

Airway Gene-Expression Classifiers for Respiratory Syncytial Virus (RSV) Disease Severity in Infants

Lu Wang

University of Rochester Medical Center

Chin-Yi Chu

University of Rochester Medical Center

Matthew N McCall

University of Rochester Medical Center

Christopher Slaunwhite

University of Rochester Medical Center

Jeanne Holden-Wiltse

University of Rochester Medical Center

Anthony Corbett

University of Rochester Medical Center

Ann R. Falsey

University of Rochester Medical Center

David J. Topham

University of Rochester Medical Center

Mary T. Caserta

University of Rochester Medical Center

Thomas J. Mariani (✉ Tom_Mariani@urmc.rochester.edu)

University of Rochester Medical Center

Edward E. Walsh

University of Rochester Medical Center

Xing Qiu

University of Rochester Medical Center <https://orcid.org/0000-0002-2330-3544>

Research article

Keywords: respiratory syncytial virus, respiratory severity score, gene expression, RNA-seq, classification

Posted Date: September 21st, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-63436/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published on February 25th, 2021. See the published version at <https://doi.org/10.1186/s12920-021-00913-2>.

Airway Gene-Expression Classifiers for Respiratory Syncytial Virus (RSV) Disease Severity in Infants

Lu Wang^{1†}, Chin-Yi Chu^{2†}, Matthew N. McCall¹, Christopher Slaunwhite², Jeanne Holden-Wiltse¹, Anthony Corbett¹, Ann R. Falsey^{3,5}, David J. Topham⁴, Mary T. Caserta², Thomas J Mariani^{2*}, Edward E. Walsh^{3,5*} and Xing Qiu^{1*}

¹Department of Biostatistics and Computational Biology, ²Department of Pediatrics, ³Department of Medicine, ⁴Department of Microbiology and Immunology, University of Rochester School Medicine; and ⁵Department of Medicine, Rochester General Hospital, Rochester, NY, USA

†These authors contributed equally to this work.

Corresponding authors:

Thomas J. Mariani

Email: Tom_Mariani@urmc.rochester.edu

Phone: 585-276-4616

Fax: 585-276-2643

Edward E. Walsh

Email: Edward.walsh@rochesterregional.org

Phone: 585-922-4331

Fax: 585-922-5168

Xing Qiu

Email: xing_qiu@urmc.rochester.edu

Phone: 585-275-0666

Fax: 585-273-1031

Keywords: respiratory syncytial virus; respiratory severity score; gene expression; RNA-seq; classification

Abstract word count: 197.

Text word count: 3273

1 **Abstract**

2 Background: A substantial number of infants infected with RSV develop severe
3 symptoms requiring hospitalization. We currently lack accurate biomarkers that are
4 associated with severe illness.

5 Method: We defined airway gene expression profiles based on RNA sequencing from
6 nasal brush samples from 106 full-term previously healthy RSV infected subjects during
7 acute infection (day 1-10 of illness) and convalescence stage (day 28 of illness). All
8 subjects were assigned a clinical illness severity score (GRSS). Using AIC-based model
9 selection, we built a sparse linear correlate of GRSS based on 41 genes (NGSS1). We
10 also built an alternate model based upon 13 genes associated with severe infection
11 acutely but displaying stable expression over time (NGSS2).

12 Results: NGSS1 is strongly correlated with the disease severity, demonstrating a naïve
13 correlation (ρ) of $\rho=0.935$ and cross-validated correlation of 0.813. As a binary classifier
14 (mild versus severe), NGSS1 correctly classifies disease severity in 89.6% of the
15 subjects following cross-validation. NGSS2 has slightly less, but comparable, accuracy
16 with a cross-validated correlation of 0.741 and classification accuracy of 84.0%.

17 Conclusion: Airway gene expression patterns, obtained following a minimally-invasive
18 procedure, have potential utility for development of clinically useful biomarkers that
19 correlate with disease severity in primary RSV infection.

20 **Introduction**

21 Respiratory Syncytial Virus (RSV) is the most important cause of respiratory illness in
22 infants and young children, accounting for more than 57,000 bronchiolitis and
23 pneumonia hospitalizations in the US annually.[1] Worldwide, 33.1 million acute lower
24 respiratory infections and 3.2 million hospitalizations in children under 5 years of age are
25 attributed to RSV each year.[2] In the US ~1-2% of newborns are hospitalized during
26 their first winter, with rates greatest in the first two months of life (25.9 per 1000).[3] Risk
27 factors for severe disease include gestational age < 29 weeks, bronchopulmonary
28 disease and symptomatic congenital cardiac disease, while less well defined risks
29 include lack of breast feeding, and exposure to tobacco smoke. However, the majority of
30 hospitalized infants are full-term infants whose only risk factor is young age at the time of
31 infection.[3]

32

33 A number of severity scores using clinical parameters, including cutaneous oximetry,
34 have been used to grade illness severity for use in management and as an outcome in
35 therapeutic, or potentially, vaccine trials. [4-13] However, none of the clinically based
36 severity scores have been universally adopted.[14] Reasons may include heterogeneity
37 in the scope and purpose of the score, the ages to which it is applied and concerns
38 about inter-observer variability and subjectivity in interpreting clinical signs, including
39 oximetry, that often are temporally dynamic over short intervals. Identification of an
40 objective biomarker that accurately correlates with, or potentially predicts, disease
41 severity could be highly useful.[15, 16]

42

43 We and others have reported a relationship between disease severity and host gene
44 expression in peripheral blood cells and nasal swab samples during infection.[17-20]

45 These results suggest such an approach may allow development of biomarkers to

46 accurately categorize RSV disease severity. As part of the AsPIRES study[21] we
47 recently reported on the feasibility of measuring gene expression of airway cells
48 collected by nasal swab in healthy infants in order to study RSV disease
49 pathogenesis.[22] However, in this manuscript, we describe the use of this gene
50 expression data during RSV infection to develop two airway gene expression-based
51 classifiers that are highly correlated with clinical disease severity. This represents a first
52 step in developing a biomarker using gene expression responses capable of accurately
53 classifying clinical severity in primary RSV-infection that could be used in conjunction
54 with clinical evaluation.

55 **Methods**

56 Study Subjects: Subjects included RSV infected infants enrolled in the AsPIRES study at
57 the University of Rochester Medical Center (URMC) and Rochester General Hospital
58 (RGH).[21] RSV-infected infants came from three cohorts during three winters (October
59 2012 through April 2015); one cohort included infants hospitalized with RSV, a second
60 cohort was recruited at birth and followed through their first winter for development of
61 RSV infection, and the third cohort was RSV infected infants seen in pediatric offices
62 and emergency departments and managed as outpatients. All subjects were full-term
63 infants undergoing a primary RSV infection during their first winter season. Nasal
64 samples were collected from the inferior nasal turbinate, by gentle brushing with a
65 flocked swab as previously described [22], during the acute illness visit (visit 1) and at a
66 convalescent visit ~28 after illness onset (visit 2). Illness severity was graded from 0-10
67 using a Global Respiratory Severity Score (GRSS), that uses nine parameters (age
68 adjusted respiratory rate, chest retractions, wheezing, rales/rhonchi, apnea, cyanosis,
69 room air oxygen saturation, lethargy and poor feeding) as previously described.[23] We

70 defined a GRSS >3.5 as severe disease as it is highly correlated with illness requiring
71 hospitalization.

72 Nasal RNA processing:

73 The process for nasal RNA recovery was previously described.[22] Briefly, following
74 flushing of the nares with 5 milliliters of saline to remove mucus and cellular debris, a
75 flocked swab was used to recover cells at the level of the turbinates. The swab was
76 immediately placed in RNA stabilizer (RNAprotect, Qiagen, Germantown, MD) and
77 maintained at 4 °C. Cells were recovered by filtering through a 0.45 µm membrane filter.
78 Recovered cells were lysed and homogenized using the AbsolutelyRNA Miniprep kit
79 (Agilent, Santa Clara, CA) according to the manufacturer's instructions. 1 ng of total
80 RNA was amplified using the SMARTer Ultra Low amplification kit (Clontech, Mountain
81 View, CA) and libraries were constructed using the NexteraXT library kit (Illumina, San
82 Diego, CA). Libraries were sequenced on the Illumina HiSeq2500. Sequences were
83 aligned against human genome version of hg19 using STARv2.5, counted with HTSeq,
84 and normalized by Fragments Per Kilobase of transcript per Million mapped reads
85 (FPKM). A total of 6,844 transcription profiles (genes) were reported after quality
86 assurance analysis and preprocessing. Additional technical details on data
87 preprocessing can be found in Supplementary Text.

88 Statistical methods:

89 Descriptive statistics are reported in Table 1. Discrete variables are summarized in
90 percentages, and continuous variables were summarized as Mean (STD). For
91 continuous variables, we performed two-sample Welch t-tests to check the equality
92 between the mild and severe groups; for categorical variables, Fisher's exact test was
93 used instead. The nasal gene-expression severity scores we developed in this study
94 were based on multivariate regression analysis with bi-directional stepwise model
95 selection based on Akaike Information Criterion (AIC). Technical details of model

96 development and cross-validation (CV) can be found in Supplementary Material. All
97 analyses were conducted using SAS 9.3 (SAS Institute Inc., Cary, NC, USA) and the R
98 programming language (version 3.5, R Foundation for Statistical Computing, Vienna,
99 Austria).

100 **Results**

101 Of the 139 RSV-infected infants enrolled in the AsPIRES study, nasal samples were
102 available from 119 subjects during acute infection (day 1-10 of illness) and 81 subjects
103 during convalescence (day 28 of illness). Among these 200 samples, 175 samples (106
104 acute samples and 69 convalescent samples) met sufficient quality to be used for
105 subsequent analyses. Demographic and clinical information for these 106 subjects are
106 provided in Table 1. The clinical severity score (GRSS) for these subjects ranged from 0
107 to 10, with 42 subjects considered to have mild disease (GRSS ≤ 3.5 ; mean \pm SE GRSS
108 of 1.63 ± 0.15) and 64 to have severe disease (GRSS > 3.5 ; mean GRSS of 6.13 ± 0.22).
109 There were no significant differences between the mild and severe groups in gender,
110 race, delivery type, breast feeding, or exposure to tobacco smoke. There also was no
111 difference in age at time of infection or in duration of illness at the time of evaluation.

112 **Nasal gene expression correlates of clinical severity during acute illness**

113 The 6,844 genes remaining after data preprocessing and filtering were subjected to the
114 Pearson correlation test to select genes that were significantly correlated with GRSS
115 during acute infection. After controlling the false discovery rate (FDR) at the 0.05 level,
116 66 significant genes were identified.[24] Using these genes, we applied model selection
117 procedures (see Supplementary Text for more details) to select an initial multivariate
118 regression model for GRSS (Model 1), which was comprised of 39 genes and had
119 relatively good predictive power (77.4% accuracy, or 24 misclassifications) for the

120 dichotomous clinical outcome (mild vs. severe illness) in leave-one-out cross-validation
121 (LOOCV).

122 Not unexpectedly, there is a strong correlation among the 66 genes, which might
123 reduce the diagnostic performance of Model 1. Using a novel method based on principal
124 component analysis (PCA), we identify ten supplementary genes as additional features
125 to model GRSS (see Supplementary Text for more details). With these additional
126 features and using the same model selection strategy, we developed two additional
127 models: Model 2 comprised of 41 genes and Model 3 comprised of 42 genes. The
128 performance of these models was evaluated by LOOCV (Table 2). We found that the
129 incorporation of the supplementary genes into Model 2 (CV prediction accuracy of
130 89.6%; 11 misclassifications) significantly improved the accuracy compared to Model 1
131 (24 misclassifications) and Model 3 (23 misclassifications). Of note, Model 2 contained 5
132 supplementary genes, and we defined it as **NGSS1** (nasal gene expression severity
133 score 1). As shown in Figure 1, NGSS1 is highly associative with GRSS (naïve $\rho=0.935$;
134 CV $\rho=0.813$). For the population of subjects in the AsPIRES study, the sensitivity and
135 specificity for identifying severe disease were high (sensitivity 90.1%, specificity 88%)
136 which would translate to a positive predictive value (PPV) of 92% and a negative
137 predictive value (NPV) of 86%.

138 **Validation of NGSS1 at the Convalescence Phase**

139 NGSS1 was trained exclusively from data collected at the acute phase (visit 1). For a
140 subset ($n=54$) of subjects, we also had their nasal transcriptome profiles at the
141 convalescence phase (day 28 after illness onset), a time when most infants had
142 completely recovered from their illness. If NGSS1 is a valid surrogate for disease
143 severity, we hypothesized that NGSS1 calculated from the severely ill subjects at visit 2
144 would converge to those of the mildly ill subjects. Compared with the acute visit, the
145 calculated NGSS1 at the convalescent visit predicted a significantly lower mean severity

146 score for severe subjects (n=29, 6.22 vs. 2.82, $p < .001$). In contrast, there was no
147 significant difference in NGSS1 between the two visits for the mildly ill group (n=25, 1.96
148 vs. 2.31, $p = 0.45$), nor between the severe and mild groups at visit 2 (2.82 vs. 2.31,
149 $p = 0.40$). These results are illustrated in Figure 2A.

150 **Exploratory Association Analysis Based on Stable Nasal Genes**

151 In the process of developing NGSS1 we observed that a large number of genes had
152 expression levels that remained stable between the acute and convalescent visits. We
153 speculated that a NGSS based on stable genes that were correlated with GRSS could
154 potentially be predictive of disease severity prior to illness onset. Thus, we next
155 developed an NGSS based on genes displaying stable expression across acute illness
156 and convalescence in the 54 subjects with samples from both time points. Specifically,
157 we included only genes whose mean expression levels correlated with disease severity
158 during acute illness, and whose expression did not change significantly from the acute to
159 convalescent stage.

160 We identified 2127 genes in subjects with mild illness and 1531 genes in subjects with
161 severe illness, based on paired two sample t-test ($p > 0.5$) and fold change increases or
162 decreases within 10%. Of the total 3658 genes, 689 stable genes were common in both
163 groups (Figure 3A). A quality assurance analysis based on IQR showed that a small
164 subset (n=14) of these genes had relatively small dynamic range in the combined
165 dataset, and were excluded. We applied marginal screening based on Pearson
166 correlation with GRSS to the remaining 675 stable genes and identified 44 marginally
167 significant genes. As in developing NGSS1, we added 5 supplementary genes with
168 strong marginal associations with GRSS. Model selection identified 13 genes as Model 4
169 (designated as **NGSS2**). The performance of NGSS2 is provided in Table 2 and
170 illustrated in Figure 3B. NGSS2 showed a significant correlation with GRSS ($\rho = 0.741$),
171 and a CV accuracy of 84% (17 misclassifications out of 106 cases, Table 2). Of note,

172 NGSS1 and NGSS2 do not contain any commonly selected gene, which is expected due
173 to different screening criteria. Figure 2B shows that on average, NGSS2 did not change
174 between visit 1 and visit 2, which is the key difference between these two classifiers. A
175 full list of genes used in NGSS1 and NGSS2, as well as their estimated linear
176 coefficients in the models, are listed in Supplementary Tables E2 and E3.

177 **Discussion**

178 Several approaches have been proposed for quantifying RSV disease severity in young
179 infants.[4-13] A variety of clinical parameters have been included in several described
180 severity scores, with incomplete agreement on the optimal factors to select.[14] One
181 reason is that many clinical signs of RSV infection in young infants, including cutaneous
182 oximetry, can fluctuate frequently and rapidly during the course of illness, making
183 consistent assessment difficult. In fact, even the direct measurement of RSV viral load in
184 serum is not significantly correlated with disease severity in the AsPIRES study [25] --
185 similar phenomenon was also reported by several other similar studies.[26-28] An
186 objective biomarker reliably correlated with clinical severity could prove useful for clinical
187 management and as a classifier and/or an outcome measure in vaccine or therapeutic
188 trials.

189

190 Transcriptomic analysis of host cells has proven informative in the study of several
191 respiratory viral infections, including RSV, with the emphasis on disease
192 pathogenesis.[17-20] Unlike this report that focuses on nasal epithelial cell samples,
193 most reports have described gene expression correlates of disease severity in peripheral
194 blood mononuclear cells during infection since RSV pathogenesis is thought to be
195 closely linked to the host's immune response.[29] In two publications from the same
196 group, RSV infection was associated with overexpression of innate immunity genes

197 (neutrophil and interferon genes) and suppression of adaptive T and B cell genes. [17,
198 19] The investigators used the results to develop a gene-expression based illness score
199 (designated Molecular Distance to Health [MDTH]) that was significantly correlated with
200 a clinical disease severity score, duration of hospitalization and need for supplemental
201 oxygen. Recently, Jong et al described an 84 gene signature that was highly predictive
202 of RSV disease severity in infants.[16] Similarly, we reported that gene expression
203 patterns in purified blood CD4 T cells during infection were correlated with clinical
204 disease severity.[18] Gene expression results from nasal swabs collected from
205 hospitalized infants during RSV infection have also been recently reported by another
206 group, with differentially expressed genes correlated with clinical severity.[20]

207

208 In this report we describe the use of RNAseq analysis of gene expression data from
209 nasal specimens collected during RSV infection to develop two nasal gene-expression
210 severity scores (NGSS1 and NGSS2) that are highly correlated with a clinically derived
211 disease severity score (GRSS). Although the nasal brush samples from the AsPIRES
212 study were collected to investigate molecular pathways and disease mechanisms
213 involved in pathogenesis (presented in a separate manuscript [30]), we also considered
214 that the data could be useful for the development of a gene based biomarker of RSV
215 severity. We used marginal screening of all genes followed by PCA analysis and step-
216 wise model selection to develop NGSS1, a multivariate linear classifier of severity. In CV
217 analysis, NGSS1 was strongly correlated with GRSS and was a relatively accurate
218 classifier of binary disease severity. Furthermore, the score tracked well with clinical
219 improvement 28 days after illness onset. Of particular note, we found that including
220 uncorrelated supplementary genes enhances the accuracy of the models, and
221 recommend this approach as a routine for future classification/prediction analyses based
222 on high-throughput data with substantial correlation. As noted, in the population enrolled

223 in our study the operating characteristics of NGSS1, including sensitivity, specificity,
224 PPV and NPV, were quite good. However, it should be recognized that the proportion of
225 mildly ill to severely ill subjects was determined by the recruiting strategy used, and that
226 the PPV and NPV would vary depending on the population to which NGSS1 was
227 applied.[21] If mildly ill subjects are increased by a factor of 3-5 this would reduce the
228 PPV to 40-70% although the NPV would remain >90%.

229

230 Although the aim of this report is not to describe molecular mechanisms operative during
231 RSV infection, it should be noted that the 41 NGSS1 genes include cytokines
232 (TNFSF10, IL6, and CXCL2), extracellular matrix proteins (VIM, MMP19, RPS15A,
233 FKBP1A, and VCAN), inflammation regulators (CXCL2, CD163), and components of
234 various signaling processes (GNS, HAVCR2, PTPRC, CTSL, INHBA, IL6, MMP19,
235 CXCL2, SLC39A8, CCDC80, VCAN, CD163). Some genes are only known to be
236 involved in fundamental biological processes and are therefore novel in RSV research,
237 including ST3GAL1 (a type II membrane protein) and ATP10B (ATPase Phospholipid
238 Transporting 10B). Note that only two genes (TNFSF10, RABGAP1L) have been
239 associated with disease severity in our recent study based on purified CD4 T cells.[18] In
240 addition, IL-6 Signaling is the only significant canonical pathway identified from the CD4
241 T cells that contains an NGSS1 gene (IL6).

242

243 A unique and very preliminary result from our analysis is the development of NGSS2
244 using differentially expressed genes associated with GRSS that did not change between
245 the acute and the convalescent time points. It is possible that these genes may simply
246 be slow to return to baseline expression levels, in contrast to those genes selected for
247 NGSS1. Although speculative, it occurred to us that “stable” genes might possibly be
248 predictive of severity regardless of when a nasal sample was obtained, thus raising the

249 possibility of infants at risk prior to or early in infection. While NGSS2 is slightly less
250 accurate than NGSS1 in predicting GRSS during acute illness, the association between
251 NGSS2 and GRSS is still relatively strong. Interestingly, the 13 NGSS2 genes were
252 broadly related to cytoplasmic activities (EXOSC10, PLK2, PPIC, CLDN10, MAP3K13,
253 MT1G, PXN), ATP binding (SEPHS2) and phosphoprotein regulation (BCKDK, PLK2,
254 MAP3K13); activities that may be less directly responsive to acute RSV infection. These
255 observations suggest that the best nasal transcriptome predictors of respiratory
256 symptoms are not necessarily limited to those genes that directly regulate the immune
257 response to RSV infection.

258

259 The use of nasal brush specimens for development of a severity biomarker in infants is
260 attractive for a number of reasons. Nasal respiratory epithelial cells are the first cells
261 infected and directly initiate early innate immune responses to RSV. The mucosa is also
262 the site of migration of both innate and adaptive immune cells during infection.
263 Importantly, we have shown that gene expression in nasal respiratory epithelial cells is
264 highly concordant with published gene expression in lower respiratory tract epithelial
265 cells, and thus should be a reasonable proxy for lung responses to RSV infection.[22] Of
266 practical importance, collection of nasal epithelial cells is relatively non-invasive and
267 simple to perform with minimal discomfort.

268

269 There are several important limitations to our study and conclusions. First, we do not
270 have an independent cohort to validate our findings; the only publically available nasal
271 gene expression data during RSV infection used microarray technology that did not
272 identify many of the genes we identified by RNAseq. Due to the lack of independent
273 samples for validation, we applied cross-validation techniques to prevent model

274 overfitting and validate the accuracy of prediction for both NGSS1 and NGSS2 at the
275 acute visit. CV estimator for prediction accuracy is known to be asymptotically unbiased
276 [31] under very weak statistical assumptions, namely, the training and testing data are
277 independent and identically distributed (which can even be relaxed further, see [32, 33]).
278 Additionally, we further validated the NGSS1 trained at the acute visit with the
279 convalescence data, and the results conformed with our prediction remarkably well.
280 Although the NGSS1 declined for the severely ill infants when clinical symptoms had
281 resolved, it would be useful to determine if NGSS1 tracked closely over the full course of
282 an illness. However, validation of our findings with an independent prospective cohort
283 will be required. In addition, the results may not be valid for infants older than 10 months
284 of age when infected with RSV, nor for infants with prematurity or other underlying
285 medical conditions.

286 Another possible limitation is that all data used in these analyses were generated on the
287 same technical platform and processed by the same team, therefore the validation
288 results do not reflect the impact of “artifacts” in transcriptomic studies such as batch
289 effects and platform differences, which can be reduced but not entirely eradicated by
290 advanced normalization methods.[34-37]

291 Importantly, speculation that NGSS2 might predict disease severity prior to infection
292 demands careful prospective validation. Finally, to extend the utility of time-intensive
293 gene expression assays beyond a research tool and use it as a clinically useful
294 biomarker of RSV disease severity, will require translation of these results to a rapid
295 readily performed multiplex reverse transcription polymerase chain reaction (RT-PCR)
296 assay, similar to those that have recently been developed for microbial diagnostics in
297 respiratory secretions.[38]

298 In conclusion, we demonstrate that analysis of gene expression data obtained from an
299 easily and safely obtained nasal brush specimen in young infants with acute RSV
300 infection shows promise for development of composite molecular biomarkers that closely
301 correlate with clinical severity score. Further studies to refine and validate the potential
302 of predictive gene expression data from readily collected nasal samples are needed.

303 **Funding**

304 This study is supported in part by Respiratory Pathogens Research Center (NIAID
305 contract number HHSN272201200005C), and the University of Rochester CTSA award
306 number UL1 TR002001 from the National Center for Advancing Translational Sciences
307 of the National Institutes of Health. The content is solely the responsibility of the authors
308 and does not necessarily represent the official views of the National Institutes of Health.

309

310 **Availability of data and materials**

311 The transcriptional data described in this manuscript are available in dbGaP
312 (phs001201.v2.p1).

313

314 **Author information**

315 **Affiliations**

316 **Department of Biostatistics and Computational Biology, University of Rochester**
317 **School Medicine, Rochester, New York, U.S.A.**

318 Lu Wang, Matthew N. McCall, Jeanne Holden-Wiltse, Xing Qiu

319 **Department of Pediatrics, University of Rochester School Medicine, Rochester,**
320 **New York, U.S.A.**

321 Chin-Yi Chu, Christopher Slaunwhite, Mary T. Caserta, Thomas J Mariani

322 **Department of Medicine, University of Rochester School Medicine, Rochester,**
323 **New York, U.S.A.**

324 Ann R. Falsey, Edward E. Walsh

325 **Department of Microbiology and Immunology, University of Rochester School**
326 **Medicine, Rochester, New York, U.S.A.**

327 David J. Topham

328 **Department of Medicine, Rochester General Hospital, Rochester, New York, U.S.A.**

329 Edward E. Walsh

330

331 **Contributions**

332 XQ, TJM, and EEW conceptualized the study. TJM, EEW, MTC, and CC designed the

333 experiments. EEW, MTC, ARF and DJT developed the cohort, and collected the

334 specimens and clinical data. LW, MNM, and XQ developed statistical models. JHW and

335 AC facilitated data organization, management and analysis. LW, CC, MNM, CS, JH-W,

336 AC, ARF, DJT, MTC, TJM, EEW, and XQ generated, analyzed and interpreted the data.

337 LW, CC, MNM, CS, JH-W, AC, ARF, DJT, MTC, TJM, EEW, and XQ wrote and/or

338 revised the manuscript. All authors read and approve the final manuscript.

339

340 **Corresponding author**

341 Correspondence to Thomas J Mariani, Edward E. Walsh, and Xing Qiu.

342

343 **Ethics declarations**

344 **Ethics approval and consent to participate**

345 This study was approved by the Institutional Review Boards of the University of

346 Rochester Medical Center (URMC) and Rochester General Hospital (RGH). All parents

347 provided written informed consent.

348

349 **Consent for publication**

350 Not applicable.

351

352 **Competing interests**

353 The authors declare that they have no competing interests.

Tables and Figures

Table 1. Demographic data of subjects. *P*-values reported in the last column were either based on Fisher's exact test (if the variable is categorical) or Welch *t*-test (if the variable is continuous). Continuous variables are reported as sample means (STD); categorical variables are reported as percentages.

	mild (n=42) ^a		severe (n=64) ^a		p value
	n	Mean (STD) or %	n	Mean (STD) or %	
Global Severity Score	42	1.63(1.00)	64	6.13(1.72)	<0.001
Visit Age (months)	42	3.52(1.99)	64	3.24(2.37)	0.5122
Gestational Age (weeks)	42	39.05(1.25)	64	38.8(1.44)	0.3437
Birth Weight (kg)	42	3.32(0.68)	64	3.36(0.57)	0.7468
Family Size	42	4.43(2.86)	64	3.98(1.73)	0.3703
Days Since Disease Onset	42	4.31(1.76)	64	4.86(1.78)	0.1209
Breast Feeding Summary	42	1.56(1.23)	63	1.53(1.25)	0.8979
Sex					
Male	23	44.23	29	55.77	0.4275
Female	19	35.19	35	64.81	
Ethnicity					
Hispanic or Latino	8	42.11	11	57.89	0.8018
Non-Hispanic or Non-Latino	34	39.08	53	60.92	
Race					
Caucasian	23	37.1	39	62.9	0.3115
Other race	19	47.5	21	52.5	
Missing	-	0	4	100	
Delivery Type					
Vaginal	29	36.71	50	63.29	0.3634
C-section	13	48.15	14	51.85	
Smoking Exposure					
Yes	14	38.89	22	61.11	1
No	28	40	42	60	
RSV Group					
A	23	38.98	36	61.02	1
B	18	39.13	28	60.87	
Missing	1	100	-	0	

^a based on GRSS ≤ 3.5 (mild) or > 3.5 (severe)

Table 2. Performance of four models used in developing NGSS1 and NGSS2. Naïve and CV RSS are the mean residual sums of squares of the predictive model in the original and cross-validation analyses, respectively. Correlation are the Pearson correlation coefficient between the predicted severity scores and the clinically defined GRSS. Prediction accuracy is the percentage of correctly predicted mild (NGSS ≤ 3.5) or severe (NGSS > 3.5) symptoms, compared with the same phenotype defined by the GRSS (mild: GRSS ≤ 3.5 ; severe: GRSS > 3.5).

	number of genes selected	Naïve RSS	Naïve Correlation	Naïve misclassified subjects (out of 106)	CV RSS	CV Correlation	CV prediction accuracy	CV misclassified subjects (out of 106)
Model 1	39 genes	1.234	0.909	15	2.743	0.797	77.4%	24
Model 2*	41 genes	0.884	0.935	9	2.681	0.813	89.6%	11
Model 3	42 genes	0.920	0.933	13	2.119	0.844	78.3%	23
Model 4**	13 genes	2.549	0.800	16	3.215	0.741	84.0%	17

* designated NGSS1. ** designated NGSS2

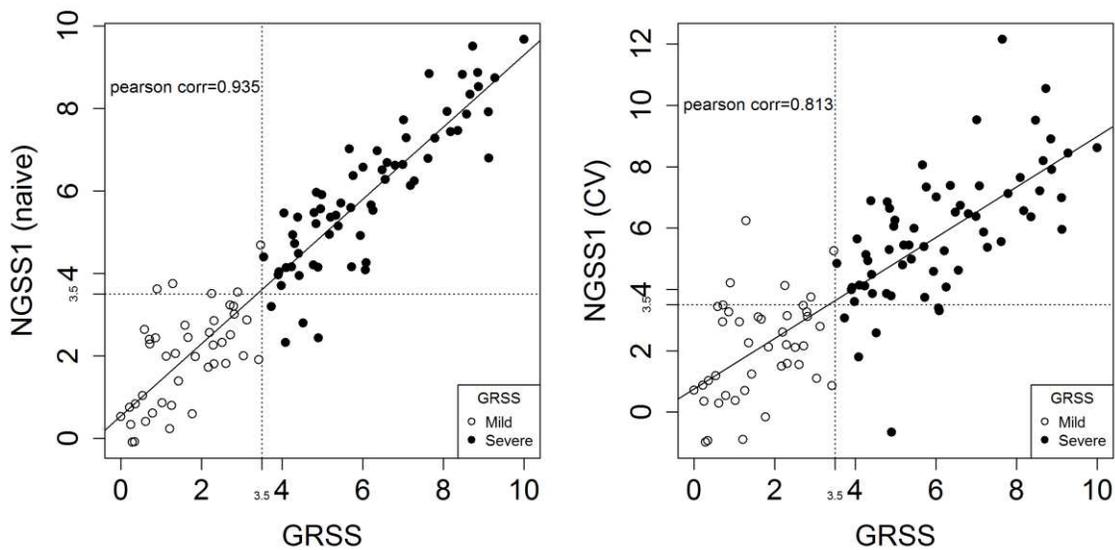


Figure 1. Correlating NGSS1 (severity score predicted by Model 2) with GRSS. Left: naïve Pearson correlation between GRSS and NGSS1 is $\rho = 0.935$. Right: cross-validated Pearson correlation between GRSS and NGSS is $\rho = 0.813$. Solid dots are subjects with severe symptoms (defined by $GRSS > 3.5$) and empty dots are those with mild symptoms ($GRSS \leq 3.5$).

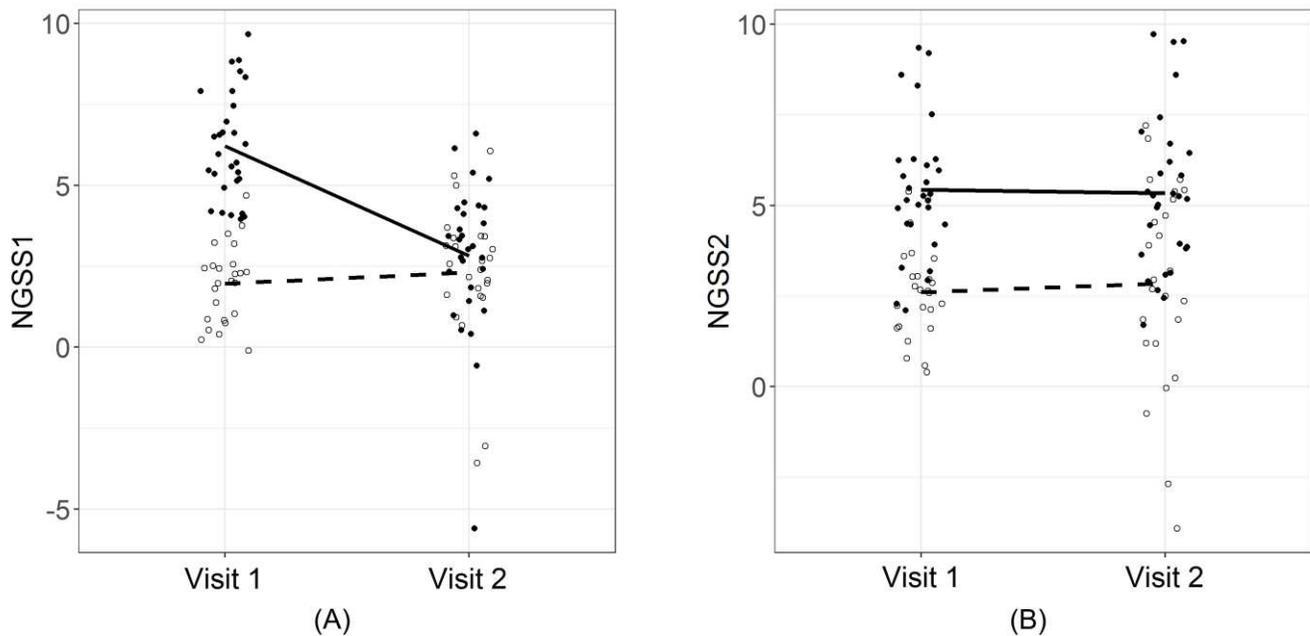


Figure 2. Paired comparisons between visit 1 and visit 2 using NGSS1 (panel (A)) and NGSS2 (panel (B)). A total of $n=54$ subjects with samples in both visits were used. Solid dots represent severe subjects and empty dots represent mild subjects. The solid line represents the mean trend of severe subjects and the broken line represents the mean trend for mild subjects. (A): At visit 1, there was a significant difference in mean NGSS1 between the severe ($n=29$) and mild ($n=25$) groups (6.22 vs. 1.96, $p<0.001$). Mean NGSS1 of the mild group was virtually unchanged between two visits (1.96 vs. 2.31, $p=0.45$). In comparison, mean NGSS1 of the severe group declined significantly at visit 2 (6.22 vs. 2.82, $p<.001$). (B): In contrast to NGSS1, the differences in NGSS2 was virtually unchanged between the two visits, due to the fact that NGSS2 were built with stable genes.

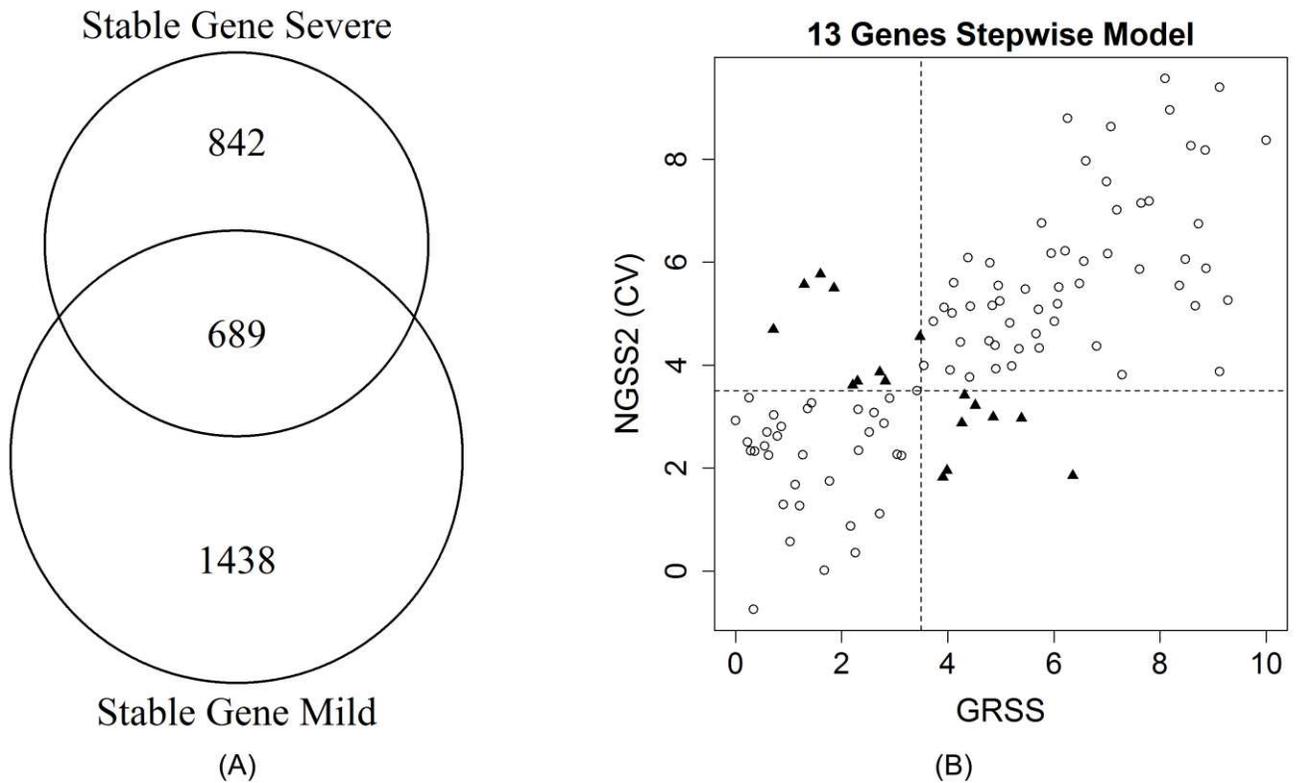


Figure 3. (A). Diagram indicating the stable genes for the mild ($GRSS \leq 3.5$) and severe ($GRSS > 3.5$) groups and the 689 intersecting stable genes common to both groups. (B). Correlating NGSS2 (severity score predicted by Model 4) with GRSS. Naïve Pearson correlation between GRSS and NGSS2 is $\rho = 0.800$. Right: cross-validated Pearson correlation between GRSS and NGSS is $\rho = 0.741$. Circles are subjects with correct cross-validated classification based on NGSS2; solid triangles are misclassified subjects.

References

1. Hall CB, Weinberg GA, Iwane MK, et al. The burden of respiratory syncytial virus infection in young children. *N Engl J Med* **2009**; 360:588-98.
2. Shi T, McAllister DA, O'Brien KL, et al. Global, regional, and national disease burden estimates of acute lower respiratory infections due to respiratory syncytial virus in young children in 2015: a systematic review and modelling study. *Lancet* **2017**; 390:946-58.
3. Hall CB, Weinberg GA, Blumkin AK, et al. Respiratory syncytial virus-associated hospitalizations among children less than 24 months of age. *Pediatrics* **2013**; 132:e341-8.
4. Bekhof J, Reimink R, Brand PL. Systematic review: insufficient validation of clinical scores for the assessment of acute dyspnoea in wheezing children. *Paediatr Respir Rev* **2014**; 15:98-112.
5. Corneli HM, Zorc JJ, Holubkov R, et al. Bronchiolitis: clinical characteristics associated with hospitalization and length of stay. *Pediatr Emerg Care* **2012**; 28:99-103.
6. Destino L, Weisgerber MC, Soung P, et al. Validity of respiratory scores in bronchiolitis. *Hosp Pediatr* **2012**; 2:202-9.
7. Duarte-Dorado DM, Madero-Orostegui DS, Rodriguez-Martinez CE, Nino G. Validation of a scale to assess the severity of bronchiolitis in a population of hospitalized infants. *J Asthma* **2013**; 50:1056-61.
8. Feldman AS, Hartert TV, Gebretsadik T, et al. Respiratory Severity Score Separates Upper Versus Lower Respiratory Tract Infections and Predicts Measures of Disease Severity. *Pediatr Allergy Immunol Pulmonol* **2015**; 28:117-20.
9. Gajdos V, Beydon N, Bommenel L, et al. Inter-observer agreement between physicians, nurses, and respiratory therapists for respiratory clinical evaluation in bronchiolitis. *Pediatr Pulmonol* **2009**; 44:754-62.
10. McCallum GB, Morris PS, Wilson CC, et al. Severity scoring systems: are they internally valid, reliable and predictive of oxygen use in children with acute bronchiolitis? *Pediatr Pulmonol* **2013**; 48:797-803.
11. Mosalli R, Abdul Moez AM, Janish M, Paes B. Value of a risk scoring tool to predict respiratory syncytial virus disease severity and need for hospitalization in term infants. *J Med Virol* **2015**; 87:1285-91.
12. Parker MJ, Allen U, Stephens D, Lalani A, Schuh S. Predictors of major intervention in infants with bronchiolitis. *Pediatr Pulmonol* **2009**; 44:358-63.
13. Fernandes RM, Plint AC, Terwee CB, et al. Validity of bronchiolitis outcome measures. *Pediatrics* **2015**; 135:e1399-408.
14. Karron RA, Zar HJ. Determining the outcomes of interventions to prevent respiratory syncytial virus disease in children: what to measure? *The Lancet Respiratory medicine* **2018**; 6:65-74.
15. Brown PM, Schneeberger DL, Piedimonte G. Biomarkers of respiratory syncytial virus (RSV) infection: specific neutrophil and cytokine levels provide increased accuracy in predicting disease severity. *Paediatr Respir Rev* **2015**; 16:232-40.
16. Jong VL, Ahout IM, van den Ham H-J, et al. Transcriptome assists prognosis of disease severity in respiratory syncytial virus infected infants. *Scientific reports* **2016**; 6:1-12.
17. de Steenhuijsen Piters WA, Heinonen S, Hasrat R, et al. Nasopharyngeal Microbiota, Host Transcriptome, and Disease Severity in Children with Respiratory Syncytial Virus Infection. *Am J Respir Crit Care Med* **2016**; 194:1104-15.
18. Mariani TJ, Qiu X, Chu C, et al. Association of Dynamic Changes in the CD4 T-Cell Transcriptome With Disease Severity During Primary Respiratory Syncytial Virus Infection in Young Infants. *J Infect Dis* **2017**; 216:1027-37.
19. Mejias A, Dimo B, Suarez NM, et al. Whole blood gene expression profiles to assess pathogenesis and disease severity in infants with respiratory syncytial virus infection. *PLoS Med* **2013**; 10:e1001549.
20. Do LAH, Pellet J, van Doorn HR, et al. Host Transcription Profile in Nasal Epithelium and Whole Blood of Hospitalized Children Under 2 Years of Age With Respiratory Syncytial Virus Infection. *J Infect Dis* **2017**; 217:134-46.
21. Walsh EE, Mariani TJ, Chu C, et al. Aims, Study Design, and Enrollment Results From the Assessing Predictors of Infant Respiratory Syncytial Virus Effects and Severity Study. *JMIR Res Protoc* **2019**; 8:e12907.
22. Chu CY, Qiu X, Wang L, et al. The Healthy Infant Nasal Transcriptome: A Benchmark Study. *Sci Rep* **2016**; 6:33994.
23. Caserta MT, Qiu X, Tesini B, et al. Development of a Global Respiratory Severity Score for Respiratory Syncytial Virus Infection in Infants. *J Infect Dis* **2017**; 215:750-6.
24. Benjamini Y, Hochberg Y. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B* **1995**; 57:289-300.

25. Walsh EE, Wang L, Falsey AR, et al. Virus-Specific Antibody, Viral Load, and Disease Severity in Respiratory Syncytial Virus Infection. *J Infect Dis* **2018**; 218:208-17.
26. Wright PF, Gruber WC, Peters M, et al. Illness severity, viral shedding, and antibody responses in infants hospitalized with bronchiolitis caused by respiratory syncytial virus. *J Infect Dis* **2002**; 185:1011-8.
27. Yan XL, Li YN, Tang YJ, et al. Clinical characteristics and viral load of respiratory syncytial virus and human metapneumovirus in children hospitalized for acute lower respiratory tract infection. *J Med Virol* **2017**; 89:589-97.
28. Piedra FA, Mei M, Avadhanula V, et al. The interdependencies of viral load, the innate immune response, and clinical outcome in children presenting to the emergency department with respiratory syncytial virus-associated bronchiolitis. *PLoS One* **2017**; 12:e0172953.
29. Collins PL, Fearn R, Graham BS. Respiratory syncytial virus: virology, reverse genetics, and pathogenesis of disease. *Curr Top Microbiol Immunol* **2013**; 372:3-38.
30. Chu C-Y, Qiu X, McCall MN, et al. Insufficiency in airway interferon activation defines clinical severity to infant RSV infection. *bioRxiv* **2019**:641795.
31. Seber GA, Lee AJ. Linear regression analysis. Vol. 329. John Wiley & Sons, **2012**.
32. Opsomer J, Wang Y, Yang Y. Nonparametric regression with correlated errors. *Statistical Science* **2001**:134-53.
33. Arlot S, Celisse A. A survey of cross-validation procedures for model selection. *Statistics surveys* **2010**; 4:40-79.
34. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* **2007**; 8:118-27.
35. Leek JT, Storey JD. Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS Genet* **2007**; 3:1724-35.
36. Rudy J, Valafar F. Empirical comparison of cross-platform normalization methods for gene expression data. *BMC Bioinformatics* **2011**; 12:467.
37. Qiu X, Hu R, Wu Z. Evaluation of bias-variance trade-off for commonly used post-summarizing normalization procedures in large-scale gene expression studies. *PLoS One* **2014**; 9:e99380.
38. Lee SH, Ruan SY, Pan SC, Lee TF, Chien JY, Hsueh PR. Performance of a multiplex PCR pneumonia panel for the identification of respiratory pathogens and the main determinants of resistance from the lower respiratory tract specimens of adult patients in intensive care units. *J Microbiol Immunol Infect* **2019**; 52:920-8.
39. Grier A, Gill AL, Kessler HA, et al. Temporal Dysbiosis of Infant Nasal Microbiota Relative to Respiratory Syncytial Virus Infection. *bioRxiv* **2020**.

Figures

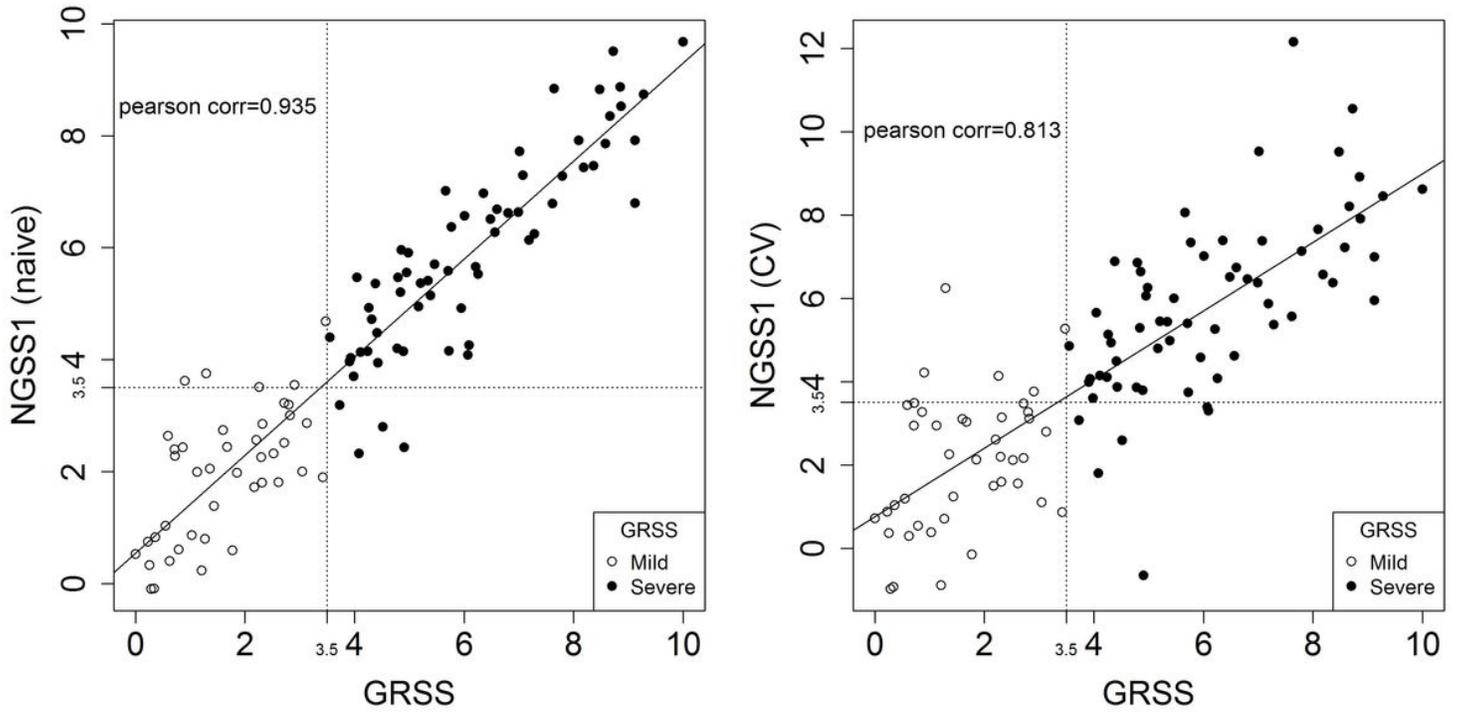


Figure 1

Correlating NGSS1 (severity score predicted by Model 2) with GRSS. Left: naive Pearson correlation between GRSS and NGSS1 is $r=0.935$. Right: cross-validated Pearson correlation between GRSS and NGSS1 is $r=0.813$. Solid dots are subjects with severe symptoms (defined by $GRSS > 3.5$) and empty dots are those with mild symptoms ($GRSS \leq 3.5$).

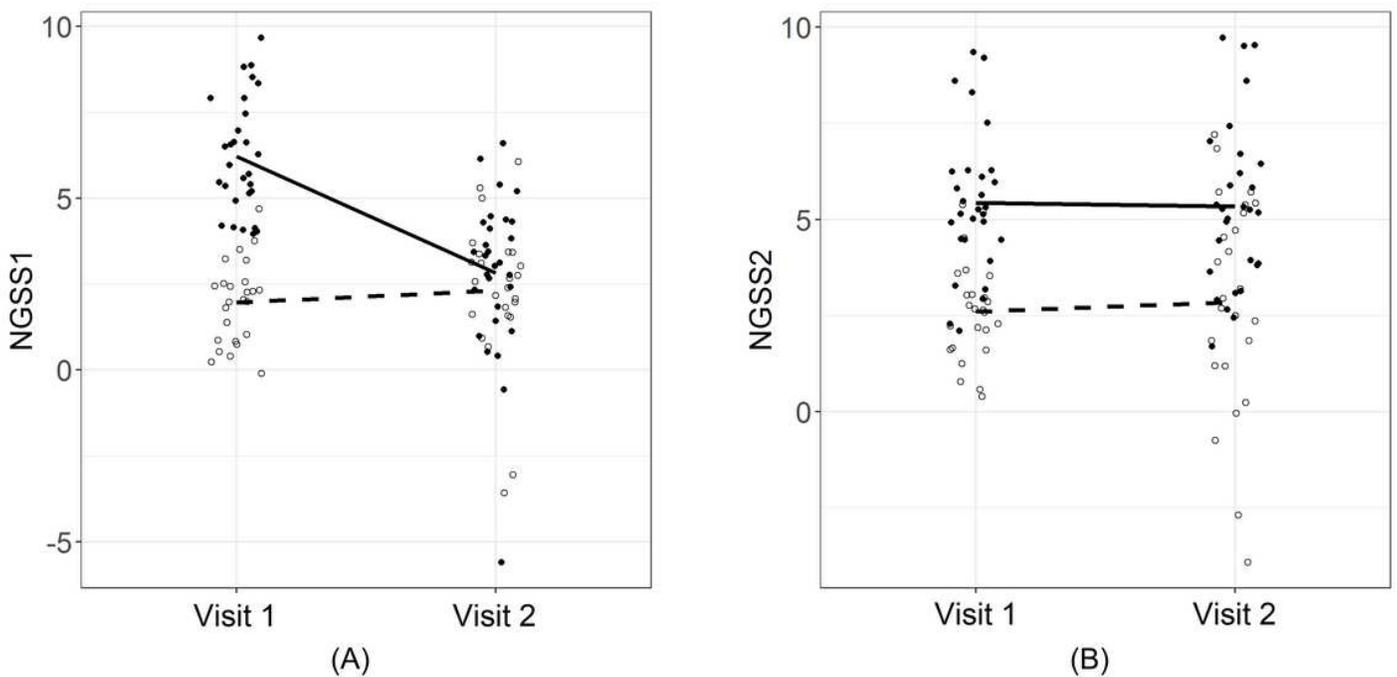


Figure 2

Paired comparisons between visit 1 and visit 2 using NGSS1 (panel (A)) and NGSS2 (panel (B)). A total of $n=54$ subjects with samples in both visits were used. Solid dots represent severe subjects and empty dots represent mild subjects. The solid line represents the mean trend of severe subjects and the broken line represents the mean trend for mild subjects. (A): At visit 1, there was a significant difference in mean NGSS1 between the severe ($n=29$) and mild ($n=25$) groups (6.22 vs. 1.96, $p<0.001$). Mean NGSS1 of the mild group was virtually unchanged between two visits (1.96 vs. 2.31, $p=0.45$). In comparison, mean NGSS1 of the severe group declined significantly at visit 2 (6.22 vs. 2.82, $p<0.001$). (B): In contrast to NGSS1, the differences in NGSS2 was virtually unchanged between the two visits, due to the fact that NGSS2 were built with stable genes.

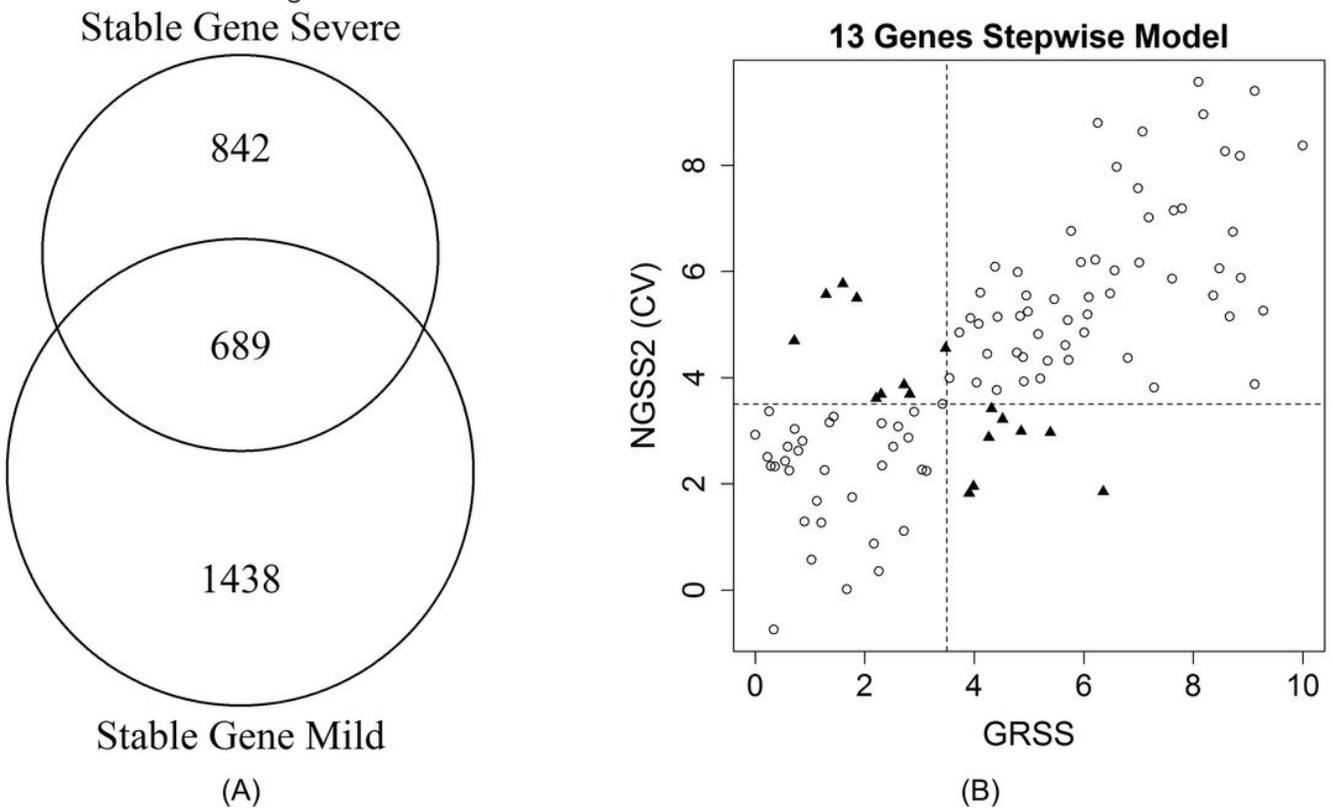


Figure 3

(A). Diagram indicating the stable genes for the mild ($GRSS \leq 3.5$) and severe ($GRSS > 3.5$) groups and the 689 intersecting stable genes common to both groups. (B). Correlating NGSS2 (severity score predicted by Model 4) with GRSS. Naïve Pearson correlation between GRSS and NGSS2 is $r=0.888$. Right: cross-validated Pearson correlation between GRSS and NGSS2 is $r=0.888$. Circles are subjects with correct cross-validated classification based on NGSS2; solid triangles are misclassified subjects.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryMaterial07282020.pdf](#)