of Medicine, [4]Division of Pulmonary and Critical Care, Department of Medicine, University of Rochester, and [5]Department of Medicine, Rochester General Hospital, Rochester, NY, USA

* Corresponding Author

Address for Correspondence:

Thomas J Mariani, PhD

Division of Neonatology and

Pediatric Molecular and Personalized Medicine Program

University of Rochester Medical Center

601 Elmwood Ave, Box 850

Rochester, NY 14642, USA.

Phone: 585-276-4616;

Fax: 585-276-2643;

E-mail: Tom_Mariani@urmc.rochester.edu.

Ann R. Falsey, M.D.

Infectious Disease Unit

1    Rochester General Hospital

2    1425 Portland Avenue

3    Rochester, NY, 14621 USA

4    FAX: 585-922-5168, Phone: 585-922-4339

5    Email: ann.falsey@rochesterregional.org

6

7    Running Title: Risk Scores for COVID-19 Severity

8

9    Total Word Count: 3477

10    Abstract word: 192

11    1 Table

12    4 Figures

13    3 Supplementary Tables

14    2 Supplementary Figures

15
16

17

1   **Abstract**

2   **Background**: The correlates of COVID-19 illness severity following infection with SARS-

3   Coronavirus 2 (SARS-CoV-2) are incompletely understood.

4   **Methods**: We assessed peripheral blood gene expression in 53 adults with confirmed SARS-

5   CoV-2-infection clinically adjudicated as having mild, moderate or severe disease. Supervised

6   principal components analysis was used to build a weighted gene expression risk score

7   (WGERS) to discriminate between severe and non-severe COVID.

8   **Results**: Gene expression patterns in participants with mild and moderate illness were similar,

9   but significantly different from severe illness. When comparing severe versus non-severe

10  illness, we identified >4000 genes differentially expressed (FDR<0.05). Biological pathways

11  increased in severe COVID-19 were associated with platelet activation and coagulation, and

12  those significantly decreased with T cell signaling and differentiation. A WGERS based on 18

13  genes distinguished severe illness in our training cohort (cross-validated ROC-AUC=0.98), and

14  need for intensive care in an independent cohort (ROC-AUC=0.85). Dichotomizing the WGERS

15  yielded 100% sensitivity and 85% specificity for classifying severe illness in our training cohort,

16  and 84% sensitivity and 74% specificity for defining the need for intensive care in the validation

17  cohort.

18  **Conclusion**: These data suggest that gene expression classifiers may provide clinical utility as

19  predictors of COVID-19 illness severity.

20

21

22  **Keywords:** COVID-19; Severity; RNA sequencing; ICU; Prediction; Molecular Markers;
23

**Introduction**

In December 2019 a novel coronavirus, SARS-CoV-2, was identified in China as a cause of severe pneumonia with explosive human-to human transmission [1]. Illness due to SARS-CoV-2 has been designated COVID-19, and on March 11, 2020, the World Health Organization officially declared SARS-CoV-2 a pandemic. To date there have been over 100 million infections and over 2 million deaths globally due to COVID-19. (Source: https://covid19.who.int/) Although most patients experience mild to moderate disease, 5-10% progress to severe or critical illness with severe pneumonia or respiratory failure [2, 3]. Early in the pandemic it became clear that certain underlying chronic medical conditions, and principally age, were key risk factors for severe disease [4, 5]. While severe disease can occur early in illness, a distinct progression to severe illness occurs in some individuals 7-12 days after symptom onset suggesting transition from a viral phase to an inflammatory phase [6]. In addition, some young individuals without co-morbidities have also developed severe illness, highlighting the incomplete understanding of disease pathogenesis due to SARS-CoV-2 infection [7].

Gene expression provides an unbiased measure of the host response to a pathogen on a cellular level. We and others have previously demonstrated the potential for peripheral blood gene expression patterns to classify the ontogeny and severity of viral respiratory illness [8, 9]. We hypothesized that analysis of gene expression in the blood of patients with SARS-CoV2-related COVID-19 might help identify those at greatest risk for severe symptoms and in need of intensive care. Gene expression analysis might also identify pathways underlying disease pathogenesis and suggest new targets amenable to potential therapeutic interventions.

1    **Methods**

2    **Acute Illness Evaluation:** Adults ≥18 years of age, either hospitalized or recruited from the

3    community with symptoms compatible with COVID-19 and documented to have SAR-CoV-2 by

4    PCR, were eligible for the study. Participants with immunosuppression or symptoms onset

5    greater than 28 days prior to admission were excluded. Hospitalized participants were assessed

6    within 24 hours of admission and outpatients were brought to the clinic within 1-2 days of being

7    identified as SARS-CoV-2 positive. Demographic, clinical, radiographic and laboratory

8    information, date of symptom onset and signs and symptoms of the illness were collected.

9    Medication use was recorded with attention to drugs (steroids, leukotriene agonists,

10   hydroxychloroquine, Kaletra, Remdesivir or biologics) that may affect transcriptional profiling.

11

12   **Clinical severity assessment**: Severity for COVID-19 participants at enrollment and

13   throughout the illness was assessed using a combination of clinical variables (symptoms,

14   physical findings, radiographic and laboratory values) as well as the National Early Warning

15   Score (NEWS) of 7 graded physiological measurements (respiratory rate; oxygen saturation;

16   oxygen supplementation; temperature; blood pressure; heart rate; level of consciousness) [10].

17   Severe illness was defined as requiring any of the following: ICU care, high flow oxygen,

18   ventilator support, presser support or evidence of new end organ failure. Non-severe illness was

19   defined as illnesses not meeting severe criteria. In addition, a panel of 4 physicians (3 infectious

20   disease and 1 pulmonary critical care) adjudicated individually, then in a live panel discussion,

21   all non-severe illnesses and categorized them as mild or moderate using the NEWS as well as

22   symptoms and physiologic parameters in the context of underlying diseases and baseline

23   oxygen requirements. Participants were followed for the duration of hospitalization and illness,

24   and outcomes were recorded as the highest level of care required or death.

25

1    **Sample Collection and Processing**: Approximately 3 ml of whole blood was collected in a

2    Tempus™ Blood RNA Tube at the time of enrollment and stored at -80C until the time of

3    processing. Following centrifugation, RNA was isolated from the pellet using the Tempus Spin

4    RNA Isolation Kit using the manufacturer recommended protocol. Total RNA was processed for

5    globin reduction using GLOBINclear Human Kit as described previously [9].

6

7    **RNA Sequencing**: Globin-reduced RNA was used for transcriptomic profiling by RNA-seq.

8    cDNA libraries were generated using 200 ng of globin-reduced total RNA. Library construction

9    was performed using the TruSeq Stranded mRNA library kit (Illumina, San Diego, CA). cDNA

10   quantity was determined with the Qubit Flourometer (Life Technologies, Grand Island, NY) and

11   quality was assessed using the Agilent Bioanalyzer 2100 (Agilent, Santa Clara, CA). Libraries

12   were sequenced on the Illumina NovaSeq6000 at a target read depth of ~20 million $1 \times 100$-

13   bp single end reads per sample. Sequences were aligned against the human genome version

14   hg38 using the Splice Transcript Alignment to a Reference (STAR) algorithm [11], and counts

15   were generated using HTSeq [12]. Raw counts were divided by participant-specific library size

16   (in millions) to yield counts per million (CPM)-normalized expression, borrowing no information

17   across participants, and gene and sample level filtering was performed to remove outlier

18   samples and low expressing genes. Normalized and filtered analytical data sets were $\log_2$-

19   transformed (after adding a pseudo-count of 1 CPM) prior to analysis. We excluded data from

20   19,861 genes with uniformly zero reads, leaving a data set comprised of 39,225 genes from 53

21   participants. Finally, we retained genes that had normalized counts exceeding 1 CPM in greater

22   than 14 participants (the smallest class size). This resulted in an analytical dataset of 14,228

23   CPM-normalized genes.

24

**Statistical Methods:** Continuous clinical variables were compared by COVID severity levels using the nonparametric Kruskal-Wallis test, and binary variables by Fisher's exact test. We tested for differential expression by COVID-19 severity using the nonparametric Wilcoxon rank sum test. To allow adjustment for important clinical covariates, we fit semi-parametric Cox proportional hazards models for normalized gene expression as a function of severe vs non-severe COVID-19, adjusted for race, sex, BMI, days since symptom onset, and library size. The Benjamini-Hochberg procedure was used to control the False Discovery Rate (FDR). Pathway analysis of significantly differentially expressed genes was performed using ENRICHR [13].

**Classifier Development and Testing:** We used a version of supervised principal components analysis to build a weighted gene expression risk score (WGERS) to discriminate between severe and non-severe COVID-19. Genes were selected based on their univariate AUC and fold-change, as estimated by Hodges-Lehmann median of all pairwise shifts in log expression, with thresholds selected within the inner loop of nested 20-fold cross-validation. The cross-validation sampling strategy was designed to efficiently approximate leaving out 1 subject from each of the 3 sub-levels of severity. The selected genes were standardized to Z-scores with mean 0 and unit variance, and their first PC score was used as the sole predictor for logistic regression. The outer loop of nested cross-validation was used to estimate the ROC and AUC of the adaptive procedure. The nested pooled AUC corresponds with the ROC curve, and compares samples across and within models. The nested stratified AUC only compares samples from the same model, and thus generally is preferable, but it does not correspond with any single ROC curve. A subsequent run of non-nested cross-validation produced the thresholds used to define the gene set for the final risk score. The WGERS is calculated as the linear combination of the standardized, $\log_2$-transformed genes that meet the chosen thresholds, with coefficients based on the first principal component loading of the genes, scaled by the coefficient from the univariate logistic regression model.

1   To perform an independent validation of our risk score, we use a dataset from Overmyer et al.

2   [14]. There were some notable differences between the validation data and our training dataset,

3   including a different definition of severity in the outcome (ICU vs non-ICU) and use of a different

4   normalization for the gene expression data (TPM). Of the 18 genes used in our risk score, 2

5   were missing in the validation data. We imputed data for these 2 genes via multiple linear

6   regression with coefficients estimated by regressing each on the 16 non-missing CPM-

7   normalized log gene expression values in the training data. We standardized the TPM-

8   normalized validation gene expression data using means and SDs estimated from the training

9   data, and then applied the risk score coefficients from the training data to construct a risk score

10  for each validation subject. Apparent miscalibration required choosing a different WGERS

11  threshold for the validation data due to gene expression measures being generally lower in the

12  validation data compared to the training data. An ROC curve with associated AUC was used to

13  assess the performance of the risk score in the validation data.

14

15

16

1   **Results**

2   Between April 30th and June 29th 2020, 58 participants with PCR documented COVID-19

3   illnesses were enrolled from inpatient and outpatient settings. Of these, 3 participants did not

4   have blood samples collected and 2 did not meet inclusion criteria, leaving 53 participants for

5   RNA sequencing analysis. Illnesses were adjudicated as 20 severe and 33 non-severe (14 mild

6   and 19 moderate). This categorization was consistent with the severity separation in the NEWS

7   (Supplemental Figure S1). Two severely ill participants received one dose of Remdesivir prior to

8   blood collection. No subject received steroids or any other experimental COVID-19 treatment

9   prior to enrollment. Five hospitalized participants had rapidly progressive hypoxemia and

10  hemodynamic instability after enrollment and required transfer to intensive care, and 3

11  subsequently were mechanically ventilated. No mildly ill outpatient illnesses progressed in

12  severity to require medical attention. There was insufficient evidence of any difference in

13  demographic characteristics or underlying conditions by disease severity, except for race and

14  time from disease onset (Table 1): white non-Hispanic comprised 93% mild vs 50-58%

15  moderate-severe (p=0.02), and median time from symptom onset to enrollment was 4 days

16  among mild, 9 days among moderate, and 6.5 days among severe (p=0.047, due to

17  heterogeneity among non-severe). The median age of participants was 62 years with 53% of

18  them being male. As expected, dyspnea, hypoxemia, the presence of infiltrates, and use of

19  supplemental oxygen were more common in moderate and severe, compared to mild illness. All

20  severely ill patients required intensive care; 15 were enrolled in the ICU and 5 were moved to

21  ICU within 48 hours of blood sampling. All severely ill participants required supplemental

22  oxygen; 12 (60%) were mechanically ventilated, one was supported with ECMO and survived,

23  13 (65%) required vasopressor support and one subject died. Median NEWS were different

24  between the 3 groups (Figure S1). Inflammatory markers were not available for most outpatients

25  but were notably elevated in those hospitalized with moderate to severe disease. (Table 1)

26

1     Blood gene expression profiling from SARS-CoV-2 positive cases (n=53) was completed by

2     standard mRNA sequencing (RNAseq) of globin mRNA-reduced RNA isolated from whole blood

3     at the time of recruitment. On average 58 ± 6 million reads were generated from each of the

4     cDNA libraries, with a mapping rate of 94.2 ± 0.6% and transcriptome coverage of

5     41.3 ± 1.3% (Supplemental Figure S2). Exploratory Principal Components Analysis suggested

6     similar patterns of gene expression might be shared by participants with mild and moderate

7     illness, but appeared distinct from those with severe illness (Figure 1A). Statistical analysis for

8     differential gene expression confirmed significant differences when comparing mild vs severe,

9     and moderate vs severe, but not mild vs moderate COVID (Figure 1B).

10

11     We next tested for differences in gene expression when comparing participants with severe

12     (n=20) vs non-severe illness (n=33), pooling the 14 mild and 19 moderate cases. We tested for

13     differential gene expression without (univariate) and with adjustment for variables potentially

14     associated with severe outcome (race, sex, BMI), the number of days since onset of symptoms,

15     and library size (Figure 1C and Supplemental Table 1). These analyses identified 6483 (46% of

16     tested) and 8435 (59% of tested) differentially expressed genes, with and without multivariate

17     adjustment, respectively.

18

19     We performed ontology analysis for the 6483 genes identified as differentially expressed in

20     severe COVID illness, focusing on the fully adjusted analysis (Figure 2). This analysis identified

21     74 pathways over-represented by genes (n=936) significantly upregulated in severe COVID,

22     and 25 pathways over-represented by genes (n=5547) significantly downregulated (Figure 2

23     and Supplemental Table S2). Activated pathways included a number associated with infectious

24     diseases as well as TNFα and NFkB signaling. Notably, there was also evidence for significant

25     upregulation of genes associated with platelet activation and coagulation. Among pathways

26     associated with downregulated genes in severe COVID were multiple pathways involved in

1    general host RNA metabolism as well as multiple pathways specifically associated with T cell

2    regulation, including Th2 and Th17 differentiation. The most significantly downregulated

3    pathway was associated with HSV1 infection.

4

5    Given the substantial number of differentially expressed genes when comparing severe vs non-

6    severe COVID, we investigated the ability of gene expression patterns to discriminate severe

7    illness. Gene-specific thresholds for univariate AUC and magnitude change were chosen via the

8    cross-validation procedure and used to produce an 18 gene weighted gene expression risk

9    score (WGERS) for severe illness. Nested cross-validation was used to estimate performance

10   via the stratified AUC (CV-AUC=0.98). The pooled CV-AUC of 0.93 corresponds with a cross-

11   validated ROC curve to graphically summarize performance (Figure 3A). The pooled CV-ROC

12   curve also was used to select a risk score threshold (-1.04) with 95% sensitivity and 88%

13   specificity, which corresponded with apparent (non-cross-validated) sensitivity of 100%,

14   specificity of 85%, and error rate of 9% (5/53), represented via the WGERS distributions for the

15   training data (Figure 4A). All 5 misclassified participants had moderate illness (Figure 4B).

16

17   We next identified an independent validation data set describing peripheral blood-based gene

18   expression profiling of COVID subjects who were either admitted (n=50) or not admitted (n=50)

19   to the ICU due to the severity of their acute illness [14]. Our 18 gene WGERS discriminated

20   between ICU and non-ICU patients with an AUC of 0.85, and thresholding at 1.77 yielded 84%

21   sensitivity and 74% specificity (Figures 3B and 4C). Furthermore, all 18 genes selected in the

22   training data were differentially expressed (FDR < 0.01) in the validation data (Supplemental

23   Table S3).

24

25

**Discussion**

SARS-CoV-2 infection causes a wide spectrum of disease ranging from minimal, often asymptomatic, respiratory illness to severe pneumonia with multisystem failure and death. Although measurements of inflammatory markers such as C-reactive protein and serum IL-6 levels are often associated with worse disease, their use to predict poor outcomes is imperfect [15-17]. Viral characteristics, such as shedding kinetics or gene sequence variation, are not reliable predictors of clinical outcome [18, 19]. Genome-wide expression profiling, a powerful and unbiased tool, can be used for multiple purposes such as relating activation or suppression of molecular pathways to clinical manifestations of disease, identification of biomarkers that may allow individual prediction of disease severity, and identification of novel gene targets for therapeutic intervention. Early predictors to identify patients that will decompensate following SARS-CoV-2 infection would be highly impactful.

In our study of 53 SARS-CoV-2 infected adults with illness ranging from very mild upper respiratory infection to acute respiratory failure, we identified >6,000 differentially expressed genes (DEGs) (FDR < 0.05) between severe and non-severe illness. The vast majority (85%) of DEGs were under-expressed, most notably with a marked effect on lymphocytes and altered function [20, 21]. Pathway analysis revealed inhibition of Th1, Th2 and Th17 cell differentiation, as well as inhibition of the T cell receptor signaling pathway. These effects are likely related to the marked lymphopenia and poor adaptive immune response in persons with severe SARS-CoV-2 infection [22]. Also notable in severe illness is the inhibition of the mRNA surveillance pathways that include the nonsense-mediated mRNA decay pathway which can degrade viral mRNA. Using a model coronavirus, murine hepatitis virus, Wada and colleagues showed viral transcription is enhanced by blocking this host cell pathway, demonstrated to be mediated by the viral nucleocapsid protein [23].

1     Several activated pathways we identified in our studies are worth comment, given what is

2     already known about SARS-CoV-2 and COVID-19. Activation of the NF-kappa B and TNF

3     signaling pathways in a setting of heightened inflammatory process is not surprising. Activation

4     of the platelet, complement, and coagulation cascade pathways are also expected, given the

5     characteristic hypercoagulable state that has been observed in severe illness [24].

6     Thrombocytopenia and activated platelets are associated with the high incidence of venous and

7     arterial clotting, while elevated levels of serum D-dimer, a fibrinogen degradation product, and

8     increased INR are all features of severe COVID-19 [25]. It is interesting that the infection-related

9     pathways most significantly activated include those principally associated with intracellular

10     bacterial (legionella, mycobacterial) and parasitic (toxoplasma, leishmania and trypanosome

11     [Chagas]) infections. These infections are associated with marked activation of macrophages,

12     and thus may be consistent with activation of the osteoclast differentiation pathway, as

13     osteoclasts and macrophages have many similarities [26, 27].

14

15     Our findings are generally consistent with the limited data currently available in the literature on

16     COVID-19 and gene expression [14, 28-31]. Specifically, Overmyer et al reported gene

17     expression and metabolomic data from 128 COVID infected and non COVID infected persons,

18     where 219 molecular features with high significance to COVID-19 status and severity were

19     discovered. [14] A number of these involved complement activation, dysregulated lipid transport,

20     and neutrophil activation. Additionally, our data is supported by the findings noted by Ouyang et

21     al in which Th17 and T cell activation and differentiation were markedly downregulated in severe

22     disease [32].

23

24     There were some novel pathways that demonstrated upregulation in our studies. A number of

25     malignancy related pathways were upregulated (acute myelogenous leukemia (AML),

26     colorectal, pancreatic and Non-small cell cancer) in the severe COVID patients. Similarly, in a

1    study by Kwan et al comparing gene expression in 45 COVID-19 cases to healthy controls, 135

2    genes were found to be differentially expressed, with enrichment for several cancer pathways

3    including viral carcinogenesis and AML [29]. Identification of mutations in cell signaling,

4    proliferation genes, and kinases such as AP2-associated protein kinase 1 (AAK1)–have led to

5    targeted treatment options for cancer patients. Baricitinib, a repurposed rheumatoid arthritis

6    drug that interferes with the Janus Kinase (JAK) pathway, was demonstrated to have efficacy

7    when combined with Remdesivir in the treatment of severe COVID  [33]. Baricitinib shows high

8    affinity for AP2 associated protein kinase 1 binding, potentially demonstrating some overlap in

9    perturbations of cell signaling pathways in malignancy and COVID-19.  Further investigation of

10    gene expression pathways differentially expressed in severely ill patients may provide clues to

11    new therapeutic targets.

12

13    Although our study was not designed to identify and validate early predictors of severe disease,

14    the data do offer a first step. Using gene expression data we were able develop and validate an

15    18 gene signature for severe disease –fully concordant with requiring ICU– with 85% AUC, 84%

16    sensitivity, and 74% specificity in an independent validation data set. In a recent paper Guardela

17    et al assessed the utility of blood transcript levels of 50 genes known to predict mortality in

18    Idiopathic Pulmonary Fibrosis patients to classify illness severity in COVID-19 [31]. A discovery

19    cohort of eight subjects was used, and then validated using a publicly available data set of 128

20    subjects [14]. The gene expression risk profile discriminated ICU admission, need for

21    mechanical ventilation, and in-hospital mortality with an AUC of 77%, 75%, and 74%,

22    respectively (p < 0.001) in a COVID-19 validation cohort.

23

24    Our current study has several limitations which are worth noting, including its relatively small

25    sample size, the non-standardized interval between symptom onset and sample collection, and

26    blood collection at one time point. The complexity of the clinical data among hospitalized

1    participants (i.e. admissions only for isolation, persons with chronic oxygen requirements,

2    COVID testing for procedures) made objective criteria to distinguish mild from moderate disease

3    difficult, necessitating the need for clinical adjudication. Lastly, certain laboratory studies were

4    not available for all subjects.

5

6

1  **Conclusions**

2  In summary, we found a large number of differentially expressed genes in the peripheral blood

3  that distinguished those with severe COVID-19 illness from those with mild or moderate

4  disease. These data could be used to identify potential targets for interventions, as well as to

5  develop predictors of disease severity. Future prospective studies are needed to follow mild to

6  moderately ill patients over time and evaluate whether any of the discriminatory genes identified

7  are affected at early stages and can serve as predicators of severity. If so, individuals with high

8  risk gene profiles might be hospitalized for observation, moved to a more closely monitored

9  setting while hospitalized, or targeted for early interventions such as monoclonal antibody

10  treatments.

11

12
13

17

1 **Footnotes**

2 **Acknowledgements**

3 We acknowledge the University of Rochester Genomics Research Center for completing RNA

4 isolation and sequencing. Christopher Slaunwhite, Mary Anne Formica, and Michael Peasley

5 provided technical support. Jeanne Holden-Wiltse and Jeffrey Williams assisted with data

6 management.

7

8 **Author Contributions**

9 DRP, EEW, ARF and TJM conceptualized the study. DPC, ARB, EEW and ARF recruited

10 patients, adjudicated clinical severity and collected samples. DRP, AMB, SB, AMC, and TJM

11 analyzed the data. DRP, AMB, SB, ARB, DPC, EEW, TJM and ARF interpreted the results. All

12 authors contributed to writing of the manuscript. All authors reviewed, edited and approved the

13 manuscript for submission.

14 The manuscript represents original work that is not currently under consideration elsewhere.

15

16 **Data Availability:** Raw and processed data from the study are currently in process of being

17 submitted to DBGAP

18

19 **Conflict of Interest**

20 ARF received grants from Janssen, Merck, Sharpe and Dohme, Pfizer, BioFire Diagnostics and

21 personal fees for serving on DSMB for Novavax. EEW has consulted for Janssen

22 Pharmaceuticals, and have received research funding from Janssen, Gilead, Medimmune,

23 Sanofi Pasteur and ADMA biologics. ARB consults for GSK and have grants from Pfizer, Merck,

24 Janssen and Cyanvac. The other authors have no competing financial interests to report.

25

26 **Funding Support:**

1    Supported by NIEHS AI137364 (TJM and ARF) and K23 ES032459 (DC).

2

3

**References**

1. Zhu N, Zhang D, Wang W, et al. A Novel Coronavirus from Patients with Pneumonia in China, 2019. N Engl J Med **2020**; 382:727-33.

2. Guan WJ, Ni ZY, Hu Y, et al. Clinical Characteristics of Coronavirus Disease 2019 in China. N Engl J Med **2020**; 382:1708-20.

3. Chen N, Zhou M, Dong X, et al. Epidemiological and clinical characteristics of 99 cases of 2019 novel coronavirus pneumonia in Wuhan, China: a descriptive study. Lancet **2020**; 395:507-13.

4. Wu Z, McGoogan JM. Characteristics of and Important Lessons From the Coronavirus Disease 2019 (COVID-19) Outbreak in China: Summary of a Report of 72314 Cases From the Chinese Center for Disease Control and Prevention. JAMA **2020**; 323:1239-42.

5. Booth A, Reed AB, Ponzo S, et al. Population risk factors for severe disease and mortality in COVID-19: A global systematic review and meta-analysis. PLoS One **2021**; 16:e0247461.

6. Siddiqi HK, Mehra MR. COVID-19 illness in native and immunosuppressed states: A clinical-therapeutic staging proposal. J Heart Lung Transplant **2020**; 39:405-7.

7. Owusu D, Kim L, O'Halloran A, et al. Characteristics of Adults Aged 18-49 Years Without Underlying Conditions Hospitalized With Laboratory-Confirmed Coronavirus Disease 2019 in the United States: COVID-NET-March-August 2020. Clin Infect Dis **2021**; 72:e162-e6.

8. Wang L, Chu CY, McCall MN, et al. Airway gene-expression classifiers for respiratory syncytial virus (RSV) disease severity in infants. BMC Med Genomics **2021**; 14:57.

9. Bhattacharya S, Rosenberg AF, Peterson DR, et al. Transcriptomic Biomarkers to Discriminate Bacterial from Nonbacterial Infection in Adults Hospitalized with Respiratory Illness. Sci Rep **2017**; 7:6548.

10. Richardson D, Faisal M, Fiori M, Beatson K, Mohammed M. Use of the first National Early Warning Score recorded within 24 hours of admission to estimate the risk of in-hospital mortality in unplanned COVID-19 patients: a retrospective cohort study. BMJ Open **2021**; 11:e043721.

1  11. Dobin A, Davis CA, Schlesinger F, et al. STAR: ultrafast universal RNA-seq aligner.

2  Bioinformatics **2013**; 29:15-21.

3  12. Anders S, Pyl PT, Huber W. HTSeq--a Python framework to work with high-throughput

4  sequencing data. Bioinformatics **2015**; 31:166-9.

5  13. Chen EY, Tan CM, Kou Y, et al. Enrichr: interactive and collaborative HTML5 gene list

6  enrichment analysis tool. BMC Bioinformatics **2013**; 14:128.

7  14. Overmyer KA, Shishkova E, Miller IJ, et al. Large-Scale Multi-omic Analysis of COVID-19

8  Severity. Cell Syst **2021**; 12:23-40 e7.

9  15. Bennett TD, Moffitt RA, Hajagos JG, et al. Clinical Characterization and Prediction of Clinical

10  Severity of SARS-CoV-2 Infection Among US Adults Using Data From the US National COVID

11  Cohort Collaborative. JAMA Netw Open **2021**; 4:e2116901.

12  16. Bentivegna M, Hulme C, Ebell MH. Primary Care Relevant Risk Factors for Adverse

13  Outcomes in Patients With COVID-19 Infection: A Systematic Review. J Am Board Fam Med

14  **2021**; 34:S113-S26.

15  17. Kalil AC, Patterson TF, Mehta AK, et al. Baricitinib plus Remdesivir for Hospitalized Adults

16  with Covid-19. N Engl J Med **2021**; 384:795-807.

17  18. Abdulrahman A, Mallah SI, Alqahtani M. COVID-19 viral load not associated with disease

18  severity: findings from a retrospective cohort study. BMC Infect Dis **2021**; 21:688.

19  19. Liu Y, Yan LM, Wan L, et al. Viral dynamics in mild and severe cases of COVID-19. Lancet

20  Infect Dis **2020**; 20:656-7.

21  20. Vigon L, Rodriguez-Mora S, Luna A, et al. Cytotoxic cell populations developed during

22  treatment with tyrosine kinase inhibitors protect autologous CD4+ T cells from HIV-1 infection.

23  Biochem Pharmacol **2020**; 182:114203.

24  21. DiPiazza AT, Graham BS, Ruckwardt TJ. T cell immunity to SARS-CoV-2 following natural

25  infection and vaccination. Biochem Biophys Res Commun **2021**; 538:211-7.

1   22. Rydyznski Moderbacher C, Ramirez SI, Dan JM, et al. Antigen-Specific Adaptive Immunity

2   to SARS-CoV-2 in Acute COVID-19 and Associations with Age and Disease Severity. Cell **2020**;

3   183:996-1012 e19.

4   23. Wada M, Lokugamage KG, Nakagawa K, Narayanan K, Makino S. Interplay between

5   coronavirus, a cytoplasmic RNA virus, and nonsense-mediated mRNA decay pathway. Proc

6   Natl Acad Sci U S A **2018**; 115:E10157-E66.

7   24. Ksander GA. Collagen coatings reduce the incidence of capsule contracture around soft

8   silicone rubber implants in animals. Ann Plast Surg **1988**; 20:215-24.

9   25. Levy JH, Iba T, Olson LB, Corey KM, Ghadimi K, Connors JM. COVID-19: Thrombosis,

10  thromboinflammation, and anticoagulation considerations. Int J Lab Hematol **2021**; 43 Suppl

11  1:29-35.

12  26. Upadhyay S, Mittal E, Philips JA. Tuberculosis and the art of macrophage manipulation.

13  Pathogens and Disease **2018**; 76.

14  27. Bogdan C. Macrophages as host, effector and immunoregulatory cells in leishmaniasis:

15  Impact of tissue micro-environment and metabolism. Cytokine X **2020**; 2:100041.

16  28. Jain R, Ramaswamy S, Harilal D, et al. Host transcriptomic profiling of COVID-19 patients

17  with mild, moderate, and severe clinical outcomes. Comput Struct Biotechnol J **2021**; 19:153-

18  60.

19  29. Kwan PKW, Cross GB, Naftalin CM, et al. A blood RNA transcriptome signature for COVID-

20  19. BMC Med Genomics **2021**; 14:155.

21  30. Li CX, Chen J, Lv SK, Li JH, Li LL, Hu X. Whole-Transcriptome RNA Sequencing Reveals

22  Significant Differentially Expressed mRNAs, miRNAs, and lncRNAs and Related Regulating

23  Biological Pathways in the Peripheral Blood of COVID-19 Patients. Mediators Inflamm **2021**;

24  2021:6635925.

25  31. Juan Guardela BM, Sun J, Zhang T, et al. 50-gene risk profiles in peripheral blood predict

26  COVID-19 outcomes: A retrospective, multicenter cohort study. EBioMedicine **2021**; 69:103439.

1   32. Ouyang Y, Yin J, Wang W, et al. Downregulated Gene Expression Spectrum and Immune

2   Responses Changed During the Disease Progression in Patients With COVID-19. Clin Infect

3   Dis **2020**; 71:2052-60.

4   33. Richardson P, Griffin I, Tucker C, et al. Baricitinib as potential treatment for 2019-nCoV

5   acute respiratory disease. Lancet **2020**; 395:e30-e1.

6

1    **Table 1 Clinical Variables**

|  | Mild (N=14) | Moderate (N=19) | Severe (N=20) | p-value |
|---|---|---|---|---|
| **Demographics** |  |  |  |  |
| Age, median (IQR) | 63.0 (41.0) | 59.0 (29.0) | 63.5 (19.5) | 0.71 |
| Male Sex, No. (%) | 5 (35.7) | 11 (57.9) | 12 (60.0) | 0.37 |
| White Non-Hispanic Race, No. (%) | 13 (92.9) | 11 (57.9) | 10 (50.0) | 0.02 |
| BMI, median (IQR) | 28.7 (9.2) | 30.4 (11.7) | 26.4 (9.5) | 0.29 |
| Days from Symptom Onset, median (IQR) | 4.0 (5.0) | 9.0 (6.0) | 6.5 (4.5) | 0.05 |
| **Underlying Conditions** |  |  |  |  |
| COPD, No. (%) | 3 (21.4) | 6 (31.6) | 2 (10.0) | 0.25 |
| CHF, No. (%) | 4 (28.6) | 1 (5.3) | 2 (10.0) | 0.17 |
| Diabetes, No. (%) | 2 (14.3) | 8 (42.1) | 5 (25.0) | 0.20 |
| Hypertension, No. (%) | 7 (50.0) | 11 (57.9) | 12 (60.0) | 0.88 |
| Asthma, No. (%) | 0 (0.0) | 1 (5.3) | 2 (10.0) | 0.77 |
| **Symptoms** |  |  |  |  |
| Cough, No. (%) | 6 (42.9) | 15 (78.9) | 11 (55.0) | 0.10 |
| Fever, No. (%) | 6 (42.9) | 15 (78.9) | 15 (75.0) | 0.08 |
| Dyspnea, No. (%) | 4 (28.6) | 12 (63.2) | 13 (65.0) | 0.08 |
| Rigors, No. (%) | 0 (0.0) | 3 (15.8) | 0 (0.0) | 0.06 |
| **Physical Findings** |  |  |  |  |
| Systolic Blood Pressure, median (IQR) | 127 (27) (N=13) | 119 (34) | 105 (33) | 0.07 |

| | | | | |
|---|---|---|---|---|
| Oxygen Saturation, median (IQR) | 96.5 (4.5) (N=12) | 90.0 (5.0) | 85.5 (11.5) | 0.0005 |
| **Laboratory Data** | | | | |
| Infiltrate on Chest Radiograph, %, (N) | (42.8) (N=7) | (88.8) (N=18) | 100 (N=20) | 0.0008 |
| C-reactive Protein, median (IQR) | 29.6 (45.5) (N=7) | 107.0 (126.9) (N=18) | 155.4 (43.7) | 0.0004 |
| D-Dimer, median (IQR) | 2551 (2891) (N=6) | 724 (420) (N=18) | 1209 (1023) | 0.002 |
| Absolute Lymphocyte count, median (IQR) | 0.9 (0.8) (N=7) | 1.2 (0.7) (N=18) | 0.5 (0.5) | 0.002 |
| **Level of Care** | | | | |
| Intensive Care, No. (%) | 0 (0) | 0 (0) | 20 (100.0) | <0.0001 |
| Pressors, No. (%) | 0 (0) | 0(0) | 13 (65.0) | <0.0001 |
| **Worst Outcome, No. (%)** | | | | <0.0001 |
| Low Flow Oxygen Supplementation (~0-10 L) | 4 (28.6) | 12 (63.2) | 0 (0) | |
| Intermediate Flow Oxygen Supplementation (~10-20L) | 0 (0) | 2 (10.5) | 1 (5.0) | |
| High Flow Oxygen Supplementation (~20-60L) | 0 (0) | 0 (0) | 5 (25.0) | |
| Non invasive positive | 0 (0) | 0 (0) | 0 (0) | |

pressure ventilation

(NIPPV)

| | | | |
|---|---|---|---|
| Mechanical Ventilation | 0 (0) | 0 (0) | 12 (60.0) |
| Extracorporeal membrane oxygenation (ECMO) | 0 (0) | 0 (0) | 1 (5) |
| Death | 0 (0) | 0 (0) | 1 (5) |
| None of the above | 10 (71.4) | 5 (26.3) | 0 (0) |

1    *data presented for laboratory variables include data from less than the entire cohort.
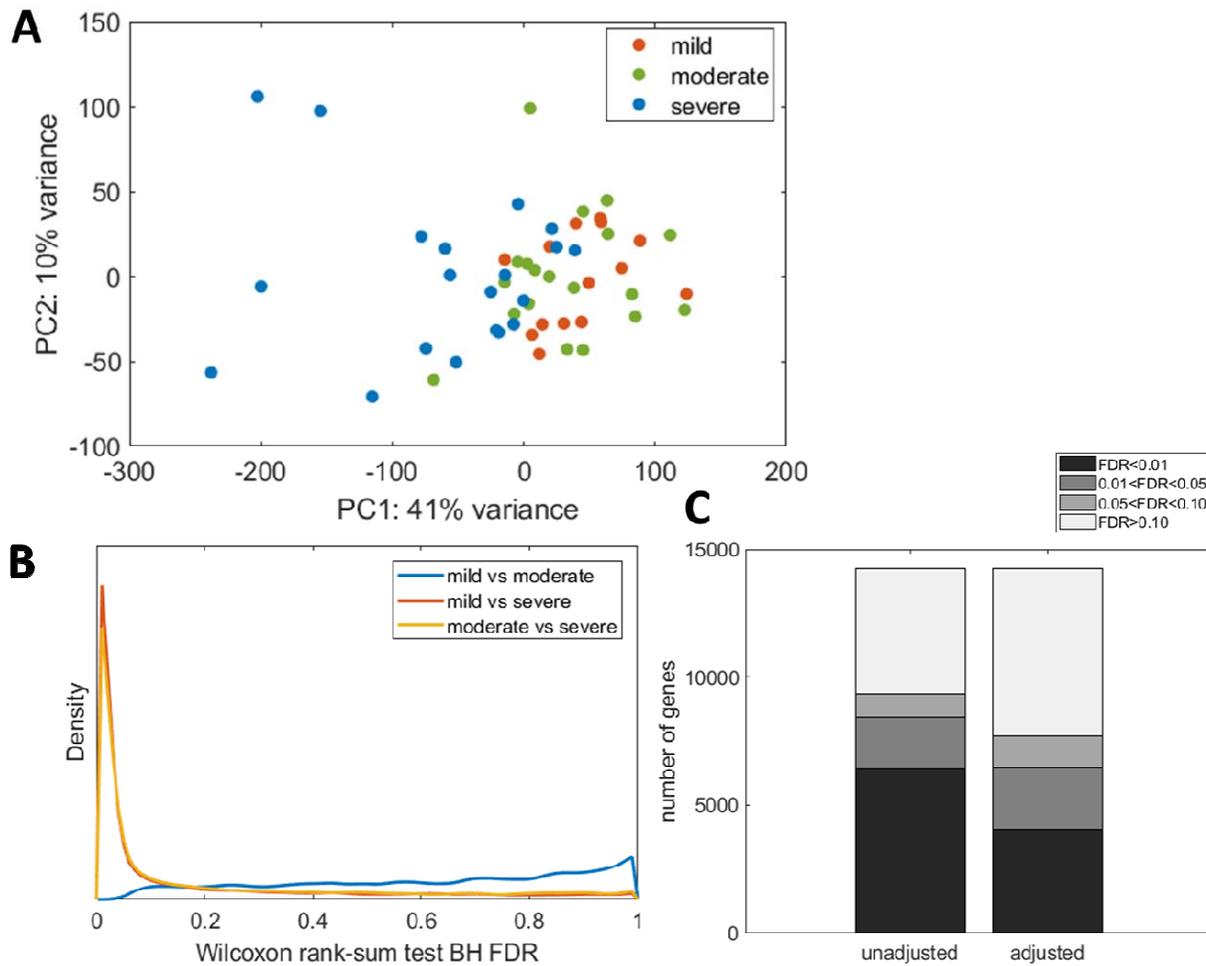
2

3

4
5

1    Figure 1. Analyses of signal levels in our dataset which was used as the training data.



2

3

4    A. Principal Components Analysis (PCA) plot for Z-score standardized CPM-normalized

5        gene expression, indexed by COVID severity.

6    B. Estimated densities of False Discovery Rates (FDR) for comparisons of COVID severity

7        levels based on nonparametric Wilcoxon tests of CPM-normalized gene expression

8        levels.

9    C. Numbers of differentially expressed genes, by FDR level, based on semiparametric Cox

10      proportional hazards models for gene expression as a function of severe vs non-severe
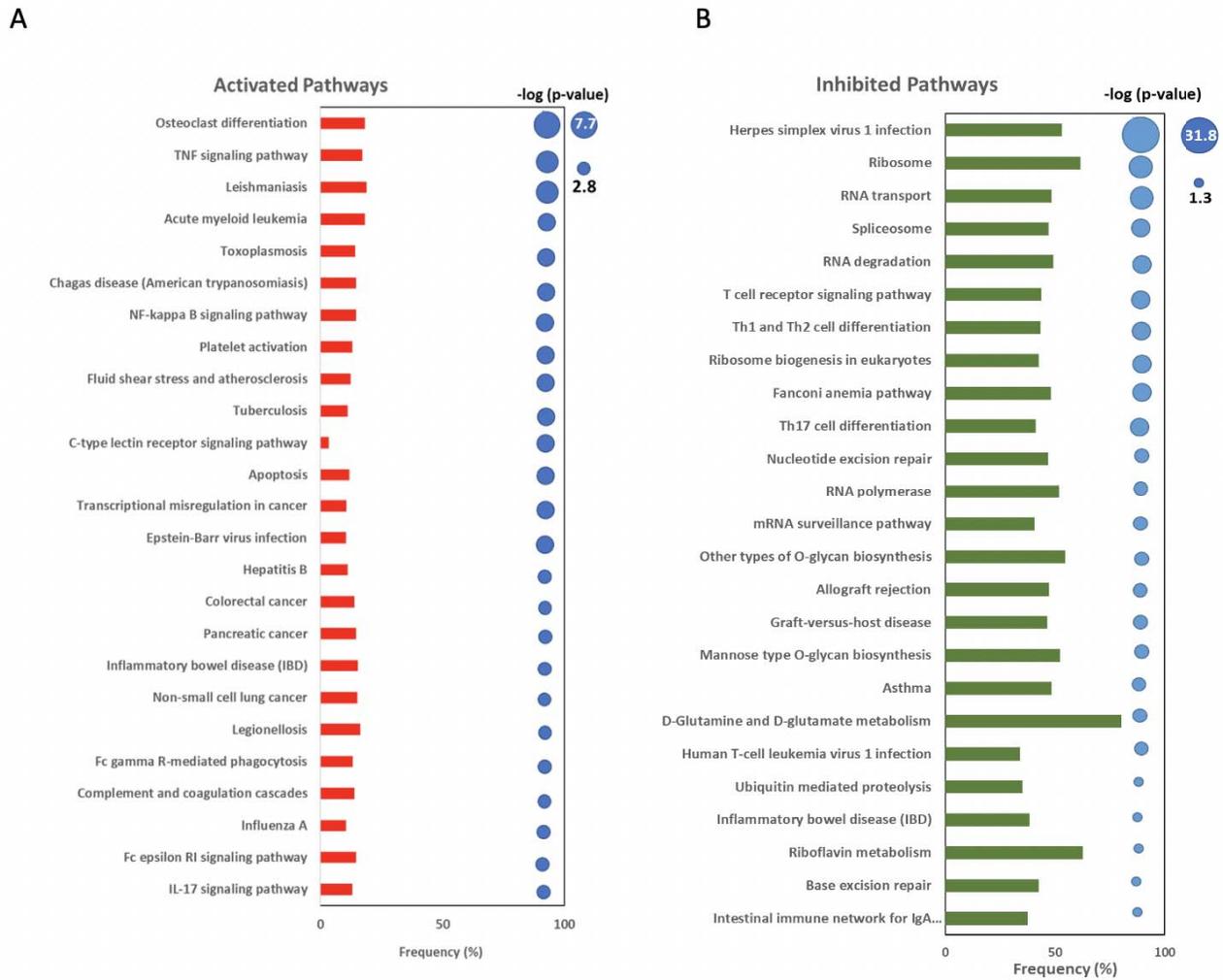
1    COVID, with and without adjustmentfor pre-specified covariates: race, sex, BMI, days
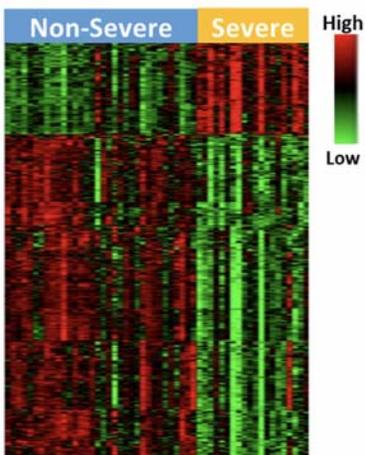
2    since symptom onset, and library size.

3

4

5

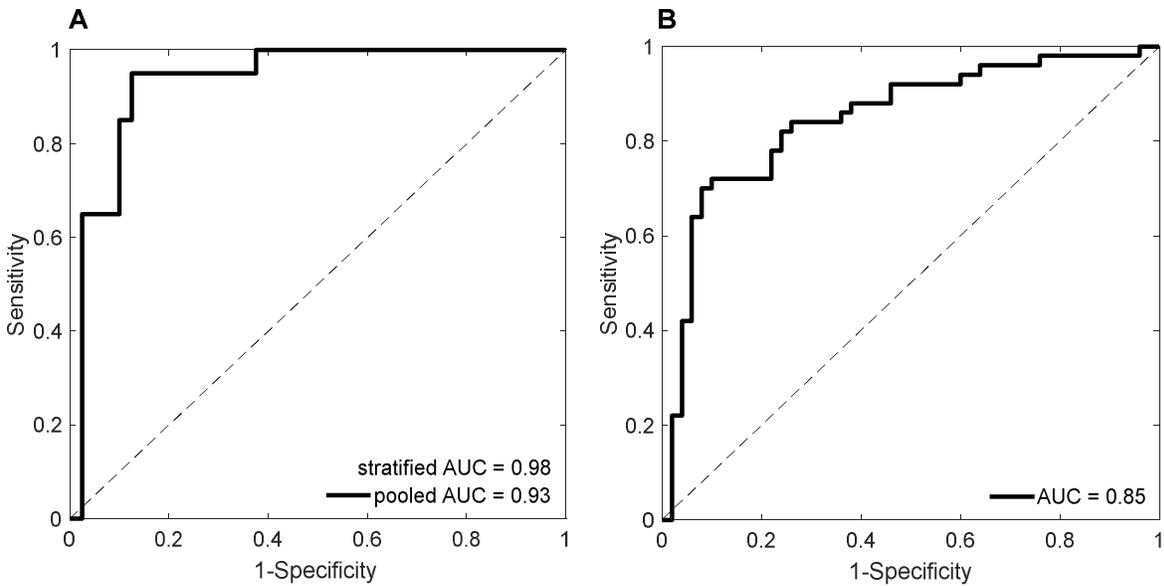Figure 2. Biological interpretation of gene expression patterns in severe versus non-severe COVID.

1    **(A, B)** Pathway analysis of genes differentially expressed between severe and non-severe

2    COVID-19 participants. Genes identified as overexpressed (n=936), and underexpressed

3    (n=5547) in severe cases of COVID-19, when compared to non-severe cases, were used for

4    pathway analysis using ENRICHR. With 936 genes overexpressed in severe COVID-19 cases,

5    ENRICHR (through KEGG Pathways database) identified 74 pathways associated with COVID-

6    19 severity, while with 5547 genes underexpressed in severe COVID-19, it identified 25

7    pathways. Shown here are the top 25 significant pathways ($p<0.05$) associated with upregulated

8    (A) and downregulated genes (B). The bar size represents the frequency of the pathway genes

9    differentially expressed in severe COVID-19; red indicates upregulation, and green indicates

10   downregulation. The size of the dots are proportional to -log(p-value). Larger dots represent

11   lower p-values. **(C)** Differential expression analysis of severe and non-severe COVID-19

12   participants identified 6483 genes as significantly different. Shown here is a heatmap of the top

13   425 differentially expressed genes, where the rows indicate genes and columns indicate

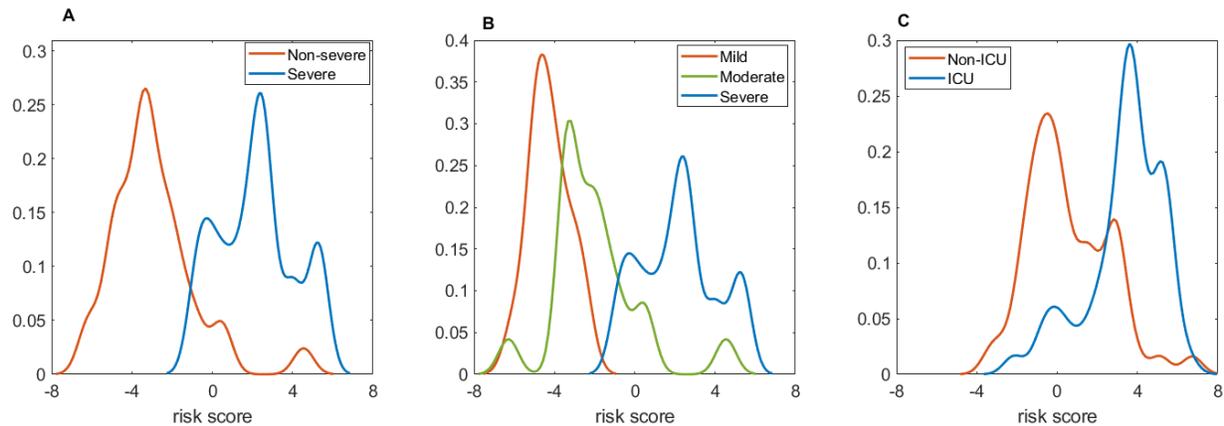14   participants. High expression is shown in red, and low expression in green.

15

1    Figure 3. Internally cross-validated and externally validated Receiver Operating Characteristic

2    (ROC) curves for the Weighted Gene Expression Risk Score (WGERS).



3

4    A. Cross-Validated (CV) ROC for severe vs non-severe COVID in the Training Data.

5        Pooled AUC corresponds with the plot (95% sensitivity and 88% specificity at a risk

6        score threshold of -1.04), which necessarily compares risk scores both within and across

7        CV folds (each with a unique fitted model). Stratified AUC corresponds with only

8        comparing risk scores within each fold of 20-fold CV (where each model is fixed), but

9        there exists no single corresponding ROC curve for this more commonly reported and

10       preferable metric.

11   B. ROC curve for ICU vs non-ICU in the Validation Data. A WGERS threshold of 1.77

12       yielded 84% sensitivity and 74% specificity.

13

14

15

16

17

1    Figure 4. Risk score distributions in the Training and Validation Data.



3    A.  Density of risk scores by non-severe (mild or moderate) vs severe COVID in the Training

4        Data. Apparent (non-cross-validated) sensitivity = 100% (20/20 severe) and specificity =

5        85% (28/33 non-severe) at the CV-optimal risk score threshold of -1.04.

6    B.  Density of risk scores by COVID severity (mild, moderate, or severe) in the Training

7        Data. Although the statistical learner was blinded to any distinction between mild and

8        moderate COVID severity, risk scores for moderate COVID participants fell between

9        those of mild and severe COVID participants, and all 5 misclassified participants had

10        moderate COVID.

11    C.  Density of risk scores by ICU vs non-ICU in the Validation Data. Sensitivity = 84% (42/50

12        ICU) and specificity = 74% (37/50 non-ICU) at a risk score threshold of 1.77.