

## Validation and Utility of ARDS Subphenotypes Identified by Machine Learning Models Using Clinical Data: An Observational Multi-Cohort Retrospective Analysis.

Manoj V. Maddali MD<sup>1,2</sup>, Matthew Churpek PhD<sup>3</sup>, Tai Pham PhD<sup>4,5</sup>, Emanuele Rezoagli MD<sup>6</sup>, Hanjing Zhuo MBBS<sup>7,8</sup>, Wendi Zhao MHI<sup>2</sup>, June He MBBS<sup>9</sup>, Kevin L Delucchi PhD<sup>10</sup>, Chunxue Wang PhD<sup>11</sup>, Nancy Wickersham BS<sup>11</sup>, J. Brennan McNeil BS<sup>11</sup>, Alejandra Jauregui BS<sup>7,8</sup>, Serena Ke BS<sup>7,8</sup>, Kathryn Vessel BS<sup>7,8</sup>, Antonio Gomez MD<sup>7,12</sup>, Carolyn M. Hendrickson MD<sup>7,12</sup>, Kirsten N. Kangelaris MD<sup>2</sup>, Aartik Sarma MD<sup>7</sup>, Aleksandra Leligdowicz PhD<sup>7,13</sup>, Prof. Kathleen D. Liu PhD<sup>8,10,14</sup>, Prof. Michael A Matthay MD<sup>7,8,15</sup>, Prof. Lorraine B. Ware MD<sup>11,16</sup>, Prof. John G. Laffey MD<sup>17,18</sup>, Prof. Giacomo Bellani PhD<sup>6,19</sup>, Prof. Carolyn S. Calfee MD<sup>7,15</sup>, Pratik Sinha PhD<sup>9,20</sup>, for the LUNG SAFE Investigators and the ESICM Trials Group

1 Division of Pulmonary, Allergy, and Critical Care Medicine, Department of Medicine, Stanford University, Stanford, CA, USA

2 Division of Hospital Medicine, Department of Medicine, University of California San Francisco, San Francisco, CA, USA

3 Division of Allergy, Pulmonary, and Critical Care, Department of Medicine, University of Wisconsin-Madison, Madison, WI, USA

4 Service de Médecine Intensive-Réanimation, AP-HP, Hôpital de Bicêtre, DMU 4 CORREVE Maladies du Cœur et des Vaisseaux, FHU Sepsis, Groupe de Recherche Clinique CARMAS, Le Kremlin-Bicêtre, France

5 Université Paris-Saclay, UVSQ, Univ. Paris-Sud, Inserm U1018, Equipe d'Epidémiologie respiratoire intégrative, CESP, 94807, Villejuif, France

6 Department of Medicine and Surgery, University of Milano-Bicocca, Monza, Italy

7 Division of Pulmonary, Critical Care, Allergy and Sleep Medicine, Department of Medicine, University of California San Francisco, San Francisco, CA, USA

8 Cardiovascular Research Institute, University of California San Francisco, San Francisco, CA, USA

9 Division of Clinical and Translational Research, Washington University School of Medicine, Saint Louis, MO, USA

10 Department of Psychiatry, University of California, San Francisco; San Francisco, CA, USA

11 Division of Allergy, Pulmonary, and Critical Care Medicine, Department of Medicine, Vanderbilt University Medical Center, Nashville, TN, USA

12 Division of Allergy, Pulmonary, and Critical Care Medicine, Department of Medicine, Zuckerberg San Francisco General Hospital and Trauma Center, San Francisco, CA, USA

13 Interdepartmental Division of Critical Care Medicine, University of Toronto, Toronto, Canada

14 Division of Nephrology, Department of Medicine, University of California San Francisco, San Francisco, CA, USA

15 Department of Anesthesia, University of California San Francisco, San Francisco, CA, USA

16 Department of Pathology, Microbiology and Immunology, Vanderbilt University Medical Center, Nashville, TN, USA

17 School of Medicine, Regenerative Medicine Institute at CÚRAM Centre for Research in Medical Devices, National University of Ireland Galway, Galway, Ireland

18 Anaesthesia and Intensive Care Medicine, Galway University Hospitals, Galway, Ireland

19 Department of Anesthesia and Intensive Care Medicine, ASST Monza-Ospedale San Gerardo, Monza, Italy

20 Department of Anesthesia, Division of Critical Care, Washington University, Saint Louis, MO, USA

### Corresponding author:

Pratik Sinha, MB ChB PhD

660 S. Euclid Ave, Campus Box 8054

St. Louis, MO 63110

Email: p.sinha@wustl.edu

Phone: 314-273-3461

**Word count:** 4,313

**Acknowledgements:** We thank Fabiana Madotto, James Anstey, and Nader Najafi for their contributions in data collection, cleaning, and analysis. We thank all patients and researchers who participated in the National Heart Lung and Blood Institute (NHLBI) ARDS Network trials from which data from this study were derived. These include the ALVEOLI, FACTT, and SAILS trials. We acknowledge the contributions of healthcare providers and research staff that enabled the successful completion of these trials. In addition, we thank the contributions of the Biological Specimen and Data Repository Information Coordinating Center of the NHLBI (BIOLINCC) that made the data and biological specimens available to do these studies. This manuscript was prepared using ALVEOLI, ARDSNET, and FACTT Research Materials obtained from the NHLBI Biologic Specimen and Data Repository Information Coordinating Center and does not necessarily reflect the opinions or views of the ALVEOLI, ARDSNET, FACTT or the NHLBI.

57 **Abstract**

58 **Background:** Two acute respiratory distress syndrome (ARDS) subphenotypes with distinct clinical and biological  
59 features and differential treatment responses have been identified using latent class analysis (LCA) in seven  
60 individual cohorts. To facilitate bedside identification of subphenotypes, clinical-classifier models using readily  
61 available clinical variables have been described in five randomized-controlled trials. Performance of these models in  
62 observational cohorts of ARDS is unknown.

63  
64 **Methods:** We evaluated the performance of machine learning clinical-classifier models for assigning ARDS  
65 subphenotypes in two observational cohorts of ARDS: EARLI (n=335) and VALID (n=452), with LCA-derived  
66 subphenotype as the gold-standard. We also assessed model performance in EARLI using data automatically  
67 extracted from the electronic health record (EHR). In LUNG SAFE (n=2813), a multinational observational ARDS  
68 cohort, we applied the model to determine the prognostic value of the subphenotypes and tested their interaction  
69 with PEEP strategy, with mortality as the dependent variable.

70  
71 **Findings:** The clinical-classifier models had an area under receiver operating characteristic curve (AUC) of 0.92  
72 (95% CI: 0.90–0.95) in EARLI and 0.88 (0.84–0.91) in VALID. Model performance was comparable when using  
73 exclusively EHR-derived predictors. In LUNG SAFE, 90-day mortality was higher in the Hyperinflammatory  
74 subphenotype (57% [414/725] vs. 33% [694/2088];  $p < 0.0001$ ). There was a significant treatment interaction with  
75 PEEP strategy and ARDS subphenotype ( $p = 0.041$ ), with lower mortality in the high PEEP group in the  
76 Hyperinflammatory subphenotype, following similar patterns to those observed in prior analyses of the ALVEOLI  
77 trial.

78  
79 **Interpretation:** Classifier models using clinical variables alone can accurately assign ARDS subphenotypes in  
80 observational cohorts. Application of these models can provide valuable prognostic information and may inform  
81 management strategies for personalised treatment, including application of PEEP, once prospectively validated.

82  
83 **Funding:** National Institutes of Health (PS: GM142992, CSC: HL140026, LBW: HL103836, HL135849), European  
84 Society of Intensive Care Medicine.

85 **Research in context**

86

87 *Evidence before this study*

88 Using latent class analysis (LCA), two biological acute respiratory distress syndrome (ARDS) subphenotypes –  
89 termed “Hypoinflammatory” and “Hyperinflammatory”– have been identified, with distinct clinical and biological  
90 features, outcomes, and differential responses to therapy. Clinical use of these subphenotypes, however, is limited  
91 by complexity and lack of point of care biomarker assays. To facilitate bedside identification of these  
92 subphenotypes, machine learning classifier models using only readily available clinical variables have been  
93 developed and validated using data from randomized controlled trials. Performance and clinical utility of these  
94 models in observational cohorts of ARDS is not known. No formal literature search was done for this study.

95

96 *Added value of this study*

97 The presented study demonstrates the validity of machine learning clinical-classifier models in accurately  
98 identifying ARDS subphenotypes in two observational cohorts. Differences in biomarkers and clinical outcomes in  
99 subphenotypes identified using these models were similar to those in LCA-derived subphenotypes. The models  
100 performed comparably when utilizing a dataset comprised of variables automatically extracted from the electronic  
101 health record (EHR), suggesting that EHR-embedded models may be feasible. When applied to a large multinational  
102 observational cohort of ARDS, the models identified patients at risk for adverse clinical outcomes. The models also  
103 identified a treatment interaction with PEEP and subphenotype, with lower mortality observed with higher PEEP in  
104 the Hyperinflammatory subphenotype, similar to patterns observed in secondary analyses of the ALVEOLI trial.

105

106 *Implications of all the available evidence*

107 LCA-derived phenotyping has recently shown promise in identifying homogenous subgroups within larger,  
108 heterogenous populations of ARDS. Clinical-classifier models using readily available clinical data can accurately  
109 identify these subphenotypes at the bedside and could facilitate prospective, subphenotype-specific trials in ARDS.  
110 Response to PEEP may differ on the basis of subphenotype.

111

112 **Introduction**

113 The acute respiratory distress syndrome (ARDS) remains a highly prevalent cause of acute respiratory failure,  
114 resulting in high morbidity and mortality.<sup>1</sup> Yet, potentially as a consequence of underlying heterogeneity, few  
115 therapeutic options have proven to be beneficial in randomized controlled trials (RCTs).<sup>2-4</sup> Two discrete biological  
116 subphenotypes have been identified using latent class analysis (LCA) in five RCTs and two observational cohorts,  
117 totaling over 4,000 patients.<sup>5-8</sup> The two subphenotypes have distinct clinical and biological features, divergent  
118 outcomes, and in three RCTs, differential treatment responses were observed.

119  
120 Although accurate parsimonious models for subphenotype identification have been developed, these models are  
121 reliant on measurement of protein biomarkers (e.g., interleukin (IL)-6, IL-8, soluble tumor necrosis factor receptor  
122 (sTNFR)-1, Protein C).<sup>9</sup> The limited availability of real-time assays for these biomarkers represents a barrier to the  
123 clinical implementation and rapid identification of the subphenotypes.<sup>5,10</sup> Recently, machine learning classification  
124 algorithms utilizing routinely available clinical variables have shown promise in identifying LCA-derived  
125 subphenotypes in RCT cohorts of ARDS.<sup>10</sup> Their performance in unselected populations of ARDS patients, where  
126 patient heterogeneity may be even greater and where comparatively higher mortality is observed, is unknown.<sup>11</sup> A  
127 critical step towards clinical application of these models is their validation in observational and representative  
128 populations of ARDS patients, particularly since it is these unselected, “real-world” patients in whom the models  
129 would be used to screen for enrollment in future RCTs.

130  
131 The primary objective of this study was to validate machine learning classifier models that use readily available  
132 clinical data in observational cohorts of ARDS. Secondary objectives were (1) to evaluate model performance in an  
133 observational cohort of patients with predictor-variables automatically extracted from the electronic health record  
134 (EHR) and (2) to evaluate the clinical utility for prognostication and seeking differential responses to positive end-  
135 expiratory pressure (PEEP) strategy of ARDS subphenotypes derived using the clinical-classifier models in a large  
136 multinational observational cohort of ARDS.

137  
138 **Methods**  
139 *Study populations*  
140 Details of the RCT cohorts used for model development are described in prior studies.<sup>10,12-15</sup> Two observational  
141 cohorts of ARDS, Early Assessment of Renal and Lung Injury (EARLI, n=335) and Validating Acute Lung Injury  
142 markers for Diagnosis (VALID, n=452), served as independent validation cohorts for the models. EARLI is an  
143 ongoing prospectively enrolled cohort of patients admitted to UCSF Medical Center and Zuckerberg San Francisco  
144 General Hospital Intensive Care Units (ICUs). Study participants were identified in the Emergency Department  
145 upon request for admission to the ICU. For this analysis, patients were selected from EARLI if they were deemed to  
146 have ARDS as defined by the American-European Consensus Conference (AECC) criteria on either Day 1 or 2 of  
147 the study and included patients recruited between 2008–2018.<sup>16</sup> Details of the study protocol have been previously  
148 published.<sup>17</sup> VALID is an ongoing prospectively enrolled cohort of patients admitted to Vanderbilt University

149 Medical Center ICU; details of the study protocol have been previously published.<sup>18</sup> Study participants were  
150 enrolled in the study on the morning of the second day of admission to a medical, surgical, trauma, or cardiovascular  
151 ICU. Patients were selected from VALID for inclusion in this analysis if they were deemed to have ARDS as  
152 defined by AECC criteria on the first or second day of ICU admission and included patients were recruited between  
153 2008–2016. Patients with trauma-related ARDS were excluded given biological and clinical differences (e.g., lower  
154 burden of inflammation and lower age-adjusted mortality) from patients with non-trauma ARDS and our previous  
155 work suggesting the subphenotyping schema is most valid in patients with non-trauma ARDS.<sup>8,19</sup> The AECC  
156 definition was used because enrollment in both cohorts started prior to development of the Berlin definition and  
157 because patients continued to be enrolled using both definitions.<sup>20</sup> The described strategy allowed capture of more  
158 patients for analysis. Further, the LCA-derived subphenotypes have been validated in these cohorts using the AECC  
159 definition with similar subphenotypes identified as when using the Berlin definition.<sup>8</sup> Both cohorts include  
160 comprehensive demographic, clinical, and biomarker data from the day (or day prior to) of ARDS diagnosis that  
161 were manually collected by trained research coordinators, as well as clinical outcome data including ventilator-free  
162 days (VFD) and hospital mortality.

163  
164 The Large Observational Study to Understand the Global Impact of Severe Acute Respiratory Failure (LUNG  
165 SAFE, n=2813) was a large, multinational, multicenter, prospectively enrolled cohort of patients admitted to 459  
166 ICUs across 50 countries from February to March 2014; details of the study protocol have been previously  
167 published.<sup>1</sup> Study participants were enrolled on the first day that acute hypoxemic respiratory failure criteria were  
168 satisfied. Patients were selected from LUNG SAFE for inclusion in this analysis if they were deemed to have ARDS  
169 as defined by Berlin criteria on the first or second day of study enrollment.<sup>20</sup> The LUNG SAFE study was conducted  
170 after the description of the Berlin definition hence its use in this cohort. This cohort includes demographic, clinical,  
171 and respiratory data from the day patients were enrolled into the study and at pre-specified intervals until ICU  
172 discharge or death (see Supplement). All study cohorts were approved by the Institutional Review Board at each  
173 participating hospital.

#### 174 175 *Model Development and Validation*

176 All models were trained to predict the Hyperinflammatory phenotype. Of the machine learning clinical-classifier  
177 models described in the original study,<sup>10</sup> we used the two best performing models to validate in this study, with a  
178 parsimonious (“vitals and labs”) model comprising only of vital signs and laboratory values serving as the primary  
179 model. As the secondary model, we used a “full feature” model comprising all the predictors in the primary model,  
180 with the addition of ventilatory variables and demographics. The “vitals and labs” model served as the primary  
181 model because it was less complex (fewer predictors), constituted exclusively of physiological predictors, was one  
182 of the most accurate in the original study, and was the most generalizable model.

183  
184 In both EARLI and VALID, due to missing predictors, the original “vitals and labs” and “full feature” models could  
185 not be validated directly. In EARLI, the “vitals and labs” model had no predictors missing and the “full feature”

186 model had one predictor missing (minute ventilation). In VALID, there was one predictor missing for the “vitals and  
187 labs” model (glucose) and three predictors missing (tidal volume, glucose, and body mass index) for the “full  
188 feature” model. To simplify the analysis, we developed new “vitals and labs” and “full feature” models comprising  
189 of common predictors available for each model in both EARLI and VALID from the original predictors. A final list  
190 of variables used in each model is described in **Table S1**.

191  
192 A schematic of the analysis plan is presented in **Figure 1**. For model development, we used a gradient-boosted  
193 machine algorithm, XGBoost: Extreme Gradient Boosting (version 1.3.2.1). In brief, gradient-boosted machines  
194 utilize an ensemble of multiple decision trees, where trees added sequentially to the model to attempt to correct the  
195 classification error of previous trees in the ensemble. We utilized 10-fold cross validation and hyperparameter  
196 tuning using a grid search to tune and optimize the models using the training set (see Supplement), recapitulating our  
197 prior approach.<sup>10</sup> All models were developed using a training set comprised of a combination of three RCT cohorts,  
198 ARMA, ALVEOLI, and FACTT (n=2022), and model performance was tested externally in SAILS (n=745). The  
199 models output a continuous probability specifying the likelihood of classification to the Hyperinflammatory  
200 subphenotype for each patient. To evaluate the validity of these two new models in relation to the models developed  
201 in the original study,<sup>10</sup> we compared the probabilities generated by corresponding new and original models using  
202 Pearson’s correlation coefficient.

203  
204 Next, performance of these models was evaluated independently in EARLI (n=335) and VALID (n=452). Validation  
205 cohorts were kept isolated from the training and testing procedures. LCA-derived subphenotypes served as the  
206 reference standard for model training, testing, and validation. The procedure for handling missing data is detailed in  
207 the Supplement and the missingness for predictors in each validation cohort are presented in **Table S2**. Overall  
208 model performance in EARLI and VALID was evaluated by (a) calculating the area under the receiver operating  
209 characteristic curve (AUC) with confidence intervals (estimated using 2000 stratified bootstrap replicates); and (b)  
210 generating calibration plots. For each model, class was assigned using a probability cutoff of 0.5 to report on  
211 accuracy, sensitivity, and specificity of subphenotype assignments. As with our prior work, we additionally  
212 performed sensitivity analysis using cutoffs of 0.3, 0.4, 0.6, and 0.7. Once patients were assigned subphenotypes,  
213 we evaluated differences in protein biomarkers and clinical outcomes (e.g., mortality and ventilator-free days).

#### 214 215 *Model validation in EHR-derived cohort*

216 Patients in the EARLI cohort were identified in the UCSF’s electronic health record, Epic (Epic Systems Corp.).  
217 Patients enrolled before implementation of the UCSF EHR in 2012 were excluded, as were patients admitted at San  
218 Francisco General Hospital, due to the Epic EHR being implemented at this institution after the study period (post  
219 2018). All vital signs and laboratory values from each participant’s admission encounter were queried using SQL  
220 and downloaded from Epic Clarity, a data warehouse and relational database that stores the majority of clinical data  
221 within Epic. Additionally, we queried usage of intravenous vasoactive agents and incorporated this into the cohort as  
222 a binary variable. The most extreme values (e.g., highest heart rate or lowest serum bicarbonate level) observed  $\pm$  12

223 hours of ARDS diagnosis were extracted. We used this “EHR-derived EARLI cohort” to identify ARDS  
224 subphenotypes using the “vitals and labs” model. Model performance was evaluated using the same procedures  
225 described above with LCA-derived subphenotype serving as the reference standard. For comparison, we evaluated  
226 model performance for the same patients identified in the EHR cohort but using vital signs and labs collected  
227 manually during the original EARLI prospective study enrollment that were used or the original LCA.

228

#### 229 *Model evaluation in LUNG SAFE*

230 In LUNG SAFE, due to limited data collection, only a small selection of predictor variables was available for  
231 modelling (**Table S1**). A custom clinical-classifier model, comprising only these variables, was developed using the  
232 same procedure described above (training: ARMA, ALVEOLI, and FACTT, testing: SAILS). As LCA-derived  
233 subphenotypes were not known in LUNG SAFE and the model contained a sparse set of predictor variables, a priori,  
234 we first sought the best probability cutoff to assign class in VALID (an observational cohort) to optimize  
235 classification accuracy. Optimal cutoff in VALID was determined based on the tradeoff between sensitivity and  
236 specificity (i.e., Youden index).<sup>21</sup> EARLI was not used to determine cutoff due to some of the LUNG SAFE  
237 variables not being available. For sensitivity analysis, we additionally evaluated model results in LUNG SAFE  
238 across a range of probability cutoffs.

239

240 Once LUNG SAFE patients were classified into subphenotypes, we compared clinical outcomes, resolution of  
241 ARDS, prevalence of underlying chronic diseases, and ventilatory/respiratory variables stratified by model-assigned  
242 subphenotype. Building on prior work showing differential subphenotype responses to PEEP, we evaluated the  
243 interaction between subphenotype allocation and PEEP in LUNG SAFE.<sup>5</sup> In order to create two groups with  
244 substantially different levels of PEEP usage, patients were classified into tertiles according to their mean PEEP over  
245 days 1–3. The top tertile was labelled as “high PEEP” and bottom tertile as “low PEEP,” with the middle tertile  
246 excluded from analysis. A logistic regression model was created with the interaction term of PEEP-group and  
247 subphenotype as an independent variable and 90-day mortality as the dependent variable. As sensitivity analyses, we  
248 tested for differences in mortality, VFDs, and PEEP treatment interaction for a range of probability cutoffs. We also  
249 tested for subphenotype treatment interaction with PEEP-groups derived when using quintiles instead of tertiles.  
250 Further sensitivity analyses included testing treatment interaction of PEEP-groups with subgroups of ARDS severity  
251 as stratified by (a) PaO<sub>2</sub>/FiO<sub>2</sub> and (b) Sequential Organ Failure Assessment (SOFA) score (see Supplement for  
252 details).

253

#### 254 *Statistical Analysis*

255 Differences in outcomes between subphenotypes were tested using Pearson’s chi-squared test. Between-group  
256 differences were tested using Student’s t-test and Wilcoxon rank-sum test, depending on variable distribution. For  
257 differences in outcomes between ARDS subphenotypes, we also computed odds ratios for mortality and rank  
258 biserial correlations for VFDs. The Wald test was used to test for significance of the interaction term in the logistic  
259 regression models. All analyses were done using R (version 4.03) and RStudio interface (version 1.4.1106). The

260 codes used for analysis can be found on our group’s GitHub page, available at [https://github.com/Calfee-Sinha-](https://github.com/Calfee-Sinha-PrecisionCriticalCareLab)  
261 PrecisionCriticalCareLab.

262

### 263 *Role of funders*

264 The funder of the study had no role in study design, data collection, data analysis, data interpretation,  
265 or writing of the report.

266

## 267 **Results**

268 Baseline patient characteristics for the training set, both validation cohorts (EARLI and VALID), and LUNG SAFE  
269 are shown in **Table S3**.

270

271 The top ten most important features for the “vital and labs” and “full feature” models in the training dataset are  
272 shown in **Figures S1A** and **S1B** respectively and in line with prior models.<sup>10</sup> In SAILS, the probabilities of  
273 subphenotype assignment generated by the new models developed for use in this study were highly correlated to  
274 probabilities generated by our previously described models: “vital and labs” model ( $r = 0.97$ ,  $p < 0.0001$ , **Figure**  
275 **S2A**) and “full feature” model ( $r = 0.95$ ,  $p < 0.0001$ ; **Figure S2B**).<sup>10</sup>

276

### 277 *Model evaluation in observational cohorts*

278 The “vitals and labs model” had an AUC of 0.92 (95% CI: 0.90 – 0.95) in EARLI and 0.88 (95% CI: 0.84 – 0.91)  
279 in VALID (**Figure 2**). Model sensitivity, specificity, and accuracy when using a probability cut-off of 0.5 are  
280 reported in **Table 1**, and over a range of probability cut-offs in **Table S4**. The calibration plot for the model in both  
281 cohorts is presented in **Figure S3A** and **S3B**.

282

283 In both cohorts, the Hyperinflammatory subphenotype identified by “vitals and labs” models had significantly  
284 higher levels of interleukin-6 (IL-6), interleukin-8 (IL-8), and soluble TNF receptor-1 (sTNFR1), and lower levels of  
285 Protein C (**Figure 3**). The Hyperinflammatory phenotype was associated with higher in-hospital mortality and fewer  
286 ventilator-free days (**Table 2**). Clinical outcomes for the subphenotypes in both cohorts over a range of probability  
287 cutoffs were similar to those using a cutoff  $\geq 0.5$  (**Table S5**).

288

289 In both EARLI and VALID, the “full feature” model had similar model performance metrics (**Figure S4**, **Tables 1**  
290 and **S4**) and differences in biomarkers (**Figure S5**) and clinical outcomes (**Table S6**) between the subphenotypes as  
291 the “vitals and labs” model.

292

### 293 *Model validation in EHR-derived cohort*

294 117 patients from the EARLI cohort were identified in the UCSF EHR. Baseline patient characteristics along with  
295 feature missingness are shown in **Table S7**. The “vitals and labs” model using EHR-derived data had an AUC of  
296 0.88 (95% CI: 0.81 – 0.94) compared to an AUC of 0.92 (95% CI: 0.88 – 0.97; **Figure S6**) using hand-curated



297 variables for the same patients. Clinical outcomes in subphenotypes assigned using EHR-derived data were similar  
298 to those derived using hand-curated data (**Table S8**).

299

#### 300 *Clinical-classifier model in LUNG SAFE*

301 When first evaluated in SAILS and VALID, the LUNG SAFE classifier model resulted in an AUC of 0.93 (0.91 –  
302 0.95) and 0.87 (0.83 – 0.90) respectively. In VALID, the model had the highest Youden index at a probability  
303 cutoff of 0.4 (**Table S9**). This probability cutoff was used to classify subphenotypes in LUNG SAFE.

304

305 Using a cutoff of 0.4, 26% (725/2813) of patients in LUNG SAFE were classified in the Hyperinflammatory  
306 subphenotype. Mortality at day 90 in the Hyperinflammatory subphenotype was 57% (414/725) compared to 33%  
307 (694/2088) in the Hypoinflammatory group ( $p < 0.0001$ ). VFDs were significantly fewer in the Hyperinflammatory  
308 subphenotype ( $p < 0.0001$ ; **Table 2**). Survival between groups diverged at day 1 that was sustained over 90 days,  
309 with a significantly lower survival in the Hyperinflammatory group (**Figure 4**). The observed differences in  
310 mortality and VFDs were consistent across a range of probability cutoffs (**Table S10**).

311

312 More patients in the Hypoinflammatory subphenotype had resolution of ARDS on day 2 (35%; 510/1447) compared  
313 to the Hyperinflammatory subphenotype (28%; 129/469;  $p = 0.0024$ ), suggesting temporal stability of ARDS  
314 diagnosis in the latter. Prevalence of underlying chronic liver disease was significantly higher in the  
315 Hyperinflammatory subphenotype, whereas prevalence of chronic obstructive pulmonary disease (COPD) was lower  
316 (**Table S11**). Difference in respiratory variables, even among those that were statistically significant, were not  
317 clinically significant between the two subphenotypes (**Figure 5**).

318

319 When stratified into tertiles based on mean day 1 to 3 PEEP, median PEEP in the top tertile (“high PEEP”;  $n=992$ )  
320 was 11 cm H<sub>2</sub>O (10 – 12) and bottom tertile (“low PEEP”;  $n=943$ ) was 5 cm H<sub>2</sub>O (5 – 6). Differences between the  
321 characteristics of the low- and high-PEEP groups can be found in **Table S12**. There was a significant interaction  
322 between PEEP subgroups and ARDS subphenotypes with 90-day mortality as the outcome; Hyperinflammatory  
323 subphenotype: “high PEEP” 54% [169/313] vs. “Low PEEP” 62% [127/205]; Hypoinflammatory subphenotype:  
324 “high PEEP” 34% [231/675] vs. “Low PEEP” 32% [233/734] ( $p = 0.041$ ; **Table 3**). The differences in outcomes  
325 and treatment interaction were significant across a range of probability cutoffs (**Table S13**). The interaction term  
326 remained significant after adjusting the model for age and PaO<sub>2</sub>/FiO<sub>2</sub> ( $p = 0.047$ ). A sensitivity analysis using  
327 quintiles to define PEEP groups (with the middle quintile eliminated) also revealed significant treatment interactions  
328 (**Table S14**). Significant interactions with PEEP groups were not observed when the population was stratified by  
329 other measures of disease severity such as PaO<sub>2</sub>/FiO<sub>2</sub> ( $p = 0.96$ ) or SOFA score ( $p = 0.30$ ; **Table S15 and S16**).

330

#### 331 **Discussion**

332 In this study, we report that machine learning classifier models, using only readily available clinical variables as  
333 predictors, can accurately assign ARDS subphenotypes in observational cohorts. Our models consistently captured

334 the rich biological information that define the LCA-derived subphenotyping schema, with marked differences in  
335 protein biomarkers between the two identified phenotypes. The models identified patients at high risk for adverse  
336 outcomes, including in the large multinational observational cohort (LUNG SAFE), where protein biomarker data  
337 were not available. Further, in LUNG SAFE, we observed differential responses to PEEP strategy by subphenotype,  
338 with higher PEEP associated with improved outcomes in the Hyperinflammatory subphenotype, similar to patterns  
339 previously identified in secondary analyses of the ALVEOLI trial.<sup>5</sup> Finally, the “vitals and labs” model performed  
340 robustly even when utilizing clinical data extracted automatically from the EHR (as opposed to values obtained  
341 manually during study enrollment). Taken together, the models presented in these studies represent a substantial step  
342 towards translating ARDS subphenotypes into the clinical workflow. Pending prospective evaluation, these models  
343 may be valuable tools for prognostication and treatment stratification in future trials.

344  
345 The utility of subphenotypes in ARDS is contingent on feasible bedside identification. Although point-of-care and  
346 real-time assays are being developed rapidly, they remain experimental.<sup>22</sup> In the interim, or as an alternative,  
347 clinical-classifier models can be a useful adjunct. Clinical-classifier models to date have been validated only in  
348 retrospective secondary analyses of relatively uniform RCTs,<sup>10</sup> which enroll typically only 5–10% of potentially  
349 eligible patients,<sup>11</sup> thus limiting their routine application. By contrast, validation of these models in observational  
350 cohorts of all-comer patients with ARDS indicates that such classification algorithms can be reliably applied to more  
351 generalizable populations and could potentially be used to screen patients for eligibility for enriched RCTs.  
352 Embedding such models into the EHR would allow for bedside screening for and enrollment into prospective  
353 clinical trials to evaluate for prognostic or therapeutic differences among patients with ARDS. Moreover, such  
354 models could more easily capture temporal trends given the rich, abundant data stream in the ICU. By demonstrating  
355 the high-performance metrics of the models with EHR-derived data, our study serves as a proof of concept that  
356 EHR-embedded machine learning models are feasible for classifying ARDS subphenotypes. If validated  
357 prospectively, such EHR-embedded models could provide on-demand decision support for clinicians and/or clinical  
358 trials, while limiting disruption to clinical workflow by automatically incorporating clinical data into the models.

359  
360 The implementation of these models in the clinical setting are, however, contingent on two factors. First, it must be  
361 prospectively demonstrated that the models can classify phenotypes robustly and consistently in real-time clinical  
362 scenarios in diverse settings. Prior to their clinical implementation, the models will need rigorous evaluation for  
363 their interaction with missing data frequently encountered in the real world setting of critical care. Second, it is  
364 imperative that a clear clinical utility of the subphenotypes is demonstrated prior to their EHR implementation.  
365 Based on its performance in our study, we would advocate the use of the “vital and labs” model for prospective  
366 evaluation in future studies. Interestingly, the clinical utility and divergent characteristics of the subphenotypes  
367 identified using the sparse model in LUNG SAFE would suggest that a model comprising of even fewer features  
368 than the “vitals and lab model” may classify with sufficient accuracy. The development and validation of such  
369 parsimonious models requires careful evaluation using the most important variables identified in the “vitals and  
370 labs” model, rather than sets of variables constrained by availability, such as in the LUNG SAFE model.

371  
372 Though model performance was comparable between both observational cohorts, model performance in EARLI was  
373 marginally better compared to VALID, potentially due a variety of factors including the timing of enrollment into  
374 the studies. Patients were enrolled on the day of hospital admission in EARLI, whereas in VALID enrollment was  
375 on day two of ICU admission. Earlier study enrollment may have captured the most extreme physiological  
376 characteristics for each patient and higher classification into the Hyperinflammatory subphenotype, but without  
377 serial protein biomarker quantification and LCA classification, the temporal kinetics of the subphenotypes remain a  
378 key knowledge gap in the field. The longitudinal model performance metrics of the clinical-classifier model requires  
379 further studies.

380  
381 Our findings in LUNG SAFE are consistent with prior studies suggesting that ARDS subphenotypes capture unique  
382 information compared to other metrics of ARDS severity, such as PaO<sub>2</sub>/FiO<sub>2</sub> or SOFA score.<sup>5,6,10</sup> In our analysis, a  
383 treatment interaction was observed between PEEP groups and the subphenotypes with differential responses.  
384 Notably, this treatment interaction was consistent with our previous secondary analysis of the ALVEOLI trial that  
385 tested the efficacy of high PEEP versus low PEEP in ARDS.<sup>5,13</sup> In that analysis, as in this study, high PEEP was  
386 associated with improved survival in the Hyperinflammatory subphenotype and worse survival in the  
387 Hypoinflammatory subphenotype, albeit the effect size in the latter was clinically insignificant in both studies. The  
388 consistent findings across both these studies suggest that there may be value in evaluating PEEP strategies more  
389 formally in subphenotype-specific trials with treatment directed by subphenotype. Specifically, in future trials  
390 testing high-PEEP strategies, inclusion of the Hypoinflammatory phenotype may lead to a dilution of the effective  
391 sample size, rendering the detection of a significant effect less likely.

392  
393 This study has several strengths. First, the models performed comparably across two observational cohorts with  
394 variable inclusion criteria, suggesting model generalizability. Second, the models were able to identify high-risk  
395 patients when utilizing inputs automatically extracted from the EHR, showing that biological “signal” can be  
396 accurately captured despite the “noise” associated with EHR-derived data. Third, the primary model performed well  
397 despite utilizing a parsimonious set of features (only vital signs and laboratory values). This approach could allow  
398 future EHR-embedded models to use the “most objective” inputs while excluding features that are dependent on  
399 epidemiological factors (e.g., race/ethnicity), or those which are harder to capture in the EHR (e.g., ARDS risk  
400 factors and ventilatory variables). The “vitals and labs” model also has the advantage of being potentially applicable  
401 in low- and middle- income countries where availability of emerging point-of-care protein biomarker quantification  
402 may not be feasible.<sup>22</sup> Fourth, unlike our prior studies, this is the first time we have tested the performance of the  
403 clinical classifier models and shown the clinical value of ARDS subphenotypes in a cohort consisting of patients  
404 derived from low and middle income countries, suggesting their generalizability across healthcare systems.

405  
406 This study also has several limitations. Due to a lack of availability of predictor variables, we were not able to  
407 validate the exact models developed in our prior study.<sup>10</sup> The strong correlation of the probabilities generated by the

408 models we presented in this study compared to models in our prior study would, however, suggest that these models  
409 are highly overlapping. It is noteworthy that the “full feature” model was trained with the race variable stratified as  
410 white and non-white, thereby limiting its generalisability and validity in populations with greater racial or ethnic  
411 diversity. However, the “vitals and labs” model without this data also performed well. The EHR-derived cohort was  
412 limited by a relatively small sample size and high missingness for some variables. In addition to limiting model  
413 validity, the observed missingness, specifically in the EHR cohort, highlights some of the challenges in applying  
414 such models prospectively and embedded in the EHR. In LUNG SAFE, fewer features were available, and the  
415 tolerance of these models for variable missingness or predictor variable parsimony requires further evaluation. There  
416 were several differences in the clinical baseline characteristics of EARLI, VALID, and LUNG SAFE. Most notably,  
417 ARDS risk factors, PaO<sub>2</sub>/FiO<sub>2</sub>, and bicarbonate levels were substantially different in LUNG SAFE, and taken  
418 together with the lack of a comparative gold-standard (LCA-derived subphenotype) to evaluate model performance,  
419 the findings of this portion of the study should be interpreted cautiously. Further, interpretation of the findings of  
420 treatment interaction with PEEP groups and subphenotypes should also be cautious given that these data are  
421 generated from observational data and level of PEEP was not randomly assigned; however, their concordance with  
422 our previous findings from randomized PEEP trials is noteworthy and suggests validity. Finally, to date, application  
423 of these models has been retrospective, and their validity in real-time clinical settings remains to be tested.

424

425 In summary, machine learning classifier models using readily available clinical data accurately assigned  
426 inflammatory subphenotypes in observational populations of ARDS. Additionally, the models performed robustly in  
427 an EHR-derived observational cohort, suggesting such models can be potentially embedded into an EHR. Finally,  
428 the model identified high-risk patients and a treatment interaction between PEEP and inflammatory subphenotype in  
429 a large observational cohort without a reference standard of LCA-derived classification, providing further support of  
430 the hypothesis that the effect of PEEP may differ in each subphenotype. Application of these models to identify  
431 subphenotypes can provide valuable prognostic information linked to distinct biological characteristics and may  
432 inform management strategies to test in future clinical trials.

433

434 **Author contributions:** MVM, PS, CSC, LBW, MAM, JGL, and GB were responsible for study conception and  
435 design. MVM, TP, PS, WZ, JH, KLD, YC, HZ, CW, NW, JBM, LBW, ER and CSC were responsible for the data  
436 cleaning and analysis. MVM, JH, and PS were responsible for data verification. All authors were responsible for  
437 data collection and/or clinical adjudication. MVM, PS, MC, CSC, LBW, JGL, and GB developed the first draft of  
438 the manuscript. All authors reviewed and edited the final version of the manuscript.

439

440 **Data Sharing:** Data from these studies can be provided to others upon reasonable request on approval of a written  
441 request to Dr Pratik Sinha. Data from the National Heart Lung and Blood Institute were accessed through the  
442 BIOLINCC public repository.

443

444 **Declaration of interests:** Dr. Churpek reports grants from NIH/NIDA (R01 DA051464), grants from  
445 DOD/PRMRP, W81XWH-21-1-0009, grants from NIH/NIA (R21 AG068720), grants from NIH/NIGMS (R01  
446 GM123193), grants from NIH/ NIDDK (R01 DK126933), grants from EarlySense (Tel Aviv, Isreal), grants from  
447 NIH/NHLBI (R01 HL157262) outside the submitted work. In addition, Dr. Churpek has a patent Patent pending  
448 (ARCD. P0535US.P2) pending to University of Chicago related to clinical deterioration risk prediction algorithms  
449 for hospitalized patients. Dr. Sarma reports grants from National Heart, Lung, and Blood Institute during the

450 conduct of the study. Dr. Matthay reports grants from Roche-Genentec, personal fees from Johnson and Johnson,  
451 personal fees from Novartis Pharmaceuticals, personal fees from Gilead Pharmaceuticals, and personal fees from  
452 Pliant Therapeutics, outside the submitted work. Dr. Ware reports grants from National Institutes of Health (US),  
453 during the conduct of the study; grants and personal fees from Boehringer Ingelheim, grants from Genentech, grants  
454 from CSL Behring, personal fees from Merck, personal fees from Citius, personal fees from Quark, and personal  
455 fees from Foresee, outside the submitted work. Dr. Laffey reports grants from European Society of Intensive Care  
456 Medicine, during the conduct of the study; personal fees from Glaxosmithkline, and personal fees from Baxter,  
457 outside the submitted work. Dr. Bellani reports grants and personal fees from Draeger Medical, personal fees from  
458 Ge Healthcare, personal fees from Hamilton Medical, and personal fees from Flowmeter SPA, outside the submitted  
459 work. Dr. Calfee reports grants from NIH, during the conduct of the study; grants and personal fees from  
460 Roche/Genentech, grants and personal fees from Bayer, personal fees from Quark Pharmaceuticals, personal fees  
461 from Gen1e Life Sciences, personal fees from Vasomune, and grants from Quantum Leap Healthcare Collaborative,  
462 outside the submitted work. The other authors report no disclosures.  
463

464 **References**

- 465  
466 1. Bellani G, Laffey JG, Pham T, et al. Epidemiology, Patterns of Care, and Mortality for Patients With Acute  
467 Respiratory Distress Syndrome in Intensive Care Units in 50 Countries. *JAMA* 2016; **315**(8): 788-800.  
468 2. Thompson BT, Chambers RC, Liu KD. Acute Respiratory Distress Syndrome. *N Engl J Med* 2017; **377**(6):  
469 562-72.  
470 3. Wilson JG, Calfee CS. ARDS Subphenotypes: Understanding a Heterogeneous Syndrome. *Crit Care* 2020;  
471 **24**(1): 102.  
472 4. Matthay MA, Arabi YM, Siegel ER, et al. Phenotypes and personalized medicine in the acute respiratory  
473 distress syndrome. *Intensive Care Med* 2020; **46**(12): 2136-52.  
474 5. Calfee CS, Delucchi K, Parsons PE, et al. Subphenotypes in acute respiratory distress syndrome: latent  
475 class analysis of data from two randomised controlled trials. *Lancet Respir Med* 2014; **2**(8): 611-20.  
476 6. Calfee CS, Delucchi KL, Sinha P, et al. Acute respiratory distress syndrome subphenotypes and differential  
477 response to simvastatin: secondary analysis of a randomised controlled trial. *Lancet Respir Med* 2018; **6**(9): 691-8.  
478 7. Famous KR, Delucchi K, Ware LB, et al. Acute Respiratory Distress Syndrome Subphenotypes Respond  
479 Differently to Randomized Fluid Management Strategy. *Am J Respir Crit Care Med* 2017; **195**(3): 331-8.  
480 8. Sinha P, Delucchi KL, Chen Y, et al. Latent class analysis-derived subphenotypes are generalisable to  
481 observational cohorts of acute respiratory distress syndrome: a prospective study. *Thorax* 2021: thoraxjnl-2021-  
482 217158.  
483 9. Sinha P, Delucchi KL, McAuley DF, O'Kane CM, Matthay MA, Calfee CS. Development and validation of  
484 parsimonious algorithms to classify acute respiratory distress syndrome phenotypes: a secondary analysis of  
485 randomised controlled trials. *Lancet Respir Med* 2020; **8**(3): 247-57.  
486 10. Sinha P, Churpek MM, Calfee CS. Machine Learning Classifier Models Can Identify Acute Respiratory  
487 Distress Syndrome Phenotypes Using Readily Available Clinical Data. *Am J Respir Crit Care Med* 2020; **202**(7):  
488 996-1004.  
489 11. Pais FM, Sinha P, Liu KD, Matthay MA. Influence of Clinical Factors and Exclusion Criteria on Mortality  
490 in ARDS Observational Studies and Randomized Controlled Trials. *Respir Care* 2018; **63**(8): 1060-9.  
491 12. Acute Respiratory Distress Syndrome N, Brower RG, Matthay MA, et al. Ventilation with lower tidal  
492 volumes as compared with traditional tidal volumes for acute lung injury and the acute respiratory distress  
493 syndrome. *N Engl J Med* 2000; **342**(18): 1301-8.  
494 13. Brower RG, Lanken PN, MacIntyre N, et al. Higher versus lower positive end-expiratory pressures in  
495 patients with the acute respiratory distress syndrome. *N Engl J Med* 2004; **351**(4): 327-36.  
496 14. National Heart L, Blood Institute Acute Respiratory Distress Syndrome Clinical Trials N, Wiedemann HP,  
497 et al. Comparison of two fluid-management strategies in acute lung injury. *N Engl J Med* 2006; **354**(24): 2564-75.  
498 15. National Heart L, Blood Institute ACTN, Truwit JD, et al. Rosuvastatin for sepsis-associated acute  
499 respiratory distress syndrome. *N Engl J Med* 2014; **370**(23): 2191-200.  
500 16. Bernard GR, Artigas A, Brigham KL, et al. The American-European Consensus Conference on ARDS.  
501 Definitions, mechanisms, relevant outcomes, and clinical trial coordination. *Am J Respir Crit Care Med* 1994; **149**(3  
502 Pt 1): 818-24.  
503 17. Kangelaris KN, Prakash A, Liu KD, et al. Increased expression of neutrophil-related genes in patients with  
504 early sepsis-induced ARDS. *Am J Physiol Lung Cell Mol Physiol* 2015; **308**(11): L1102-13.  
505 18. Ware LB, Koyama T, Zhao Z, et al. Biomarkers of lung epithelial injury and inflammation distinguish  
506 severe sepsis patients with acute respiratory distress syndrome. *Crit Care* 2013; **17**(5): R253.  
507 19. Calfee CS, Eisner MD, Ware LB, et al. Trauma-associated lung injury differs clinically and biologically  
508 from acute lung injury due to other clinical disorders. *Crit Care Med* 2007; **35**(10): 2243-50.  
509 20. Force ADT, Ranieri VM, Rubenfeld GD, et al. Acute respiratory distress syndrome: the Berlin Definition.  
510 *JAMA* 2012; **307**(23): 2526-33.  
511 21. Youden WJ. Index for rating diagnostic tests. *Cancer* 1950; **3**(1): 32-5.  
512 22. Sinha P, Calfee CS, Cherian S, et al. Prevalence of phenotypes of acute respiratory distress syndrome in  
513 critically ill patients with COVID-19: a prospective observational study. *Lancet Respir Med* 2020; **8**(12): 1209-18.  
514

515  
516

**Tables**

Cohort	Model	AUC	Accuracy	Sensitivity	Specificity	Classifier model-derived Hyperinflammatory Subphenotype	LCA-derived Hyperinflammatory Subphenotype
EARLI n=335	Vitals and labs	0.92 (0.90 – 0.95)	0.84	0.85	0.84	41% (139/335)	37% (124/335)
	Full features	0.92 (0.89 – 0.95)	0.85	0.80	0.88	37% (124/335)	
VALID n=452	Vitals and labs	0.88 (0.84 – 0.91)	0.80	0.66	0.86	30% (134/452)	30% (137/452)
	Full features	0.87 (0.84 – 0.90)	0.81	0.64	0.89	27% (123/452)	

517  
518  
519  
520  
521  
522

**Table 1 Model performance metrics of clinical classifier models.** Model validation in EARLI (n=335) and VALID (n=452) cohorts, for both the primary (“vitals and labs”) model and the secondary (“full featured”) model (demographics, vital signs, laboratory values, and ventilator parameters), using a probability cutoff of 0.5 for subphenotype assignments. Abbreviations: Area Under Receiver Operating Characteristic Curve (AUC with 95% confidence intervals), LCA: latent class analysis.

Cohort	Model	Outcome	Hypoinflammatory	Hyperinflammatory	Effect size	P value
EARLI n=335	Vitals and labs	Mortality*	29% (57/196)	58% (80/139)	3.3 (2.1 – 5.2)	<0.0001
		VFD	24 (0 – 28)	0 (0 – 24)	0.30 (0.18 – 0.41)	<0.0001
VALID n=452	Vitals and labs	Mortality*	27% (85/318)	49% (66/134)	2.7 (1.7 – 4.0)	<0.0001
		VFD	21 (5 – 25)	6 (0 – 22)	0.31 (0.20 – 0.41)	<0.0001
LUNG SAFE n=2813	Custom feature set	Mortality†	33% (694/2088)	57% (414/725)	2.7 (2.2 – 3.2)	<0.0001
		VFD	15 (0 – 23)	0 (0 – 19)	0.23 (0.18 – 0.28)	<0.0001

524 **Table 2 Clinical outcomes in the ARDS subphenotypes.** Mortality (count; percentage) and Ventilator Free Days  
525 (VFD; median and interquartile range) in the three observational cohorts of ARDS (EARLI, VALID, and LUNG  
526 SAFE). In EARLI (n=335) and VALID (n=452), outcomes are presented for the “vitals and labs” model, with a  
527 probability cutoff of 0.5 for subphenotype assignments. In LUNG SAFE (n=2813), outcomes are presented for a  
528 custom classifier model using a limited set of features and a probability cutoff of 0.4 for subphenotype assignments.  
529 Effect size was estimated using odds ratio for mortality and rank biserial correlation for VFD, with 95% confidence  
530 intervals. P-value represent the Chi-squared test for mortality and Wilcoxon-rank test for VFD. \* In Hospital  
531 Mortality; † 90-day Mortality.  
532



533

Subphenotype	Mortality in Low PEEP group	Mortality in High PEEP group	P value
Hyperinflammatory	62% (127/205)	54% (169/313)	0·041
Hypoinflammatory	32% (233/734)	34% (231/675)	

534 **Table 3 Mortality at day 90 in PEEP-groups stratified by ARDS Subphenotypes in LUNG SAFE (n=2813).**

535 ARDS Subphenotypes were assigned by a custom clinical classifier (“LUNG SAFE”) model using a probability  
536 cutoff of 0·4. PEEP subgroups were defined as “high PEEP” (n=992; median PEEP 11 cm H<sub>2</sub>O [10 – 12]) and “low  
537 PEEP” (n=943; median 5 cm H<sub>2</sub>O [5 – 6]) subgroups based on the mean PEEP over the first three days. P-value is  
538 for the interaction term of PEEP subgroups and ARDS subphenotypes with mortality as the dependent variable and  
539 was derived using the Wald test. Abbreviations: Positive End Expiratory Pressure (PEEP).

540 **Figures**

541

542 **Figure 1 Schematic of analysis plan.** The models were originally trained in ARMA, ALVEOLI and FACTT (n =  
543 2022) and tested in SAILS (n=745), which were all randomised controlled trials. The models were validated in two  
544 observational cohorts: EARLI (n=335) and VALID (n=452). A custom (“LUNG SAFE”) model using a limited set  
545 of predictor variables was developed to evaluate the clinical utility of ARDS subphenotypes in LUNG SAFE  
546 (n=2813) a large multinational observational cohort of ARDS. The optimal probability cutoff for the “LUNG  
547 SAFE” model determined by first evaluating the model in VALID (n=452).

548

549 **Figure 2 Receiver operating characteristic (ROC) curve for primary (“vital and labs”) model in EARLI**  
550 **(n=335) and VALID (n=452).** AUC = Area under the ROC curve. EARLI AUC = 0.92; VALID AUC = 0.88.

551

552 **Figure 3 Differences in protein biomarkers in ARDS subphenotypes.** ARDS subphenotypes were identified by  
553 the “vitals and labs” model using a probability cut-off of 0.5. Differences in biomarker data are presented in EARLI  
554 (n=335) and VALID (n=452). Y-axis was limited to aid better data visualization. Consequently, in EARLI, 9, 10,  
555 13, and 4 observations were censored, and in VALID, 13, 16, 17, and 3 observations were censored for Interleukin-  
556 6, Interleukin-8, Soluble tumor necrosis factor (TNF) receptor-1, and Protein C, respectively.

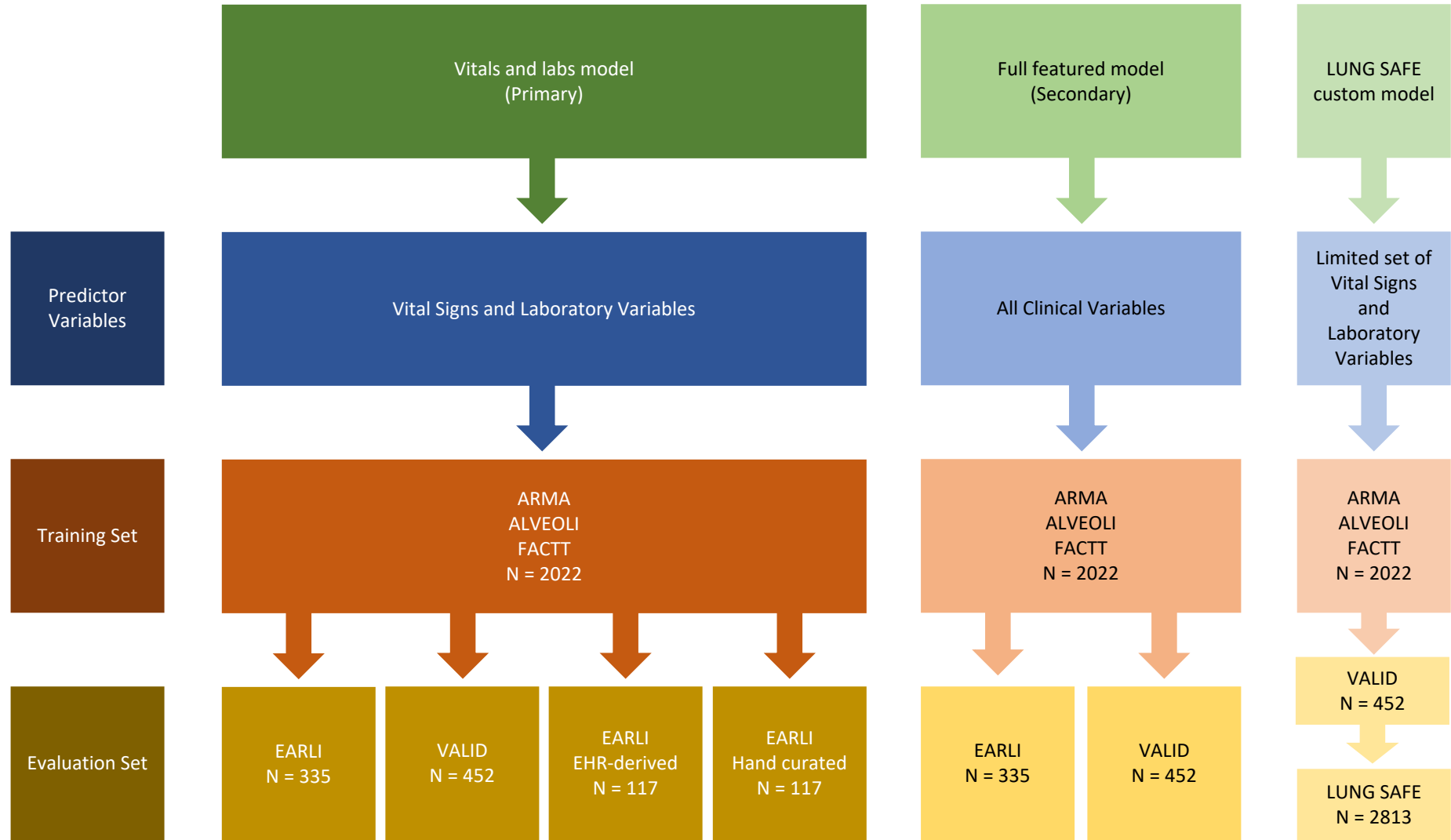
557

558 **Figure 4 Survival curves of the two ARDS subphenotypes in LUNG SAFE (n=2813).** ARDS subphenotypes  
559 were identified by a custom clinical-classifier (“LUNG SAFE”) model using a probability cutoff of 0.4 to assign  
560 class. Abbreviations: Acute Respiratory Distress Syndrome (ARDS). P-value was calculated using the log-rank test.

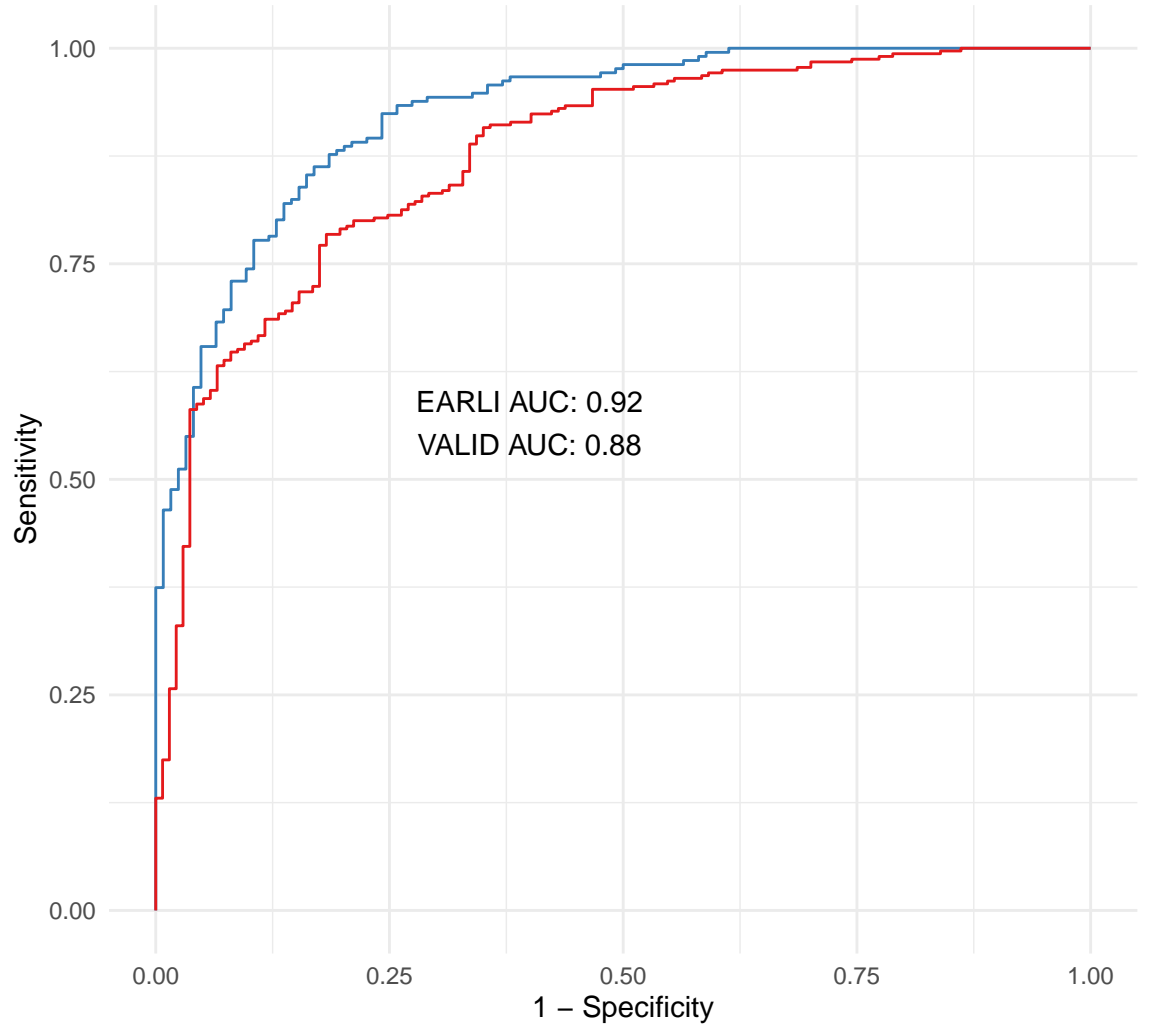
561

562 **Figure 5 Comparison of respiratory variables between the two ARDS subphenotypes in LUNG SAFE**  
563 **(n=2813).** ARDS subphenotypes were identified by a custom clinical-classifier (“LUNG SAFE”) model using a  
564 probability cutoff of 0.4 to assign class. Driving pressure is defined as the difference between plateau pressure and  
565 PEEP. Abbreviations: Positive End Expiratory Pressure (PEEP); Hyperinflammatory subphenotype (Hyper);  
566 Hypoinflammatory subphenotype (Hypo). P-value was calculated using either the t-test or Wilcoxon rank test  
567 depending on the distribution of the data.

568



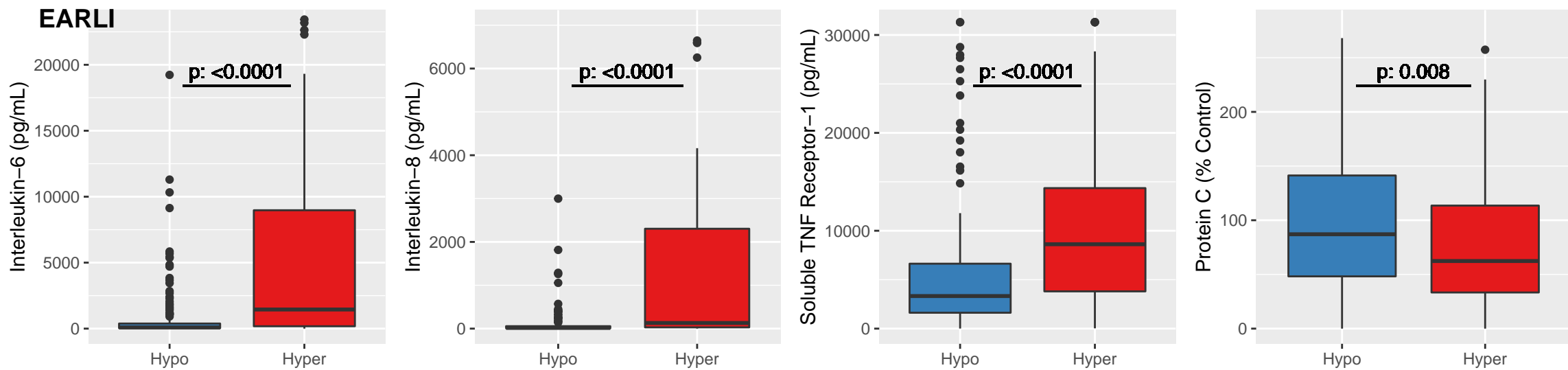
— EARLI — VALID



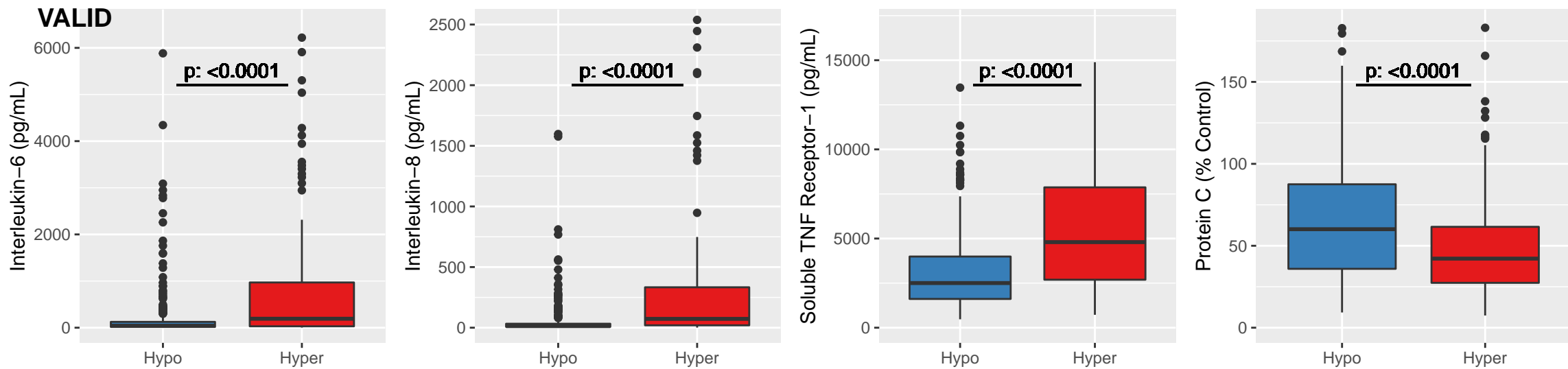
EARLI AUC: 0.92

VALID AUC: 0.88

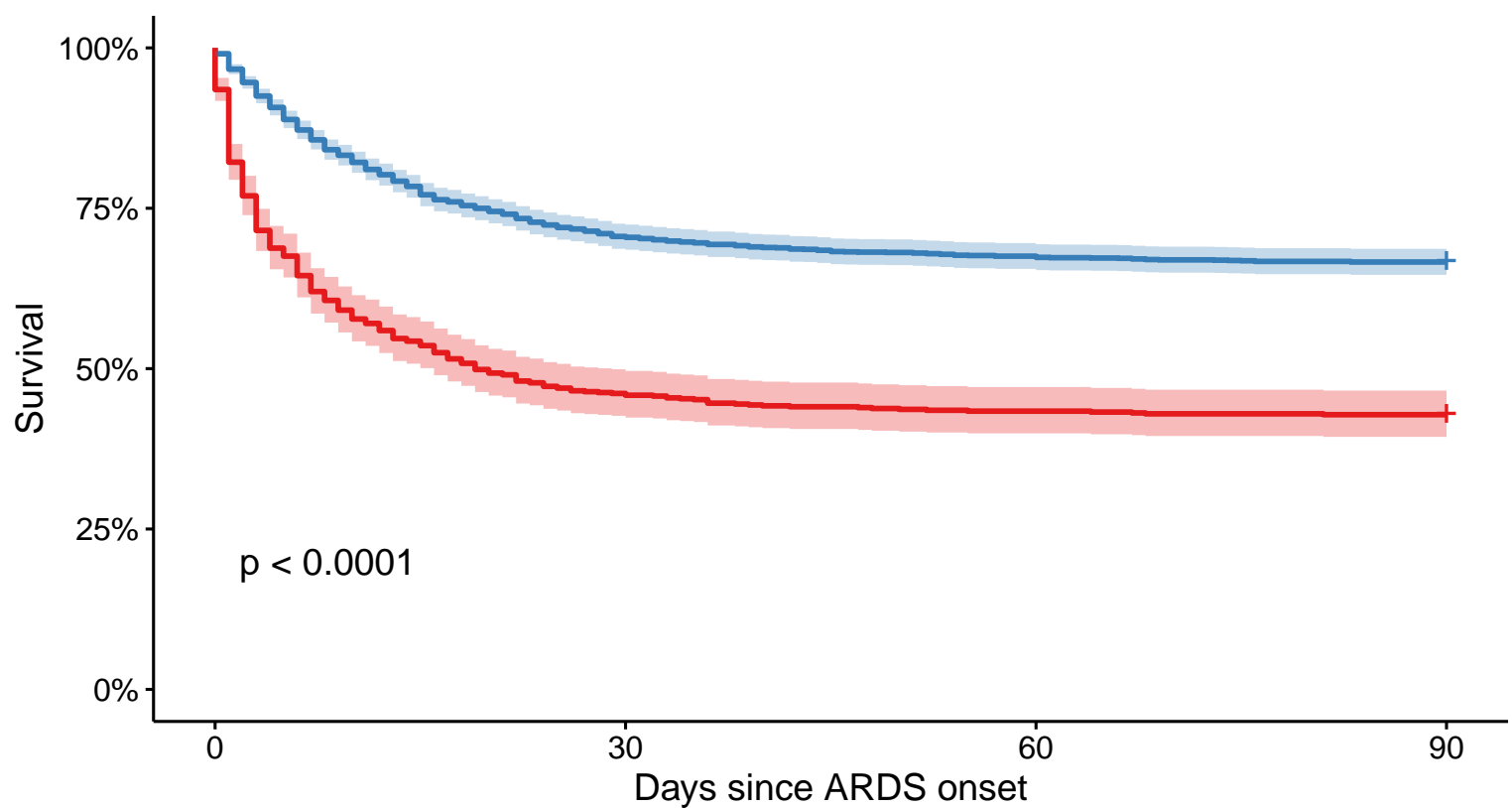
## EARLI



## VALID



Predicted Phenotype



Number at risk (number censored)

Model subphenotype classification	0	30	60	90
Hypoinflammatory	2079 (0)	1468 (0)	1404 (0)	1385 (1385)
Hyperinflammatory	724 (0)	334 (0)	314 (0)	310 (310)

Model subphenotype classification + Hypoinflammatory + Hyperinflammatory

