

**STUDY PROTOCOL**

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients through a Web-based  
Randomized Clinical Vignette Multicenter Study

**Co-Principal Investigators: Michael W. Sjoding, Jenna Wiens**

**Funded by: NIH NHLBI R01 HL158626**

**last updated:**

**8 June 2023**

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13

Principal Investigator(s)	<p>J. Wiens, PhD</p> <p>Department of Computer Science and Engineering 2260 Hayward Street Ann Arbor, MI 48109 (734) 647-4832 Email: <a href="mailto:wiensj@umich.edu">wiensj@umich.edu</a></p> <p>Michael W. Sjoding, MD</p> <p>Internal Medicine, Division of Pulmonary and Critical Care G027W Building 16 NCRC, 2800 Plymouth Road, SPC 2800, Ann Arbor, MI 48109 Phone: 734-763-1796 <a href="mailto:msjoding@umich.edu">msjoding@umich.edu</a></p>
---------------------------	--

14

15

16	Table of Contents	
17	<b>1. INTRODUCTION AND RATIONALE</b> .....	<b>8</b>
18	<b>2. OBJECTIVES</b> .....	<b>8</b>
19	<b>3. Study Design</b> .....	<b>9</b>
20	<b>4. Study Population</b> .....	<b>10</b>
21	4.1 Population (base) .....	<b>10</b>
22	4.2 Inclusion criteria .....	<b>10</b>
23	4.3 Exclusion criteria .....	<b>10</b>
24	4.4 Screen failure .....	<b>10</b>
25	4.5 Technical failure .....	<b>10</b>
26	4.6 Participant withdrawal .....	<b>11</b>
27	4.7 Sample size calculation.....	<b>11</b>
28	<b>5. TREATMENT OF SUBJECTS</b> .....	<b>12</b>
29	5.1 Intervention .....	<b>12</b>
30	5.2 Use of co-intervention (if applicable).....	<b>13</b>
31	5.3 Escape medication (if applicable) .....	<b>13</b>
32	<b>6. INVESTIGATIONAL PRODUCT</b> .....	<b>13</b>
33	6.1 Name and description of investigational product(s).....	<b>13</b>
34	6.2 Summary of findings from non-clinical studies. ....	<b>14</b>
35	6.3 Summary of findings from clinical studies.....	<b>14</b>
36	6.4 Summary of known and potential risks and benefits .....	<b>15</b>
37	6.5 Description and justification of route administration and dosage .....	<b>15</b>
38	6.6 Dosages, dosage modifications and method of administration .....	<b>15</b>
39	6.7 Preparation and labelling of Investigational Medicinal Product .....	<b>16</b>
40	6.8 Drug accountability .....	<b>16</b>
41	<b>7. METHODS</b> .....	<b>16</b>
42	<b>7.1 STUDY PARAMETERS/ENDPOINTS</b> .....	<b>16</b>
43	7.1.1 Main study parameters/endpoints .....	<b>16</b>
44	7.2 Randomization, blinding, and treatment allocation .....	<b>16</b>
45	7.3 Study procedures .....	<b>16</b>
46	7.4 Withdrawal of individual subjects .....	<b>17</b>
47	7.4.1 Specific criteria for withdrawal (if applicable).....	<b>17</b>
48	7.5 Replacement of individual subjects after withdrawal.....	<b>17</b>
49	7.6 Follow-up of subjects withdrawn from treatment.....	<b>17</b>

50	7.7 Premature termination of the study .....	17
51	<b>8. SAFETY REPORTING.....</b>	<b>17</b>
52	8.1 Temporary halt for reasons of subject safety .....	17
53	8.2 AEs, SAEs, SUSARs .....	17
54	8.2.1 Adverse events (AEs) .....	17
55	8.2.2 Serious adverse events (SAEs).....	18
56	8.3 Annual safety report .....	18
57	8.4 Follow-up of adverse events.....	18
58	8.5 Data Safety Monitoring Board (DSMB) / Safety Committee .....	18
59	<b>9. STATISTICAL ANALYSIS.....</b>	<b>18</b>
60	9.1 Primary study parameters/endpoints.....	18
61	9.2 Interim analysis (if applicable).....	18
62	9.3 Statistical analysis plan.....	18
63	<b>10. ETHICAL CONSIDERATIONS.....</b>	<b>19</b>
64	10.1 Regulation statement .....	19
65	10.2 Recruitment and consent .....	19
66	10.3 Objection by minors or incapacitated subjects (if applicable).....	19
67	10.4 Benefits and risks assessment, group relatedness.....	20
68	10.5 Compensation for injury .....	20
69	10.6 Incentives (if applicable).....	20
70	<b>11. ADMINISTRATIVE ASPECTS, MONITORING AND PUBLICATION.....</b>	<b>20</b>
71	11.1 Handling and storage of data and documents.....	20
72	11.2 Monitoring and Quality Assurance .....	20
73	11.3 Public disclosure and publication policy.....	20
74	<b>12. AMENDMENTS.....</b>	<b>20</b>
75	<b>13. REFERENCES.....</b>	<b>20</b>

76  
77

78 List of abbreviations

79

<b>AE</b>	<b>Adverse Event</b>
<b>ARF</b>	<b>Acute Respiratory Failure</b>
<b>COPD</b>	<b>Chronic obstructive pulmonary disease</b>

80

81 **SUMMARY**

82 **Rationale:**

83

84 Acute respiratory failure (ARF) develops in over 3 million patients hospitalized in the United  
85 States annually.<sup>1</sup> Pneumonia, heart failure, and/or chronic obstructive pulmonary disease  
86 (COPD) are 3 of the most common reasons for ARF,<sup>2</sup> and these conditions are among the top  
87 reasons for hospitalization in the United States.<sup>3</sup> Determining the underlying causes of ARF is  
88 critically important for guiding treatment decisions, but can be clinically challenging, as initial  
89 testing such as brain natriuretic peptide (BNP) levels or chest radiograph results can be non-  
90 specific or difficult to interpret.<sup>4</sup> This is especially true for older adults,<sup>5</sup> patients with comorbid  
91 illnesses,<sup>6</sup> or more severe disease.<sup>7</sup> Incorrect initial treatment often occurs, resulting in worse  
92 patient outcomes or treatment delays.<sup>8</sup> Artificial intelligence technologies have been proposed  
93 as a strategy for improving medical diagnosis by augmenting clinical decision-making,<sup>9</sup> and  
94 could play a role in the diagnostic evaluation of patients with ARF.

95

96 Artificial intelligence (AI) has achieved high accuracy at identifying abnormalities in clinical  
97 images, such as pneumonia from chest radiographs, diabetic retinopathy from fundus images,  
98 or skin cancer from histopathology images.<sup>10-12</sup> However, systematic bias in AI models can lead  
99 to inaccurate predictions for entire subpopulations.<sup>13-15</sup> When presented with such incorrect  
100 predictions, physician performance can be harmed<sup>16</sup> due to automation bias,<sup>17</sup> which is  
101 especially concerning in safety-critical settings. Thus, the extent to which AI can be safely  
102 integrated into clinical workflows and to support diagnostic decisions is still unknown.

103

104 This study aims to study the effectiveness of providing clinicians with image-based AI model  
105 explanations to help them catch when models are making incorrect decisions.

106

107

108 **Study design:**

109 This is web-based randomized clinical vignette study.

110

111 **Objectives**

112

113 Survey Data Collection Phase

114 *Objectives*

- 115 • What is clinician accuracy in diagnosing pneumonia, heart failure, and COPD in a patient  
116 population with ARF without any AI model input?
- 117 • How do standard AI model predictions without explanations affect clinician accuracy in  
118 diagnosing pneumonia, heart failure, and COPD in a patient population with ARF?

- 119
- How do standard AI model predictions with explanations affect clinician accuracy in  
120 diagnosing pneumonia, heart failure, and COPD in a patient population with ARF?
  - How do intentionally biased AI model predictions without explanations affect clinician  
121 accuracy in diagnosing pneumonia, heart failure, and COPD in a patient population with  
122 ARF?
  - How do intentionally biased AI model predictions with explanations affect clinician  
123 accuracy in diagnosing pneumonia, heart failure, and COPD in a patient population with  
124 ARF?
  - How does text input (always accurate) from a clinician affect clinician accuracy in  
125 diagnosing pneumonia, heart failure, and COPD in a patient population with ARF?  
126
  - How does text input (always accurate) from a clinician affect clinician accuracy in  
127 diagnosing pneumonia, heart failure, and COPD in a patient population with ARF?  
128
- 129

### 130 **Study population**

131

132 Hospitalist physicians, nurse practitioners, and physician assistants who commonly care for  
133 patients with ARF from 12 US hospitals.

134

### 135 **Intervention:**

136

137 Within the clinical vignette survey, the primary research question is to understand the impact of  
138 providing AI model explanations to clinicians. Therefore, we will randomize participants to see:  
139

- 140 (1) AI model explanations vs. no AI model explanations: When clinicians are shown AI  
141 models, they will be randomized to see AI model prediction alone each time they are  
142 shown an AI model, or randomized to see AI model predictions with an explanation each  
143 time they are shown an AI model.
- 144

145 Within the survey, they will also be randomized to see

- 146 (2) Type of bias: they type of systematically biased AI model shown in the vignette, either  
147 against age, BMI, or model preprocessing.
- 148 (3) Vignette ordering: for vignettes with standard model predictions or intentionally biased  
149 predictions (vignettes 3-8), the ordering of these will be randomized.
- 150

151 We look at the effect of standard model predictions and standard model predictions with  
152 explanations to test if such model input improves clinical diagnostic accuracy. We also look at  
153 the effect of intentionally biased AI model predictions to test if such inputs hurt diagnostic  
154 accuracy, and whether providing explanations when clinicians are shown systematically biased  
155 AI models help clinicians recover in terms of diagnostic accuracy.

156

157 **Main study parameters/endpoints: Clinician Diagnostic accuracy after reviewing clinical**  
158 **vignette. We will specifically evaluate the following:**

- 159 • Do standard model predictions improve clinician diagnostic accuracy?
- 160 • Do standard model predictions with explanations further improve clinician diagnostic  
161 accuracy?
- 162 • Do intentionally biased model predictions hurt clinician diagnostic accuracy?
- 163 • Do intentionally biased model explanations help clinicians recover from the negative  
164 effects of intentionally biased model predictions?

165  
166 **Nature and extent of the burden and risks associated with participation, benefit, and group**  
167 **relatedness:**

168  
169 Because this is a web-based study involving benign (non-harmful) behavioral interventions, no  
170 adverse events are expected during this study.

171

## 1. INTRODUCTION AND RATIONALE

Acute respiratory failure (ARF) develops in over 3 million patients hospitalized in the United States annually.<sup>1</sup> Pneumonia, heart failure, and/or chronic obstructive pulmonary disease (COPD) are 3 of the most common reasons for ARF,<sup>2</sup> and these conditions are among the top reasons for hospitalization in the United States.<sup>3</sup> Determining the underlying causes of ARF is critically important for guiding treatment decisions, but can be clinically challenging, as initial testing such as brain natriuretic peptide (BNP) levels or chest radiograph results can be non-specific or difficult to interpret.<sup>4</sup> This is especially true for older adults,<sup>5</sup> patients with comorbid illnesses,<sup>6</sup> or more severe disease.<sup>7</sup>

Incorrect initial treatment for ARF often occurs, resulting in worse patient outcomes or treatment delays.<sup>8</sup> Artificial intelligence technologies have been proposed as a strategy for improving medical diagnosis by augmenting clinical decision-making,<sup>9</sup> and could play a role in the diagnostic evaluation of patients with ARF. If integrated into clinical workflows effectively, such technologies could improve clinician diagnostic accuracy for ARF and result in better patient outcomes. We developed an artificial intelligence algorithm that can predict the underlying etiologies of ARF based on patient chest X-rays and clinical data. Theoretically, this algorithm could improve clinician's diagnostic accuracy. While promising, systematic bias in AI models can lead to inaccurate predictions, ultimately hurting physician performance. Thus, the extent to which AI can be safely integrated into clinical workflows to support diagnostic decisions is still unknown.

The model developed in this study predicts whether the patient has pneumonia, heart failure, and/or COPD based on their chest X-ray and clinical data. The model also has corresponding explanations based on the chest X-ray, which highlight the areas that the model found important for its decision. The standard model developed highlights clinically relevant regions for pneumonia (e.g., lungs), heart failure (e.g., enlarged heart), and COPD (e.g., tracheal narrowing), whereas the systematically biased models highlight clinically irrelevant findings for pneumonia (bone density for age), heart failure (body mass for BMI), and COPD (features of image preprocessing blur). In this web-based study, we will test our hypothesis that providing participants with systematically biased predictions without explanations will hurt their diagnostic accuracy, whereas providing them with systematically biased predictions with explanations will help them recover from these negative effects.

## 2. OBJECTIVES



209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225

Objectives:

- To determine clinician accuracy in diagnosing pneumonia, heart failure, and COPD in a patient population with ARF without any AI model input?
- To determine how standard AI model predictions without explanations affect clinician accuracy in diagnosing pneumonia, heart failure, and COPD in a patient population with ARF?
- To determine how standard AI model predictions with explanations affect clinician accuracy in diagnosing pneumonia, heart failure, and COPD in a patient population with ARF?
- To determine how do systematically biased AI model predictions without explanations affect clinician accuracy in diagnosing pneumonia, heart failure, and COPD in a patient population with ARF?
- To determine how do systematically biased AI model predictions with explanations affect clinician accuracy in diagnosing pneumonia, heart failure, and COPD in a patient population with ARF?

226 **3. Study Design**

227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243

We aim to include 400 participants from multiple hospital centers. Participants will be invited to participate in randomized clinical vignette survey in which participants are randomly shown 9 clinical vignettes out of 45 possible vignettes. The first two vignettes are not accompanied by AI model predictions and are used to estimate baseline participant diagnostic accuracy. The next 6 vignettes include AI model predictions, but half of the participants will be randomized to also see AI model explanations when shown the AI model predictions. These 6 vignettes include 3 vignettes with standard model predictions and 3 with systematically biased model predictions. Participants will be randomized to see one of three types of systematically biased AI model predictions, and the ordering of the 3 standard and 3 systematically biased model predictions are randomized. In the final vignette, all participants are provided a clinical consult, which is a short narrative provided by a hypothetical trusted colleague, who describes the rationale behind which diagnoses were most likely and what treatment plan they recommend. By design, the clinical consult always provides the correct diagnosis and appropriate treatment plan to provide a realistic upper bound of participant diagnostic accuracy.

244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279

## 4. Study Population

### 4.1 Population (base)

Hospitalist physicians, nurse practitioners, and physician assistants who commonly care for patients with acute respiratory failure.

### 4.2 Inclusion criteria

To be eligible to participant in this study, a participant must answer “Yes” to the following question:

“Do you hold any of the following roles on a healthcare team, or any similar roles?”

- Nurse Practitioner (NP)
- Physician Assistant
- Resident
- Fellow
- Attending Physician

### 4.3 Exclusion criteria

Any participant answering “No” to the following question will be excluded from the study:

“Do you hold any of the following roles on a healthcare team, or any similar roles?”

- Nurse Practitioner (NP)
- Physician Assistant
- Resident
- Fellow
- Attending Physician

### 4.4 Screen failure

Not applicable.

### 4.5 Technical failure

280 In the unlikely event that data is not properly recorded to the Qualtrics survey  
281 interface, the participant responses will be excluded.

282

#### 283 4.6 Participant withdrawal

284

285 Any completed vignette will be analyzed, even if the participant does not  
286 complete all 9 vignettes.

287

#### 288 4.7 Sample size calculation

289

290 A sample size of 400 will have 80% power to detect a decrease in accuracy of  
291 25% with the systematically biased AI model compared to baseline and a 10%  
292 improvement with the biased AI model with explanations compared to no  
293 explanations using a generalized linear mixed model with a 0.001 significance  
294 level.

295

296 The sample size calculation is based on the primary endpoint of clinician  
297 diagnostic accuracy for pneumonia, heart failure, and COPD. To calculate  
298 diagnostic accuracy, dichotomized responses are compared to the reference  
299 standard labels generated by a group of 5 physicians who reviewed the patients  
300 complete medical record. Each vignette's three diagnosis responses are analyzed  
301 separately within the generalized linear model.

302

303

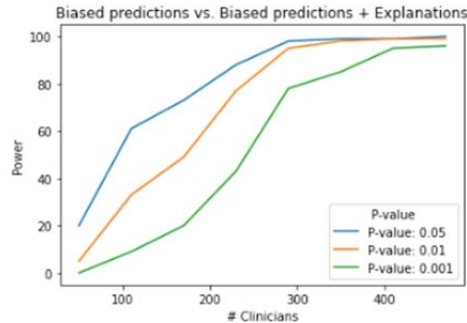
304 Sample size calculations were based on 100 simulated studies performed at  
305 sample size levels of 50 to 500, in increments of 50. The simulations were based  
306 on the assumptions that clinician diagnostic accuracy was 0.68 for pneumonia,  
307 0.72 for heart failure, and 0.82 for COPD. The systematically biased model had an  
308 accuracy of 0.33. Clinician performance was simulated such that clinicians  
309 listened to the biased model 50% of the time. Furthermore, when presented  
310 with biased model explanations, clinicians recovered by 50%.

311

312 Given these simulated data, we fit a generalized linear mixed model in R to  
313 measure if the recovery of clinician performance given the model explanation is  
314 statistically significant. We repeat this for every simulated study, and calculate  
315 power as the percentage of time the effect of the explanation is statistically  
316 significant across all simulations.

317

318 We determine that we have 80% power to detect a statistically significant effect  
319 of the explanation at a significance level of 0.001:  
320  
321



322 Figure 1. Sample Size Simulation Power Plot.  
323  
324

325 **Details sample size calculation:**  
326

327 **Expected loss of data:**  
328

329 If data is not properly stored on the Qualtrics server, the participants responses  
330 will be withdrawn.  
331

332 **5. TREATMENT OF SUBJECTS**

333  
334 **5.1 Intervention**  
335

336 The purpose of the study is to understand the effect of providing AI model explanations  
337 in addition to AI model predictions on clinicians’ diagnostic and treatment decisions  
338 when diagnosing the underlying causes of acute respiratory failure. In this study, we  
339 investigate the use of gradCAM heatmaps as an image-based explanation of the AI  
340 model’s decision.<sup>18</sup>  
341

342 GradCAM heatmaps are a commonly used model explanation tool by AI model  
343 developers.<sup>19</sup> It is used to highlight the regions of an image used by an AI model to make  
344 its prediction. For example, a gradCAM heatmap generated from a model trained to  
345 predict heart failure based on a patient’s chest X-ray might highlight the patient heart.  
346 Testing the usefulness of gradCAM heatmaps means presenting these heatmaps with AI

347 model predictions, when the AI model predicts that a patient has disease. This does not  
348 induce any harm or risk to the patients in the vignettes or participants in the study.

349  
350  
351 The AI models provide a score for each diagnosis (pneumonia, heart failure, and COPD)  
352 on a scale of 0-100, with a score above 50 corresponding to a positive diagnosis. In  
353 general, when the standard AI model predicts a positive diagnosis, the explanation is  
354 expected to highlight the relevant region of the chest X-ray (e.g., lung infiltrate). When  
355 participants are shown a systematically biased AI models, they are randomized to 1 of 3  
356 intentionally biased AI models based on patient age (predicting pneumonia if age  $\geq$  80  
357 years), BMI (predicting heart failure if BMI  $\geq$  30), or chest X-ray preprocessing (predicting  
358 COPD if a blur was applied to the X-ray). Explanations associated with the systematically  
359 biased AI models were generated based on models trained to predict age, BMI, and  
360 preprocessing parameters, and highlighted areas of the X-ray corresponding to age, BMI,  
361 or preprocessing (e.g., low bone density, soft tissue).

362  
363 We aim to test our hypothesis that providing participants with systematically biased  
364 predictions without explanations will hurt their diagnostic accuracy, whereas providing  
365 them with systematically biased predictions with explanations will help them recover  
366 from these negative effects. While image-based explanations have been studied in  
367 various settings, this will be the first to test the use in the diagnosis of acute respiratory  
368 failure at this scale.

369  
370 **5.2 Use of co-intervention (if applicable)**

371 Not applicable.

372  
373 **5.3 Escape medication (if applicable)**

374 Not applicable.

375

## 376 6. INVESTIGATIONAL PRODUCT

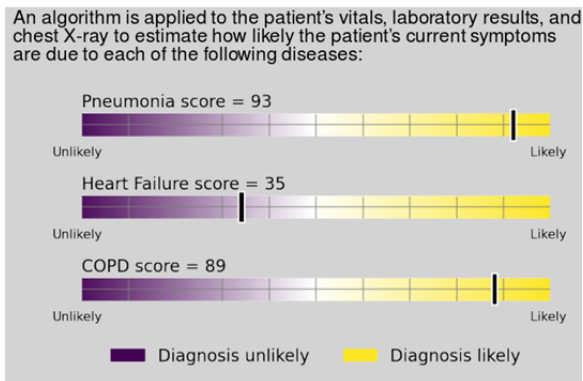
377 The AI model evaluated in the clinical vignette study is based on Jabbour et al.<sup>20</sup>.

### 378 6.1 Name and description of investigational product(s)

379  
380 This model takes as input the patient's clinical data and chest X-ray at the time of ARF  
381 and outputs three separate probabilities that the patient has pneumonia, heart failure,  
382 and COPD.

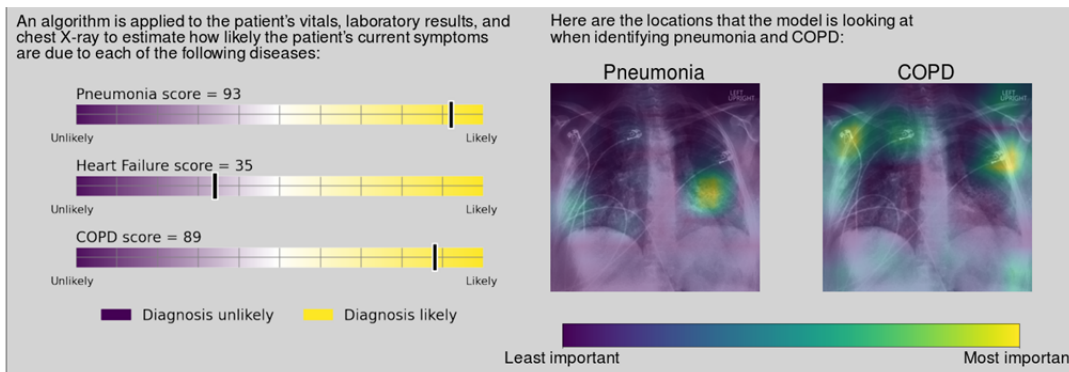
383

384 The model provides a score of 0-100 for each of the diagnoses and presents them on a  
385 color bar to indicate the likelihood of each disease:  
386



387  
388 Figure 2. Model scores for each disease.

389  
390 When the participant is randomized to see model explanations, the model also presents  
391 an explanation for each diagnosis, if the score for the diagnosis is greater than 50  
392 (indicating a positive diagnosis).



393  
394 Figure 3. Model scores each disease and corresponding explanations when the model predicts a  
395 positive diagnosis.

## 396 6.2 Summary of findings from non-clinical studies.

397  
398  
399 These models were not tested in non-clinical settings.

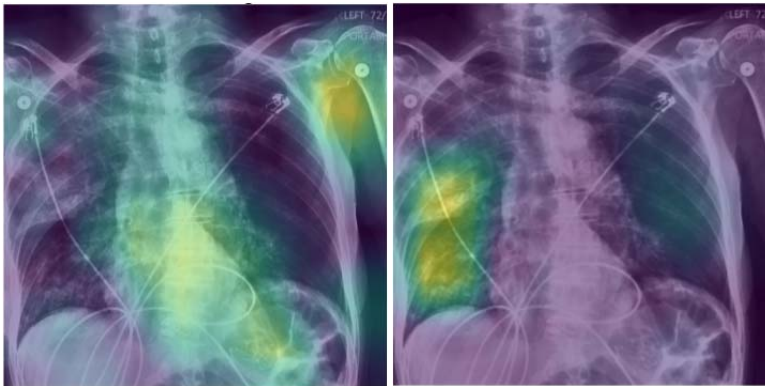
## 400 401 6.3 Summary of findings from clinical studies.

402  
403 The AI model developed in this study was based on prior work by Jabbour et al.<sup>20</sup> The  
404 study trained machine learning models to predict pneumonia, heart failure, and COPD  
405 using chest radiographs and clinical data from the electronic health record, and applied  
406 the models to an internal cohort at Michigan Medicine and an external cohort from Beth

407 Israel Deaconess Medical Center. They showed that a model combining chest  
408 radiographs and EHR data outperformed models based on each data modality alone and  
409 can accurately differentiate between common causes of acute respiratory failure.

410  
411 However, AI models trained on clinical data are known to make biased predictions due to  
412 spurious correlations present in training data.<sup>13</sup> For example, deep learning models can  
413 learn to predict patient age, sex, or BMI based on chest X-rays alone. If trained using  
414 standard approaches, these models could use these features in their predictions. To  
415 date, computational approaches can only partially mitigate the use of these shortcuts,  
416 and might fail in settings where the shortcut is not known in advance. However,  
417 heatmaps that highlight the regions of a chest X-ray that a model focuses on might signal  
418 that a model is taking a shortcut, such as highlighting the features of patient age (e.g.,  
419 osteoporosis) instead of clinically relevant features.

420



421  
422 Figure 4. Left: A model highlighting the features of patient age rather than clinically relevant  
423 features of disease. Right the model highlighting clinically relevant features in lungs.

424

#### 425 6.4 Summary of known and potential risks and benefits

426

427 Because this is a web-based survey study, there are minimal potential risks to patients  
428 and participants and no adverse events are expected during this study. The participant  
429 taking the survey is free to exit the survey at any point.

430

#### 431 6.5 Description and justification of route administration and dosage

432 Not applicable

433

#### 434 6.6 Dosages, dosage modifications and method of administration

435 Not applicable

436

437 6.7 Preparation and labelling of Investigational Medicinal Product

438 Not applicable

439

440 6.8 Drug accountability

441 Not applicable

442

## 443 7. METHODS

444

445

### 446 7.1 STUDY PARAMETERS/ENDPOINTS

447

#### 448 7.1.1 Main study parameters/endpoints

449

450 We compare participants when presented with AI model predictions to participants  
451 when presented with AI model predictions and explanations.

452

453 • Diagnostic accuracy for pneumonia, heart failure, and COPD

454 • Treatment accuracy in selecting antibiotics, diuretics, and steroids.

455

#### 456 7.2 Randomization, blinding, and treatment allocation

457

458 Participants were randomly shown 9 clinical vignettes. The first two vignettes are not  
459 accompanied by AI model predictions and are used to estimate baseline participant  
460 diagnostic accuracy. The next 6 vignettes include AI model predictions, but some  
461 participants are randomized to also see model explanations. These 6 vignettes include 3  
462 vignettes with standard model predictions and 3 with systematically biased AI model  
463 predictions. Participants are randomized to see one of three types of systematically  
464 biased AI model predictions, and the ordering of the 3 standard and 3 systematically  
465 biased model predictions was randomized. In the final vignette, all participants are  
466 provided a clinical consult, which is a short narrative provided by a hypothetical trusted  
467 colleague, who describes the rationale behind which diagnoses are most likely and what  
468 treatment plan they recommend. By design, the clinical consult always provided the  
469 correct diagnosis and appropriate treatment plan to provide a realistic upper bound of  
470 participant diagnostic accuracy.

471

#### 472 7.3 Study procedures

473



474 *Study population:*  
475 Hospitalist physicians, nurse practitioners, and physician assistants who commonly care  
476 for patients with acute respiratory failure.

477  
478 *Data collection:*  
479 Data collection will occur through Qualtrics, as it is approved for HIPAA data storage.  
480 Data collection starts when the participant clicks the survey link and stops when the  
481 participant either completes the survey, or two weeks after they exit the survey. This  
482 allows the participant to re-enter the survey to continue working on it. The survey data  
483 is saved and anonymized. It will be extracted and stored on HIPAA-aligned servers only  
484 accessible by the study team named at the University of Michigan. To preserve  
485 anonymity, participants will be redirected to another survey that is not linked to their  
486 responses to collect their contact information for payment purposes.

487  
488 **7.4 Withdrawal of individual subjects**  
489 Participants can exit the study at any time for any reason if they wish to do so without  
490 any consequences.

491 **7.4.1 Specific criteria for withdrawal (if applicable)**  
492 Not applicable

493  
494 **7.5 Replacement of individual subjects after withdrawal**  
495 Not applicable

496  
497 **7.6 Follow-up of subjects withdrawn from treatment**  
498 Not applicable

499  
500 **7.7 Premature termination of the study**  
501 We do not expect any serious adverse events directly related to this study. Therefore, we  
502 do not expect to have to terminate this study prematurely.

503

## 504 **8. SAFETY REPORTING**

505

506 **8.1 Temporary halt for reasons of subject safety**

507

508 **8.2 AEs, SAEs, SUSARs**

509 **8.2.1 Adverse events (AEs)**

510

511 Adverse events are defined as any undesirable experience occurring to a subject  
512 during the study, whether or not considered related to the intervention. All  
513 adverse events reported by the participants or observed by the investigator or  
514 their staff will be recorded.

515

#### 516 8.2.2 Serious adverse events (SAEs)

517

518

519 Due to the nature of this study, which was deemed as minimal risk, we will not  
520 be directly working with patients and do not anticipate any SAEs. However, **the**  
521 **investigator will report all SAEs to the sponsor without undue delay after**  
522 **obtaining knowledge of the events.**

523

524

#### 525 8.3 Annual safety report

526 Not applicable

527

#### 528 8.4 Follow-up of adverse events

529 All AEs will be followed until they have abated, or until a stable situation has been  
530 reached.

531

#### 532 8.5 Data Safety Monitoring Board (DSMB) / Safety Committee

533 Because the study was deemed to be minimal risk survey based study, a DSMB was not  
534 formed for the study.

535

536

## 537 9. STATISTICAL ANALYSIS

### 538 9.1 Primary study parameters/endpoints

- 539 • Diagnostic accuracy for pneumonia, heart failure, and COPD
- 540 • Treatment accuracy in selecting antibiotics, diuretics, and steroids.

541

### 542 9.2 Interim analysis (if applicable)

543 Not applicable

544

### 545 9.3 Statistical analysis plan

546

547 The study aims to recruit 400 participants based on a sample size to measure a decrease  
548 in accuracy of 25% with the systematically biased AI model compared to baseline and a  
549 10% improvement with the biased AI model with explanations compared to no  
550 explanations.

551  
552 Completed vignettes were included in the analysis regardless of whether a participant  
553 completed all 9 vignettes. Diagnostic accuracy and treatment decision accuracy will be  
554 compared using a generalized linear mixed-effects models, accounting for the nested  
555 data structure of repeated measures and controlling for individual-related variables,  
556 where individual diagnostic responses are nested within participants. After fitting the  
557 model, to aid in model interpretation, marginal effects and predictive margins will be  
558 reported. Statistical analyses will be performed in R. Statistical significance was based on  
559 a p-value < 0.05.  
560

## 561 10. ETHICAL CONSIDERATIONS

### 562 10.1 Regulation statement

563 This study has been approved by the UM IRB HUM00180745

### 564 10.2 Recruitment and consent

565 We will recruit hospitalist physicians, nurse practitioners, and physician assistants who  
566 commonly care for patients with ARF from US hospitals. We will identify hospitalist site  
567 champions who will send out an invitation email with the study information to  
568 hospitalist clinicians at their respective institutions. Consent will be obtained prior to  
569 participant randomization. Specifically, once participants click on the study link, they will  
570 be shown a page to screen for their eligibility. If eligible, they will be redirected to an  
571 introduction page that informs them about the study. This includes that the study will be  
572 completely anonymous, it will take 25 minutes to complete, and that they can stop the  
573 study at any time and come back to the point where they leave off. They are also told  
574 that some details of the study's purpose will be withheld and that they will receive a \$50  
575 Amazon.com gift card upon completion. If they agree to these terms, they can click  
576 forward and are then randomized. If not, they can click out of the survey at this, or any  
577 other point.  
578

### 579 10.3 Objection by minors or incapacitated subjects (if applicable)

580 Not applicable.

581 10.4 Benefits and risks assessment, group relatedness  
582 Not applicable.

583 10.5 Compensation for injury  
584 Not applicable.

585 10.6 Incentives (if applicable)  
586 Participants who complete the study will receive a \$50 amazon gift card.

## 587 11. ADMINISTRATIVE ASPECTS, MONITORING AND PUBLICATION

588 11.1 Handling and storage of data and documents  
589 The data will be handled confidentially. The participant responses will be retrieved from  
590 the Qualtrics interface, which is only accessible to a subset of the study team members.  
591 The data will then be stored on HIPAA approved servers at the University of Michigan for  
592 subsequent analyses. All data will be anonymized. Any publication arising from this study  
593 will not contain data that can be traced to a specific participant.

594  
595 11.2 Monitoring and Quality Assurance  
596

597 As data is collected throughout the study, it will be downloaded and checked to ensure  
598 that the randomization is set for each participant and all data is recorded as expected.  
599

600 11.3 Public disclosure and publication policy  
601 Results of this study will be submitted for publication in a peer reviewed scientific  
602 medical journal.  
603

## 604 12. AMENDMENTS

605

## 606 13. REFERENCES

- 607 1. Kempker JA, Abril MK, Chen Y, Kramer MR, Waller LA, Martin GS. The Epidemiology of  
608 Respiratory Failure in the United States 2002-2017: A Serial Cross-Sectional Study. *Crit Care*  
609 *Explor.* Jun 2020;2(6):e0128. doi:10.1097/cce.000000000000128  
610 2. Stefan MS, Shieh MS, Pekow PS, et al. Epidemiology and outcomes of acute respiratory  
611 failure in the United States, 2001 to 2009: a national survey. *J Hosp Med.* Feb 2013;8(2):76-82.  
612 doi:10.1002/jhm.2004

- 613 3. HCUP Fast Stats. Healthcare Cost and Utilization Project (HCUP). April 2021. Agency for  
614 Healthcare Research and Quality, Rockville, MD. [www.hcup-  
615 us.ahrq.gov/faststats/national/inpatientcommondiagnoses.jsp?year1=2018](http://www.hcup-<br/>615 us.ahrq.gov/faststats/national/inpatientcommondiagnoses.jsp?year1=2018).
- 616 4. Roberts E, Ludman AJ, Dworzynski K, et al. The diagnostic accuracy of the natriuretic  
617 peptides in heart failure: systematic review and diagnostic meta-analysis in the acute care  
618 setting. *BMJ*. Mar 2015;350:h910. doi:10.1136/bmj.h910
- 619 5. Lien CT, Gillespie ND, Struthers AD, McMurdo ME. Heart failure in frail elderly patients:  
620 diagnostic difficulties, co-morbidities, polypharmacy and treatment dilemmas. *Eur J Heart Fail*.  
621 Jan 2002;4(1):91-8. doi:10.1016/s1388-9842(01)00200-8
- 622 6. Daniels LB, Clopton P, Bhalla V, et al. How obesity affects the cut-points for B-type  
623 natriuretic peptide in the diagnosis of acute heart failure. Results from the Breathing Not  
624 Properly Multinational Study. *Am Heart J*. May 2006;151(5):999-1005.  
625 doi:10.1016/j.ahj.2005.10.011
- 626 7. Levitt JE, Vinayak AG, Gehlbach BK, et al. Diagnostic utility of B-type natriuretic peptide  
627 in critically ill patients with pulmonary edema: a prospective cohort study. *Crit Care*.  
628 2008;12(1):R3. doi:10.1186/cc6764
- 629 8. Zwaan L, Thijs A, Wagner C, van der Wal G, Timmermans DR. Relating faults in diagnostic  
630 reasoning with diagnostic errors and patient harm. *Acad Med*. Feb 2012;87(2):149-56.  
631 doi:10.1097/ACM.0b013e31823f71e6
- 632 9. National Academies of Sciences EaM. *Improving diagnosis in health care*. National  
633 Academies Press; 2015.
- 634 10. Tschandl P, Rinner C, Apalla Z, et al. Human-computer collaboration for skin cancer  
635 recognition. *Nature Medicine*. 2020;26(8):1229-1234.
- 636 11. Gulshan V, Peng L, Coram M, et al. Development and validation of a deep learning  
637 algorithm for detection of diabetic retinopathy in retinal fundus photographs. *Jama*.  
638 2016;316(22):2402-2410.
- 639 12. Van der Laak J, Litjens G, Ciompi F. Deep learning in histopathology: the path to the  
640 clinic. *Nature medicine*. 2021;27(5):775-784.
- 641 13. Jabbour S, Fouhey D, Kazerooni E, Sjoding MW, Wiens J. Deep Learning Applied to Chest  
642 X-Rays: Exploiting and Preventing Shortcuts. PMLR; 750-782.
- 643 14. Gichoya JW, Banerjee I, Bhimireddy AR, et al. AI recognition of patient race in medical  
644 imaging: a modelling study. *The Lancet Digital Health*. 2022;4(6):e406-e414.
- 645 15. Obermeyer Z, Powers B, Vogeli C, Mullainathan S. Dissecting racial bias in an algorithm  
646 used to manage the health of populations. *Science*. 2019;366(6464):447-453.
- 647 16. Gaube S, Suresh H, Raue M, et al. Do as AI say: susceptibility in deployment of clinical  
648 decision-aids. *NPJ digital medicine*. 2021;4(1):31.

- 649 17. Goddard K, Roudsari A, Wyatt JC. Automation bias: a systematic review of frequency,  
650 effect mediators, and mitigators. *Journal of the American Medical Informatics Association*.  
651 2012;19(1):121-127.
- 652 18. Selvaraju RR, Cogswell M, Das A, Vedantam R, Parikh D, Batra D. Grad-cam: Visual  
653 explanations from deep networks via gradient-based localization. 618-626.
- 654 19. Bhatt U, Xiang A, Sharma S, et al. Explainable machine learning in deployment. 648-657.
- 655 20. Jabbour S, Fouhey D, Kazerooni E, Wiens J, Sjoding MW. Combining chest X-rays and  
656 electronic health record (EHR) data using machine learning to diagnose acute respiratory failure.  
657 *Journal of the American Medical Informatics Association*. 2022;29(6):1060-1068.  
658