1 **Title: Integrating a host transcriptomic biomarker with a large language model for**
2 **diagnosis of lower respiratory tract infection**
3
4 **Authors:** [†]Hoang Van Phan[1], [†*]Natasha Spottiswoode[1], Emily C. Lydon[1], Victoria T. Chu[2,3],
5 Adolfo Cuesta[1], Alexander D. Kazberouk[4], Natalie L. Richmond[1], Padmini Deosthale[1], Carolyn
6 S. Calfee[5], Charles R. Langelier[1,3]
7
8 [†]equal contributions
9 *Corresponding author: natasha.spottiswoode@ucsf.edu
10
11 **Affiliations:**
12 [1] Department of Medicine, Division of Infectious Diseases, University of California San Francisco
13 [2] Department of Pediatrics, Division of Infectious Diseases and Global Health, University of
14 California San Francisco
15 [3] Chan Zuckerberg Biohub San Francisco
16 [4] Department of Medicine, University of California San Francisco
17 [5] Department of Medicine, Division of Pulmonary, Critical Care, Allergy and Sleep Medicine,
18 University of California San Francisco
19
20

21 **Description**. We present the novel use of a host transcriptional biomarker combined with
22 artificial intelligence analysis of electronic medical record data to diagnose lower respiratory
23 tract infections in a derivation cohort of critically ill adults, then the validation of this approach in
24 a second, fully independent, cohort. This approach demonstrated high diagnostic accuracy
25 compared to a gold standard of post-hoc multi-physician adjudication.

## Abstract

**BACKGROUND**

Lower respiratory tract infections (LRTIs) are a leading cause of mortality worldwide and can be difficult to diagnose in critically ill patients, as non-infectious causes of respiratory failure can present with similar clinical features.

**METHODS**

We developed a LRTI diagnostic method combining the pulmonary transcriptomic biomarker *FABP4* with electronic medical record (EMR) text assessment using the large language model Generative Pre-trained Transformer 4 (GPT-4). We evaluated this approach in a prospective cohort of critically ill adults with acute respiratory failure from whom tracheal aspirate *FABP4* expression was measured by RNA sequencing. Patients with LRTI or non-infectious conditions were identified using retrospective, multi-physician clinical adjudication. We then confirmed our findings by applying this method to an independent validation cohort of 115 adults with acute respiratory failure.

**RESULTS**

In the derivation cohort, a combined classifier incorporating *FABP4* expression and GPT-4– assisted EMR analysis achieved an AUC of 0.93 (±0.08) and an accuracy of 84%, outperforming *FABP4* expression alone (AUC 0.84 ± 0.11) and GPT-4–based analysis alone (AUC 0.83 ± 0.07). By comparison, the primary medical team's admission diagnosis had an accuracy of 72%. In the validation cohort, the combined classifier yielded an AUC of 0.98 (±0.04) and an accuracy of 96%.

**CONCLUSIONS**

Integrating a host transcriptional biomarker with EMR text analysis using a large language model may offer a promising new approach to improving the diagnosis of LRTIs in critically ill adults.

**INTRODUCTION**

Lower respiratory tract infections (LRTIs) are a leading cause of death worldwide, yet remain challenging to diagnose[1]. This is especially true in the intensive care unit (ICU), where non-infectious acute respiratory illnesses often have similar clinical manifestations. Further complicating accurate diagnosis is the failure to identify a causative pathogen in most clinically recognized cases of LRTI[2]. The resulting diagnostic uncertainty drives the overuse of empiric antibiotics, leading to adverse outcomes ranging from *Clostridioides difficile* infection to the development of antimicrobial resistance[3,4].

Host transcriptional biomarkers are a promising modality for LRTI diagnosis that overcome several limitations of traditional microbiologic tests[5,6]. By offering a more direct and dynamic measure of the host immune response, they can enable earlier and more accurate identification of infection, and differentiate between bacterial and viral causes of pneumonia, even in cases where pathogen detection is challenging[1,2]. Single gene biomarkers are particularly amenable to clinical translation, as they can be readily incorporated into nucleic acid amplification platforms widely used in healthcare settings.

Pulmonary *FABP4*, for instance, was recently identified as a novel LRTI diagnostic biomarker in critically ill patients with acute respiratory failure, achieving an area under the receiver operating characteristic curve (AUC) of 0.85 ± 0.12 in adults[7]. Despite better performance characteristics than existing clinical protein biomarkers such as C-reactive protein[8] or procalcitonin[9], *FABP4*, like most pneumonia diagnostic biomarkers, may not yet provide the accuracy needed to enable confident clinical decisions regarding antimicrobial use in ICU patients with acute respiratory failure.

Given that large language models (LLMs) such as Generative Pre-trained Transformer 4 (GPT-4) have shown promise in a diversity of medical applications[10], we considered the possibility that they could be leveraged to improve host biomarker-based LRTI diagnosis. LLMs have demonstrated remarkable performance for image interpretation[11], patient risk

79    stratification[12], and assisting with clinical reasoning[13-15], although their utility for diagnosing LRTI

80    or other critical illness syndromes has not been assessed. Here, we address this gap by

81    building a diagnostic classifier combining *FABP4* with GPT-4 analysis of electronic health record

82    (EHR) data. We find that this combination affords remarkably accurate LRTI diagnosis,

83    suggesting a promising new approach to improve the care of critically ill patients.

84

85    **METHODS**

86    **COHORTS AND ADJUDICATION OF LRTI STATUS**

87          We studied patients from two prospective observational cohorts of critically ill adults with

88    acute respiratory failure enrolled at the University of California San Francisco (UCSF) Medical

89    Center (**Figure 1, Table 1**). All patients were enrolled within 72 hours of intubation under UCSF

90    Institutional Review Board protocols #10-02701 (derivation cohort[16], N=202; enrolled 10/2013-

91    01/2019) or #20-30497 and #10-02852 (validation cohort; N=115; enrolled 04/2020-12/2023).

92          Adjudication of LRTI status was performed retrospectively following ICU discharge by

93    two or more physicians using all available information in the EMR, and based on the U.S.

94    Centers for Disease Control and Prevention (CDC) PNEU1 criteria[17]. Patients with a clear

95    alternative reason for their acute respiratory failure besides pulmonary infection, representing

96    the clinically relevant control group, were also identified (No LRTI group). Any adjudication

97    discrepancies were resolved by a third physician, and patients with indeterminate LRTI status

98    were excluded.

99

100    **EXTRACTION OF EMR DATA**

101          The primary medical or ICU team's clinical note from the day prior to study enrollment

102    and the chest X-ray (CXR) read from the day of enrollment were extracted from the EMR. If no

103    note was written on the day prior to enrollment, a note from two days prior was substituted

104    (**Table 1**). If no CXR was performed on the day of enrollment, the next closest CXR read prior to
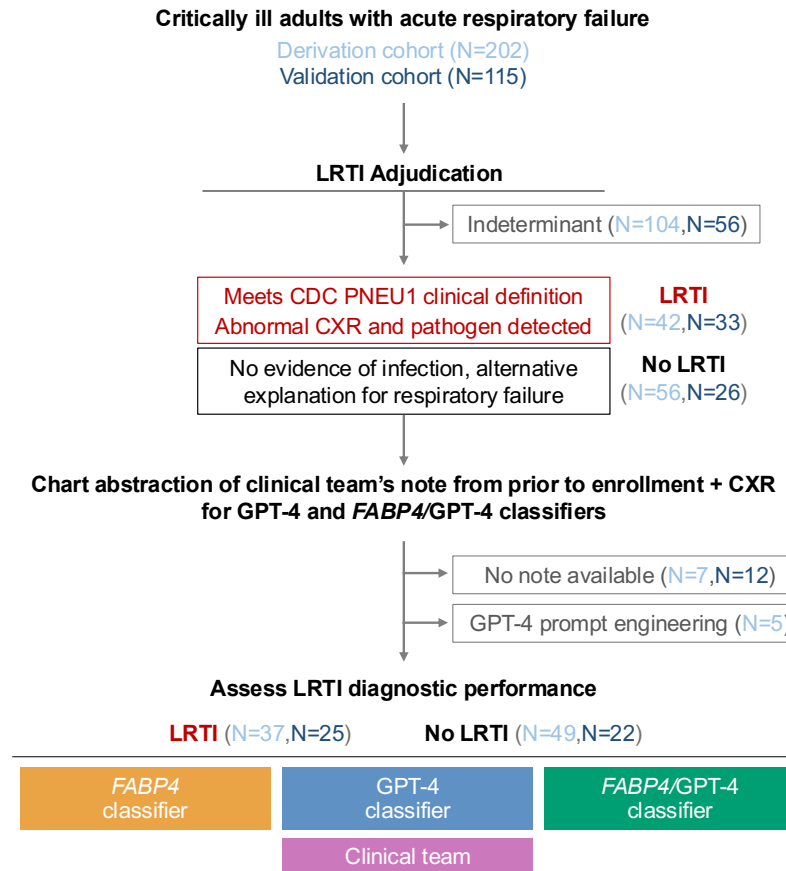
**Figure 1.** Study flow diagram and overview. Abbreviations: LRTI = lower respiratory tract infection; RNA-seq = RNA sequencing; CXR = chest X ray, FABP4 = gene encoding fatty acid binding protein 4; CDC = U.S. Centers for Disease Control and Prevention; GPT-4 = Generative Pre-trained Transformer 4.

105    the date of enrollment was used instead. Patients with no clinical notes available prior to study

106    enrollment were excluded (N=7 derivation cohort, N=12 validation cohort). The clinical treatment

107    team's ICU admission LRTI diagnosis was extrapolated based on administration of antibiotics

108    for empiric treatment of LRTI within one day of study enrollment, excluding antibiotics for

109    established non-pulmonary infections.

110

111    **RNA SEQUENCING**

112         RNA was extracted from tracheal aspirate collected on the day of enrollment and

113    underwent rRNA depletion followed by library preparation using the NEBnext Ultra 2 kit on a

114    Beckman-Coulter Echo liquid handling instrument, as previously described[7]. Finished libraries

115    underwent paired-end sequencing on an Illumina NovaSeq.

116

117    *FABP4* **DIAGNOSTIC CLASSIFIER**

118        *FABP4* expression was normalized using the varianceStabilizingTransformation function

119    from DESeq2 package (v1.42.1)[18], and used to train a logistic regression classifier. In each

120    iteration of k-fold cross-validation, both training and test sets were filtered to retain only genes

121    with at least 10 counts across 20% of the samples in the training set. The test fold's *FABP4*

122    expression level was normalized using varianceStabilizingTransformation and the dispersions of

123    the training folds, and input to the trained logistic regression classifier to assign LRTI or No LRTI

124    status for each patient in the test fold. The performance and receiver operating characteristic

125    (ROC) curve for each of the five folds was evaluated using the package pROC v1.18.5[19]. The

126    mean AUC and standard deviation were calculated from the average AUC derived from each

127    test fold. The sensitivity and specificity at the Youden's index were extracted for each test fold

128    separately using the function coords from the pROC package, and the average and standard

129    deviation was calculated across the cross-validation folds.

130

131    **GPT-4 INPUT, SCORING, AND PROMPT ENGINEERING**

132        We used the GPT-4 turbo model with 128k context length and a temperature setting of

133    0.2, implemented in Versa, a University of California San Francisco (UCSF) Health Insurance

134    Portability and Accountability Act-compliant model. For each patient, compiled clinical notes and

135    CXR reads were input into the GPT-4 chat interface. Prompt engineering was initially carried out

136    by iterative testing on clinical notes and CXR reads from five randomly selected patients in the

137    derivation cohort, who were excluded from further analyses. We employed a chain-of-thought

138    prompt strategy[20] that involved asking GPT-4 to analyze the note and CXR step-by-step. The

139    validation cohort included patients enrolled during the height of the COVID-19 pandemic and

140   thus we redacted the terms "SARS-CoV-2" or "COVID-19" from their notes to avoid biasing the

141   GPT-4 analysis. In our final version of the prompt (**Supplement, Appendix 1**), we asked GPT-4

142   to choose either LRTI or no LRTI, as exemplified in (**Supplement, Appendix 2**). For each

143   patient, GPT-4 was asked to diagnose LRTI in three separate chat sessions. A per-patient GPT-

144   4 score was calculated based on the total number of LRTI-positive diagnoses made by GPT-4.

145

146   **INTEGRATED CLASSIFIER**

147   The integrated classifier's performance was tested using 5-fold cross-validation in the

148   derivation cohort. Because of the smaller sample size, 3-fold cross-validation was used in the

149   validation cohort. For each test fold, a logistic regression classifier was trained on the remaining

150   training folds using both normalized *FABP4* expression and the GPT-4 score. The performance

151   and ROC curve for each fold was evaluated as described above. The sensitivity, specificity and

152   accuracy were calculated based on whether the out-of-fold predicted probability of LRTI is

153   greater than or equal to 50%.

154

155   **COMPARING GPT-4 TO PHYSICIANS PROVIDED THE SAME DATA**

156   We compared LRTI diagnosis by GPT-4 against LRTI diagnosis made by three

157   physicians trained in internal medicine (ADK) or additionally subspecializing in infectious

158   diseases (AC, NLR). The physicians were provided identical information and prompt as GPT-4,

159   and they were asked to assign each patient as either LRTI or No LRTI. The comparison

160   physician group score (0 to 3) was calculated based on the total number of LRTI-positive

161   diagnoses made by the comparison physicians.

162

163   **DATA AND CODE AVAILABILITY**

164       The gene count data are available at https://github.com/infectiousdisease-langelier-

165   lab/LRTI_FABP4_GPT4_classifier. The code and required source data are available at

166   https://github.com/infectiousdisease-langelier-lab/LRTI_FABP4_GPT4_classifier.

167

168   **RESULTS**

169       We evaluated the performance of four different diagnostic approaches (*FABP4*, GPT-4,

170   integrated *FABP4/*GPT-4 classifier, and admission diagnosis by the primary medical team)

171   against a gold-standard of retrospective LRTI adjudication performed by two or more

172   physicians. In the derivation cohort, this adjudication process identified 42 patients with LRTI

173   and 56 with no evidence of infection and a clear alternative explanation for respiratory failure

174   (No LRTI group) (**Figure 1**). In the validation cohort, 33 LRTI and 26 No LRTI patients were

175   identified.

176       We provided GPT-4 with practical clinical summary information from the EMR that would

177   be available to a treating physician on the day of ICU care: a CXR radiology report from the day

178   of enrollment, and a note written by the medical team from the day prior. In the derivation

179   cohort, notes and radiology reports from five patients were utilized for GPT-4 prompt

180   engineering and optimization (Methods) and seven lacked a clinical note from the day prior to

181   study enrollment, leaving a total of 37 LRTI and 49 No-LRTI cases available for analysis.

182       We first compared the accuracy of the primary medical team's ICU admission diagnosis

183   against the gold-standard retrospective LRTI adjudication. The medical team correctly identified

184   36/37 (97%) of true LRTI cases but incorrectly called LRTI in 23/49 (47%) of patients in the No

185   LRTI group, equating to an accuracy of 72% (**Figure 2A, Table 1**). We next assessed the

186   diagnostic performance of *FABP4* and found that it achieved an AUC of 0.84 ± 0.11 (mean ±

187   standard deviation) by five-fold cross validation (**Figure 2B**). We then assessed the

188   performance of GPT-4 to diagnose LRTI, with three independent diagnoses per patient. A
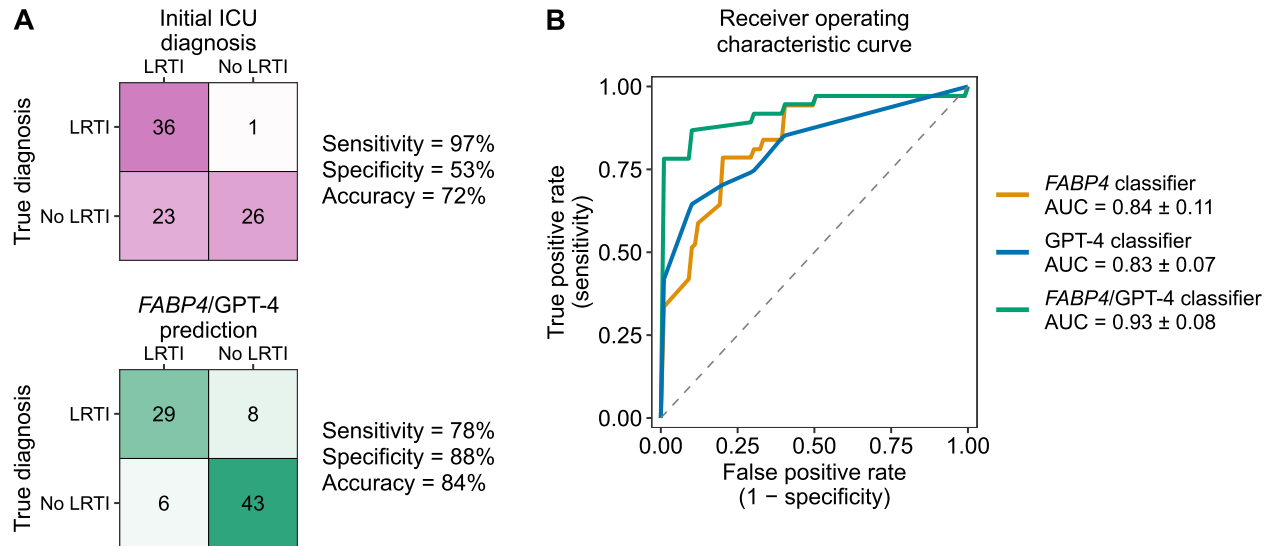
189    logistic regression classifier based on the GPT-4 score achieved an AUC of 0.83 ± 0.07 (**Figure**

190    **2B**).

191         We then combined *FABP4* and GPT-4 in a single logistic regression model and found

192    that this integrated classifier achieved an AUC of 0.93 ± 0.08 (**Figure 2B**), outperforming both

193    *FABP4* (P = 0.002, one-sided paired t-test) and GPT-4 alone (P = 0.008, one-sided paired t-

194    test). Considering an out-of-fold probability of 50% as LRTI-positive, the integrated *FABP4/*GPT-

195    4 classifier had a sensitivity of 78%, specificity of 88%, and accuracy of 84% (**Figure 2A**).

196    Assessment of the integrated classifier's performance at the Youden's index within each cross-

197    validation fold demonstrated an average sensitivity of 86%, specificity of 98%, and accuracy of

198    93%.

199         Next, we assessed the validation cohort (**Figure 1**), in which the primary medical team

200    correctly identified 25/25 (100%) of LRTI cases but unnecessarily treated for LRTI in 7/22 (32%)

201    of patients in the No LRTI group, equating to an accuracy of 85% **(Figure 2C)**. The integrated

202    *FABP4/*GPT-4 classifier achieved a sensitivity of 96%, specificity of 95%, and accuracy of 96%,

203    again outperforming either *FABP4* (accuracy 79%) or GPT-4 alone (accuracy 79%). In the

204    validation cohort, the integrated classifier achieved an AUC of 0.98 ± 0.04 using 3-fold cross-

205    validation, as compared to *FABP4* (0.86 ± 0.06) or GPT-4 (0.90 ± 0.01) alone (P = 0.08 and P =

206    0.02, respectively, one-sided paired t-test) (**Figure 2D).**

207         To gain insight into how GPT-4 returns diagnoses based on limited information, we

208    compared the LLM against the decision making of three comparison physicians provided

209    identical input. From the same limited EMR data and prompt provided to GPT-4, we asked the

210    comparison physicians to assign a diagnosis of LRTI or no evidence of LRTI for each patient in

211    the derivation cohort. Considering a threshold of at least one LRTI diagnosis per patient across

212    the three physicians as LRTI-positive, we found a sensitivity of 78%, specificity of 88%, and

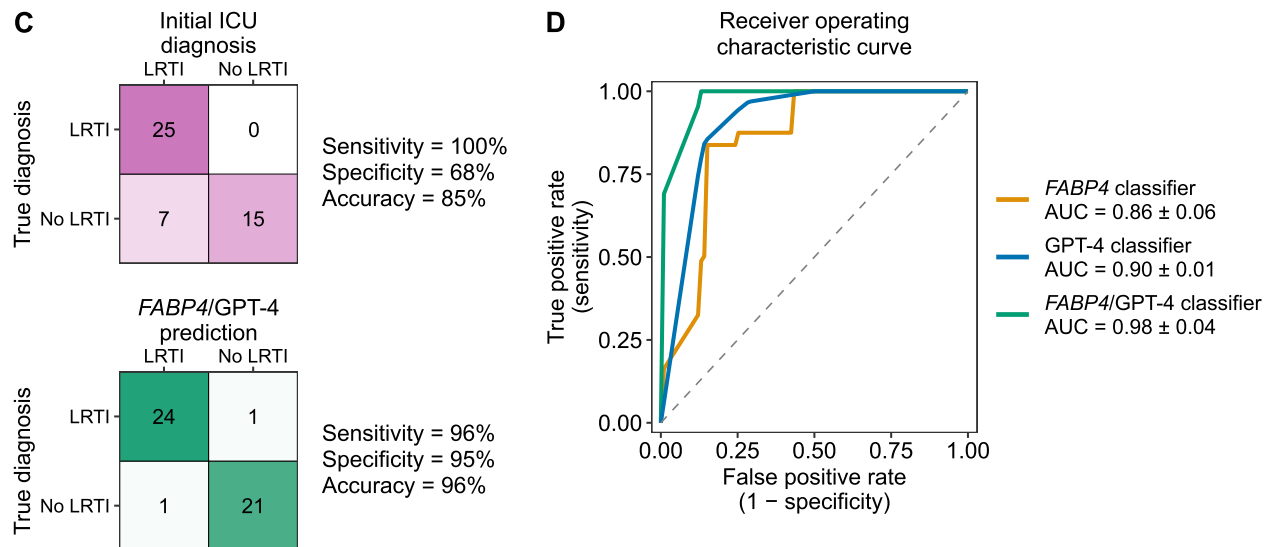213    accuracy of 84% **(Figure 3A).**

**Figure 2. Performance of *FABP4,* GPT-4 and integrated LRTI diagnostic classifiers in the derivation and validation cohorts. A)** Confusion matrices for initial ICU diagnosis and the integrated *FABP4/*GPT-4 classifier in the derivation cohort. **B)** Receiver operating characteristic curves from GPT-4 classifier, *FABP4* classifier, and integrated *FABP4/*GPT-4 classifier in the derivation cohort. **C)** Confusion matrices for initial ICU diagnosis and the integrated *FABP4/*GPT-4 classifier in the validation cohort. **D)** Receiver operating characteristic curves from GPT-4 classifier, *FABP4* classifier, and integrated *FABP4/*GPT-4 classifier in the validation cohort. In panels A and C, the classifiers output an LRTI diagnosis if the patients had a predicted out-of-fold LRTI probability of 50% or higher. In panels B and D, the area under the curves (AUCs) are presented as mean ± standard deviation.

215        Considering a threshold of at least one LRTI diagnosis per patient across the three

216    physicians as LRTI-positive, we found a sensitivity of 78%, specificity of 88%, and accuracy of

217    84% **(Figure 3A).** Finally, we sought to identify potential biases in GPT-4 diagnoses by

218    comparing GPT-4 results to those of the comparison physicians (**Figure 3B**), focusing on cases

219    with two or more discordant LRTI diagnoses. Of the nine patients more frequently diagnosed

220    with LRTI by GPT-4 versus the comparison physicians (**Figure 3B**), six had clinical notes with

221    no mention of LRTI, but explicit concern for LRTI in the CXR report. This suggested that GPT-4

222    may have placed more weight on CXR reads relative to physicians. Of the two patients

223    disproportionately diagnosed with LRTI by comparison physicians versus GPT-4 (**Figure 3B**),

224    one had a final diagnosis of e-cigarette/vaping associated lung injury, and the other had LRTI
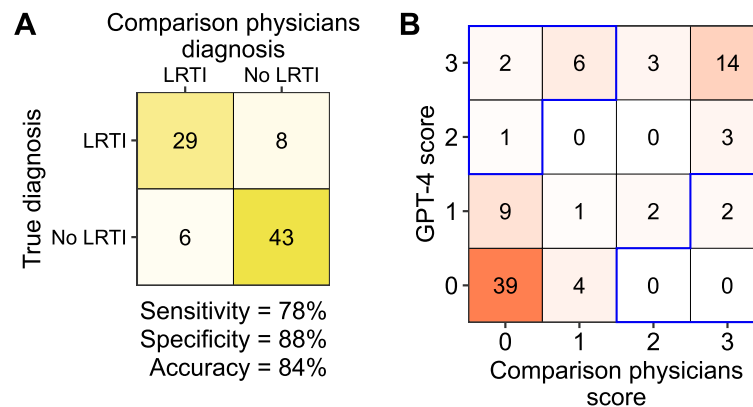
225    attributed to rhinovirus.

226



**Figure 3. Comparison of GPT-4 performance to physicians provided the same EMR data from the LRTI derivation cohort. A)** Confusion matrix of diagnosis by three GPT-4 comparison physicians who received the same prompt and data as GPT-4. **B)** Comparison of GPT-4 LRTI scores as compared to physicians. In Panel B, X-axis depicts the number of times GPT-4 diagnosed LRTI out of 3, Y-axis shows the number of times the physicians called LRTI out of 3. Blue boxes indicate instances in which GPT-4 diagnoses were most discordant with comparison physicians (the scores differ by 2 or more).

227 **DISCUSSION**

228       Our findings demonstrate that the combination of a host transcriptomic biomarker with AI

229 assessment of EMR text data can improve LRTI diagnosis in critically ill patients. We found that

230 an integrated *FABP4/*GPT-4 classifier achieved higher LRTI diagnostic accuracy than *FABP4*

231 alone, GPT-4 alone, or the treating medical team. In our study population, we found that the

232 initial treating physicians unnecessarily prescribed antibiotics in a third to half of patients

233 ultimately found to have non-infectious causes of acute respiratory failure. Had our integrated

234 classifier results been theoretically available at time of ICU admission, we estimate that

235 inappropriate antibiotic use could have been prevented in >90% of No LRTI patients who were

236 unnecessarily treated. Acute respiratory illness is a leading reason for inappropriate antibiotic

237 use[21], and our results suggest a potential role for biomarker/AI classifiers in antimicrobial

238 stewardship, a major goal of the U.S. CDC[22] and the World Health Organization[23].

239       Previous studies have found that GPT-4 is influenced by the precise language used in a

240 prompt, leading to a need for prompt engineering[14]. By iterating our prompt on a subset of

241 patients, and through direct comparison to physicians provided identical EMR data, we identified

242 possible blind spots of GPT-4 and gained insights that may help guide future optimization of

243 LLMs for infectious disease diagnosis.

244       A primary strength of this study is the novel combination of a host transcriptional

245 biomarker with AI interpretation of EMR text data to advance infectious disease diagnosis. We

246 address one of the most common and challenging diagnostic dilemmas in the ICU, leverage

247 deeply characterized cohorts, and employ a rigorous post-hoc LRTI adjudication approach

248 incorporating multiple physicians. Importantly, clinicians with access to a HIPAA-compliant GPT-

249 4 interface can readily use our prompt without any prior bioinformatics expertise. Weaknesses

250 of this study include a relatively small sample size, assessment of a biomarker not yet

251 commonly used in clinical practice, and restriction to mechanically ventilated patients.

252        Future work can test whether GPT-4 can improve the marginal performance of widely

253    available clinical biomarkers such as C-reactive protein, assess *FABP4/*GPT-4 classifier

254    performance in larger independent cohorts, and evaluate these methods for the diagnosis of

255    other critical illness syndromes such as sepsis.

256     **Corresponding Author.** Natasha Spottiswoode, MD, DPhil (natasha.spottiswoode@ucsf.edu),
257     Department of Medicine, Division of Infectious Diseases, University of California San Francisco,
258     513 Parnassus Avenue, Suite S380, San Francisco CA 94143

259     **Author Contributions.** Drs. N. Spottiswoode and C. R. Langelier had full access to all of the
260     data in the study and take responsibility for the integrity of the data and the accuracy of the data
261     analysis. Drs. H. V. Phan and N. Spottiswoode are co-first authors.
262
263     *Concept and design:* Phan, Spottiswoode, Langelier
264     *Acquisition, analysis, or interpretation of data:* All authors
265     *Drafting of the manuscript:* Phan, Spottiswoode, Langelier
266     *Critical review of the manuscript for important intellectual content:* Lydon, Chu, Calfee
267     *Statistical analysis:* Phan, Spottiswoode, Langelier
268     *Administrative, technical, or material support:* UCSF
269     *Supervision:* Calfee*,* Langelier
270

280     **Data Sharing Statement:** The gene count data are available at
281     https://github.com/infectiousdisease-langelier-lab/LRTI_FABP4_GPT4_classifier. The code and
282     required source data are available at https://github.com/infectiousdisease-langelier-
283     lab/LRTI_FABP4_GPT4_classifier.

**References**

1. The Top 10 Causes of Death. World Health Organization, 2021.
   (https://www.who.int/news-room/fact-sheets/detail/the-top-10-causes-of-death).

2. Jain S, Self WH, Wunderink RG, et al. Community-Acquired Pneumonia Requiring
   Hospitalization among U.S. Adults. N Engl J Med 2015;373(5):415-27. DOI:
   10.1056/NEJMoa1500245.

3. Langford BJ, Soucy JR, Leung V, et al. Antibiotic resistance associated with the COVID-
   19 pandemic: a systematic review and meta-analysis. Clin Microbiol Infect
   2023;29(3):302-309. DOI: 10.1016/j.cmi.2022.12.006.

4. 2022 Special Report: COVID-19 US Impact on Antimicrobial Resistance. Centers for
   Disease Control and Prevention, 2022. (https://stacks.cdc.gov/view/cdc/1190252022).

5. Tsalik EL, Henao R, Nichols M, et al. Host gene expression classifiers diagnose acute
   respiratory illness etiology. Sci Transl Med 2016;8(322):322ra11. DOI:
   10.1126/scitranslmed.aad6873.

6. Mick E, Tsitsiklis A, Kamm J, et al. Integrated host/microbe metagenomics enables
   accurate lower respiratory tract infection diagnosis in critically ill children. J Clin Invest
   2023;133(7). DOI: 10.1172/JCI165904.

7. Lydon EC, Phan HV, Mick E, et al. Pulmonary FABP4 Is an Inverse Biomarker of
   Pneumonia in Critically Ill Children and Adults. Am J Respir Crit Care Med
   2024;210(12):1480-1483. DOI: 10.1164/rccm.202403-0516RL.

8. van der Meer V, Neven AK, van den Broek PJ, Assendelft WJ. Diagnostic value of C
   reactive protein in infections of the lower respiratory tract: systematic review. BMJ
   2005;331(7507):26. DOI: 10.1136/bmj.38483.478183.EB.

9. Self WH, Wunderink RG, Jain S, Edwards KM, Grijalva CG, Etiology of Pneumonia in
   the Community Study I. Procalcitonin as a Marker of Etiology in Adults Hospitalized With

309        Community-Acquired Pneumonia. Clin Infect Dis 2018;66(10):1640-1641. DOI:

310        10.1093/cid/cix1090.

311   10.   OpenAI. Introducing ChatGPT. (https://openai.com/index/chatgpt/).

312   11.   Zhou Y, Ong H, Kennedy P, et al. Evaluating GPT-V4 (GPT-4 with Vision) on Detection of

313        Radiologic Findings on Chest Radiographs. Radiology 2024;311(2):e233270. DOI:

314        10.1148/radiol.233270.

315   12.   Beaulieu-Jones BK, Yuan W, Brat GA, et al. Machine learning for patient risk

316        stratification: standing on, or looking over, the shoulders of clinicians? NPJ Digit Med

317        2021;4(1):62. DOI: 10.1038/s41746-021-00426-3.

318   13.   Maillard A, Micheli G, Lefevre L, et al. Can Chatbot Artificial Intelligence Replace

319        Infectious Diseases Physicians in the Management of Bloodstream Infections? A

320        Prospective Cohort Study. Clin Infect Dis 2024;78(4):825-832. DOI: 10.1093/cid/ciad632.

321   14.   Lee P, Bubeck S, Petro J. Benefits, Limits, and Risks of GPT-4 as an AI Chatbot for

322        Medicine. Reply. N Engl J Med 2023;388(25):2400. DOI: 10.1056/NEJMc2305286.

323   15.   Goh E, Gallo R, Hom J, et al. Large Language Model Influence on Diagnostic

324        Reasoning: A Randomized Clinical Trial. JAMA Netw Open 2024;7(10):e2440969. DOI:

325        10.1001/jamanetworkopen.2024.40969.

326   16.   Langelier C, Kalantar KL, Moazed F, et al. Integrating host response and unbiased

327        microbe detection for lower respiratory tract infection diagnosis in critically ill adults. Proc

328        Natl Acad Sci U S A 2018;115(52):E12353-E12362. DOI: 10.1073/pnas.1809700115.

329   17.   CDC/NHSN Surveillance Definition for Specific Types of Infections.  Centers for Disease

330        Control and Prevention,

331        2021.( https://www.cdc.gov/nhsn/pdfs/pscmanual/17pscnosinfdef_current.pdf)

332   18.   Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for

333        RNA-seq data with DESeq2. Genome Biol 2014;15(12):550. DOI: 10.1186/s13059-014-

334        0550-8.

335    19.    Robin X, Turck N, Hainard A, et al. pROC: an open-source package for R and S+ to

336           analyze and compare ROC curves. BMC Bioinformatics 2011;12:77. DOI: 10.1186/1471-

337           2105-12-77.

338    20.    Wei J, Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q.V., Zhou,

339           D. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models.  Advances

340           in Neural Information Processing Systems 35 (NeurIPS 2022): NeurIPS Proceedings

341           2022.

342    21.    Merenstein DJ, Barrett B, Ebell MH. Antibiotics Not Associated with Shorter Duration or

343           Reduced Severity of Acute Lower Respiratory Tract Infection. J Gen Intern Med 2024.

344           DOI: 10.1007/s11606-024-08758-y.

345    22.    Antibiotic Resistance Threats In the United States. U.S. Department of Health and

346           Human Services, 2019. (https://www.cdc.gov/antimicrobial-resistance/media/pdfs/2019-

347           ar-threats-report-508.pdf).

348    23.    Antimicrobial stewardship interventions: a practical guide. World Health Organization,

349           2021. (https://www.who.int/europe/publications/i/item/9789289056267).

350

351 **Table 1. Clinical and demographic features of cohorts.**

| | Derivation Cohort | | | Validation Cohort | | |
|---|---|---|---|---|---|---|
| | **LRTI** | **No LRTI** | *P* | **LRTI** | **No LRTI** | *P* |
| **N** | 37 | 49 | | 25 | 22 | |
| **Age, years** (Median, IQR) | 65.0 (51.0 – 75.0) | 62.0 (53.0-73.0) | 0.81 | 53.9 (49.4 – 66.9) | 58.5 (57.5 - 63.9) | 0.69 |
| **Female Sex** (No., %) | 11 (30) | 30 (61) | 0.0074 | 12 (48) | 8 (36) | 0.61 |
| **Race** (No., %) | | | 0.91 | | | 0.10 |
| White | 18 (49) | 24 (49) | - | 5 (20) | 10 (46) | - |
| Black/African American | 4 (11) | 5 (10) | - | 2 (8) | 1 (5) | - |
| Asian | 8 (22) | 9 (18) | - | 3 (12) | 4 (18) | - |
| Native Hawaiian/Pacific Islander | 1 (3) | 0 (0) | - | 1 (4) | 2 (9) | - |
| Other/Unknown | 6 (16) | 11 (22) | - | 14 (56) | 5 (23) | - |
| **Hispanic ethnicity** (No., %) | 5 (14) | 11 (22) | 0.44 | 10 (40) | 5 (23) | 0.40 |
| **Comorbidities** (No., %) | 36 (97) | 45 (92) | 0.54 | 24 (96) | 17 (77) | 0.14 |
| **Immunosuppressed** (No., %) | 8 (22) | 6 (12) | 0.38 | 1 (4) | 3 (14) | 0.51 |
| **Microbiologic diagnosis** (No., %) | | | | | | |
| Bacterial | 28 (76) | | | 3 (12) | | |
| Viral | 4 (11) | | | 1 (4) | | |
| SARS-CoV-2 | 0 (0) | | | 13 (52) | | |
| Fungal | 1 (3) | | | 0 (0) | | |
| Multiple | 4 (11) | | | 8 (32) | | |
| **No-LRTI group cause of respiratory failure** (No., %) | | | - | | | - |
| Surgery | | 14 (29) | - | | 6 (27) | - |
| Neurologic | | 12 (25) | - | | 7 (32) | - |
| Cardiovascular | | 8 (16) | - | | 4 (18) | - |
| Non-LRTI infection | | 5 (10) | - | | 6 (27) | - |
| Other | | 11 (22) | - | | 2 (9) | - |
| **ICU admission diagnosis** (No., %) | | | - | | | - |
| LRTI | 36 (97) | 23 (47) | - | 25 (100) | 7 (32) | - |
| No LRTI | 1 (3) | 26 (53) | - | 0 (0) | 15 (68) | - |
| **Clinical team writing note** (No., %) | | | 0.20 | | | <0.0001 |
| Medicine | 15 (41) | 9 (18) | - | 1 (4) | 5 (23) | - |
| Critical Care | 5 (14) | 7 (14) | - | 24 (96) | 6 (27) | - |
| Neurosurgery | 3 (8) | 9 (18) | - | 0 (0) | 4 (18) | - |
| Cardiology | 3 (8) | 6 (12) | - | 0 (0) | 2 (9) | - |
| Other | 11 (30) | 18 (37) | - | 0 (0) | 5 (23) | - |
| **Note to enrollment** (No., %) | | | 0.080 | | | 0.11 |
| 1 day | 34 (92) | 49 (100) | - | 21 (84) | 22 (100) | - |
| 2 days | 3 (8) | 0 (0) | - | 4 (16) | 0 (0) | - |
| **CXR to enrollment** (No., %) | | | 0.24 | | | 0.11 |
| 0 days | 30 (81) | 33 (67) | - | 19 (76) | 11 (50) | - |
| 1 day | 7 (19) | 13 (27) | - | 4 (16) | 9 (41) | - |
| 2 days | 0 (0) | 3 (6) | - | 2 (8) | 2 (9) | - |

352 Chi-squared test used for all categorical variables except time from note to enrollment for which Fisher exact test
353 was used. Wilcoxon rank-sum test used for continuous variables (age). IQR = interquartile range. One No LRTI
354 patient in the derivation cohort, and three No LRTI patients in the validation cohort, were adjudicated as having ≥1
355 etiology of respiratory failure. Four LRTI patients in the derivation cohort were adjudicated as having more than

356    one microbiologic diagnosis (three with viral-bacterial co-infection, and one with viral-fungal co-infection) and eight

357    LRTI patients in the validation cohort were adjudicated as having more than one microbiologic diagnosis (all with

358    SARS-CoV-2-bacterial co-infection).