Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized Vignette-Based Multicenter Study

1

| TRIAL FULL TITLE | Measuring the Impact of AI in the Diagnosis of Hospitalized Patients through a Web-based Randomized Survey Vignette Multi-Center Study |
|---|---|
| SAP VERSION | 2 |
| SAP VERSION DATE | 6/22/2023 |
| TRIAL PRINCIPAL INVESTIGATOR | Jenna Wiens, PhD<br>Michael Sjoding, MD |
| SAP AUTHOR(s) | Sarah Jabbour, Michael Sjoding |

2
3
4
5
6
7
8
9
10

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized Vignette-Based Multicenter Study

**Table of Contents**

Statistical Analysis Plan

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized
Vignette-Based Multicenter Study

72

73

74
75
76
77
78

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized
Vignette-Based Multicenter Study

79 **Abbreviations and Definitions**

80

81

| AE | Adverse Event |
|---|---|
| ARF | Acute Respiratory Failure |
| CRF | Case Report Form |
| COPD | Chronic Obstructive Pulmonary Disease |
| IMP | Investigational Medical Product |
| SAP | Statistical Analysis Plan |

82

83 **1    Introduction**

84 **1.1    Preface**

85 Artificial intelligence (AI) has achieved high accuracy at identifying abnormalities in clinical images,

86 such as pneumonia from chest radiographs, diabetic retinopathy from fundus images, or skin cancer

87 from histopathology images.[10-12] However, systematic bias in AI models can lead to inaccurate

88 predictions for entire subpopulations.[13-15] When presented with such incorrect predictions, physician

89 performance can be harmed[16] due to automation bias,[17] which is especially concerning in safety-

90 critical settings. Thus, the extent to which AI can be safely integrated into clinical workflows and to

91 support diagnostic decisions is still unknown.

92

93 This study aims to provide insight into the effectiveness of providing clinicians with image-based AI

94 model explanations to help them catch when models are making incorrect decisions.

95 **1.2    Scope of the analyses**

96 These analyses will primarily assess the extent to which showing clinicians systematically biased AI

97 model predictions and explanations improves their diagnostic accuracy after reviewing clinical

98 vignettes of patients with acute respiratory failure and determining the patient's likely diagnosis

99 compared to the setting where clinicians are shown biased AI model predictions without

100 explanations.

101 **2    Study Objectives and Endpoints**

102 **2.1    Study Objectives**

103 <u>Survey Data Collection Phase</u>

104 *Objectives*

105    • To determine clinician accuracy in diagnosing pneumonia, heart failure, and chronic

106       obstructive pulmonary disease (COPD) after reviewing clinical vignettes of patients with

107       acute respiratory failure (ARF) without any AI model input.

108    • To determine how AI model predictions without explanations affect clinician accuracy in

109       diagnosing pneumonia, heart failure, and COPD in a patient population with ARF.

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized
Vignette-Based Multicenter Study

110   • To determine how standard AI model predictions with explanations affect clinician accuracy
111      in diagnosing pneumonia, heart failure, and COPD in a patient population with ARF.
112   • To determine how intentionally biased AI model predictions without explanations affect
113      clinician accuracy in diagnosing pneumonia, heart failure, and COPD in a patient population
114      with ARF?
115   • To determine how do intentionally biased AI model predictions with explanations affect
116      clinician accuracy in diagnosing pneumonia, heart failure, and COPD in a patient population
117      with ARF.

118   **2.2   Endpoints**
119   *Primary endpoints*
120   • Clinician diagnostic accuracy for identify the cause of ARF after reviewing clinical vignettes
121      during following settings
122      o Clinicians provided no AI model predictions
123      o Clinicians provided standard AI model predictions without explanations
124      o Clinicians provided standard AI model predictions with explanations
125      o Clinicians provided biased AI model predictions without explanations
126      o Clinicians provided biased AI model predictions with explanations
127
128   *Secondary endpoints*
129   • Accuracy of treatment selection in the above settings

130   **3   Study Methods**

131   **3.1   General Study Design and Plan**
132
133   • Study configuration and experimental design: This study is a block randomized web-based
134      survey clinical vignette study
135   • Type of Comparison: Clinician diagnostic accuracy when provided AI model predictions with
136      explanations versus AI model without explanation
137   • Type of control(s): no AI model, AI model with predictions alone.
138   • Level and method of blinding (e.g. double-blind): Single blind study (clinicians are unaware
139      they are randomized to see AI model with or without predictions)
140   • Method of treatment assignment: Survey participant level randomization
141   • At what point in time subjects are randomized relative to treatments, events and study
142      periods: Participants are randomized after survey initiation
143   • Sequence and duration of all study periods: The survey is anticipated to take an average of
144      20 minutes to complete.
145

146   **3.2   Inclusion-Exclusion Criteria and General Study Population**
147   (ICH E3;9.3. ICH E9;2.2.1)
148
149   *Inclusion Criteria*

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized
Vignette-Based Multicenter Study

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

151

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

153  • Nurse Practitioner (NP)

154  • Physician Assistant

155  • Resident

156  • Fellow

157  • Attending Physician

158

159  **3.3    Randomization and Blinding**

160

161  **Overview**

162



163

164  **Figure 1. Survey flow and randomization** After confirming study eligibility and consent, participants
165  will complete two baseline clinical vignettes where they review patient clinical data and then
166  determine whether the patient has heart failure, pneumonia, and/or COPD without any AI model
167  predictions (Vignettes 1-2). All participants are then randomized to (1; green arrows) AI model
168  predictions with or without model explanations, (2; orange arrows) one of three types of biased AI
169  models shown (biased based on age, BMI, or preprocessing features), and (3; purple boxes) the
170  ordering of the 6 clinical vignettes where 3 standard model predictions and 3 systematically biased
171  model predictions were shown with the clinical vignette. All participants are shown a ninth vignette
172  (vignette 9), which features a clinical consult. The clinical consult provides includes a short block of
173  text providing a prediction and explanation for the patient's likely diagnosis from a hypothetical
174  trusted colleague.

175

176  Details of block randomization.

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized
Vignette-Based Multicenter Study

178 Block randomization was used to determine the specific patient vignettes and the order of vignettes
179 that subjects would see during the survey. Block randomization was performed in blocks of 90 to
180 achieve all three randomizations described in Figure 1. This ensured that all 45 patient vignettes
181 would be evenly assigned across the first two baseline vignettes and the last clinical consult vignette,
182 and to ensure that a subject would only see a clinical vignette once during the survey. Within the
183 blocks of 90, 30 subjects were randomly assigned to each of the three AI model bias types (age, BMI,
184 or preprocessing features). There were 6 specific clinical vignettes where the AI model displayed the
185 biased behavior for each bias type. Therefore, for the 30 subjects randomly assigned to a specific
186 bias types, 3 of the 6 specific vignettes where the AI model displayed the specific biased behavior
187 was randomly assigned to the subject. An additional 3 vignettes from the 45 total vignettes were
188 randomly selected to be shown with standard model predictions. These 6 patient vignettes (3 with
189 standard AI model, 3 with biased AI model) were then displayed in random order. After the
190 randomization blocks of 90 subjects were generated, carefully tested was performed to ensure all
191 specifications were met.

### 3.4 Study Assessments

195 The study is designed to take on average 20 minutes per participant. The participant can exit out of
196 the survey and return within two weeks to continue. After the two weeks, the survey is closed and
197 the participant can no longer continue the survey.

199 Participants will be asked to rate the independent likelihoods that pneumonia, heart failure, and
200 COPD are contributing to the patient's ARF on a scale of 0-100. They will be instructed that patients
201 can have one, more than one, or none of these diagnoses. Clinician diagnostic accuracy will be
202 determined by comparing their answer to an independent assessment of each patients likely
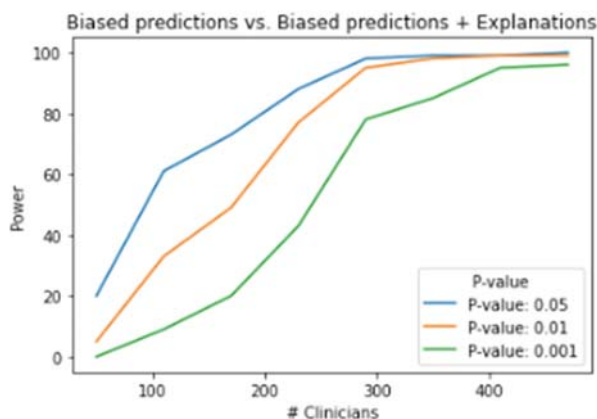203 diagnosis performed by an panel of clinician reviewers.

### 4 Sample Size

206 The sample size calculation is based on the primary endpoint of clinician diagnostic accuracy for
207 pneumonia, heart failure, and COPD. We performed sample size calculations to ensure we would
208 have adequate power to detect both a reduction in diagnostic accuracy when clinicians were shown
209 a biased model, assuming they would follow the biased model's recommendations 50% of the time,
210 and adequate power to detect an improvement in accuracy when clinicians were shown a biased
211 model and explanations, assuming they would follow the biased model recommendation 25% of the
212 time when also shown the explanation. These assumptions would translate into decrease in
213 diagnostic accuracy by 20% when clinicians were shown a biased model and a 10% improvement
214 when shown a biased model and explanation. We used a generalized linear mixed model with a
215 0.001 significance level. Given the simulated data generated as further described below, we fit a
216 generalized linear mixed model in R to measure if the recovery of clinician diagnostic accuracy when
217 shown the model explanation was significantly different compared to the clinician diagnostic
218 accuracy when shown a biased model alone. We performed 100 simulated studies at each sample
219 size level, and calculated power as the percentage of time a statistically significant difference was
220 measured. We found that the study would have very high power to detect a difference in diagnostic
221 accuracy when comparing clinician baseline diagnostic accuracy and clinician accuracy when shown a

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized Vignette-Based Multicenter Study

222 biased AI model. The power calculation illustrated in the figure below describes the sample size
223 needed to detect a difference in diagnostic accuracy when clinicians shown a biased AI model alone
224 and when clinicians are shown a biased AI model with explanation.
225



Biased predictions vs. Biased predictions + Explanations

226
227
228 **Detailed sample size calculation:**
229
230 For our simulation, we model the likelihood that a study subject gets a diagnosis correct as a
231 combination of their baseline diagnostic accuracy ($b$), the difficulty of the patient case ($d$), the skill
232 of the clinician ($c$), and the effect of either being shown an AI model prediction alone ($\beta_1$) or being
233 shown an AI prediction with an explanation ($\beta_2$), where $AI\_Alone$ and $AI\_Explanation$ are
234 indicator variables and $\sigma(.)$ denotes the sigmoid function. Details of each of the variables
235 represented in the equation are described in more detail below.
236
237 $$p = sigmoid(b + d + c + \beta_1 AI\_Alone + \beta_2 AI\_Explanation)$$
238
239 Then, during the simulation, whether a clinician obtains the correct diagnosis is determined by
240 drawing a random variable from a Bernoulli distribution of probability p, with probability determined
241 based on the above data generation model.
242
243 Sample size for the study was determined by performing by 100 simulations at participant sample
244 sizes of 50 to 550, in increments of 50, using the above equation to model the data generating
245 process in the survey. In each simulation, a clinician is shown 2 vignettes with no AI model input and
246 then shown either 3 vignettes with systematically biased AI model without explanations or 3
247 vignettes with systematically biased AI model with explanations. We simulate 100 of these studies at
248 each sample size level.
249
250 During a simulated study, we generate blocks of vignettes to assign to hypothetical subjects by the
251 combinations of (1) whether they were shown an AI model explanation, and (2) the type of
252 systematic bias seen. This generates 6 possible assignments:
253
254    1. AI model with Pneumonia bias, no explanation
255    2. AI model with Pneumonia bias, explanation
256    3. AI model with Heart failure bias, no explanation

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized Vignette-Based Multicenter Study

257    4.  AI model with Heart failure bias, explanation
258    5.  AI model with COPD bias, no explanation
259    6.  AI model COPD bias, explanation
260
261    For every hypothetical clinician in the study, we then assign them one of the above conditions (in
262    order) and then generate the likelihood that the participant gets the diagnosis correct based on the
263    data generation model. For example, the first clinician (assigned to 1), is shown 2 clinical vignettes
264    without an AI model and 3 vignettes of a biased AI model with pneumonia bias and no explanation.
265    The clinician's diagnostic accuracy for each of these clinical vignettes was determined using the data
266    generating model.
267
268    Details of the generative model parameters:
269
270    1. Baseline diagnostic accuracy
271
272    Average baseline diagnostic accuracy for clinicians was assumed to be 0.7 across all three diagnoses
273    but then updated after calculating baseline accuracy at an interim analysis (see 8.5). Accuracy at the
274    interim analysis was determined to be:
275         a.  Pneumonia: 0.68
276         b.  Heart Failure: 0.72
277         c.  COPD: 0.82
278    These probabilities are transformed to log odds for the data generation model, i.e., logit(x).
279
280    •  If the participant is randomized to see the Pneumonia bias, then
281    $$b = logit(0.68) = 0.75$$
282
283    •  If the participant is randomized to see the heart failure bias, then
284    $$b = logit(0.72) = 0.94$$
285
286    •  If the participant is randomized to see the COPD bias, then
287    $$b = logit(0.82) = 1.5$$
288
289    2. Draw Clinician skill $c_i$
290
291    We assumed variation in clinician skill was a normally distributed random variable with mean
292    $\mu_{clinician} = 0$. We assumed the best clinician, who was 2 standard deviations better than average
293    clinician, got the average case right 90% of the time, then $\sigma_{clinician} = \frac{logit(t) - logit(x)}{2}$.
294
295    For each clinician $c_i$, their skill level is drawn:
296
297    $$c_i \sim N(0, \sigma_{clinician}) \text{ for i = 1,2,…, n; where n is the number of clinicians in the simulation}$$
298
299
300    3. Draw clinical vignette simplicity $d_j$
301    We assumed variation exists in clinical vignette diagnostic difficulty, such that cases that are 1 std.
302    easier to diagnosis than the average vignette are answered correctly 90% of the time. Case
303    diagnostic difficulty was assumed to be a normally distributed random variable with mean $\mu_{case} = 0$

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized Vignette-Based Multicenter Study

304 and std. $\sigma_{case} = logit(z) - logit(x)$ , where z = 0.9, and x = {0.68, 0.72, and 0.82} for the three
305 diagnoses.
306
307       For each case $d_j$, the case simplicity is drawn:
308
309 $$d_j \sim N(0, \sigma_{case})$$
310
311
312 <u>3. Draw the Impact of a systematically biased AI model</u>
313 We assumed that a systematically biased AI model prediction would reduce clinician diagnostic
314 accuracy by $a$. Therefore, if average diagnostic accuracy was x, the impact of a systematically biased
315 AI model on accuracy is $(x - a)\%$. We assumed participants would listen to the AI model 50% of the
316 time, which meant that x - $a = x * 0.5 + 0.33 * 0.5$. In the data generation model, an indicator
317 variable was included indicating whether clinicians were shown a systematically biased AI model
318 with coefficient $\beta_1$ . When $\beta_1 < 0$, this variable represents a decrease in the likelihood that the
319 clinician will get a case correct. It is the difference between the likelihood that the participant gets
320 the case correct with the AI input minus the likelihood that the participant gets the case correct
321 without AI model input: $\beta_1 = logit(x - a) - logit(x)$.
322
323 In the simulation, if the participant was shown the AI model for the vignette, then
324
325     •   If the participant is randomized to see the Pneumonia bias, then
326           $\beta_1 = logit(0.68 * .5 + 0.33 * 0.5) - logit(0.68) = -0.73$
327
328     •   If the participant is randomized to see the heart failure bias, then
329           $\beta_1 = logit(0.72 * .5 + 0.33 * 0.5) - logit(0.72) = -0.84$
330
331     •   If the participant is randomized to see the COPD bias, then
332           $\beta_1 = logit(0.82 * .5 + 0.33 * 0.5) - logit(0.82) = -1.2$
333
334 <u>Impact of a systematically biased AI model with explanation</u>
335 We assumed that providing an AI model explanation helps clinicians recover diagnostic accuracy by $r$
336 when shown a biased AI model that reduce their accuracy by $a$. Therefore, if accuracy on an average
337 case was x, the impact of showing a biased AI model and explanation is ((x-a) + r)%. We assumed
338 that participants would recover 50% back to their baseline diagnostic accuracy, which means ((x-a) +
339 r)% = (x − a + 0.5a)% = $x * 0.75 + 0.33 * 0.25$. In the data generation model, an indicator variable
340 was included indicating whether clinicians were shown a systematically biased AI model explanation
341 with coefficient $\beta_2$. $\beta_2$ represents the change in likelihood that the clinician gets the case correct.
342 When $\beta_2 > 0$, this variable represents an increase in the likelihood that the clinician gets the case
343 correct: $\beta_2 = logit(x - a + r) - logit(x - a)$.
344
345 In the simulation, if the participant was shown the AI model for the vignette, then
346
347     •   If the participant is randomized to see the Pneumonia bias, then
348           $\beta_2 = logit(0.68 - 0.25 + 0.1) - logit(0.68) = -0.38$
349
350     •   If the participant is randomized to see the heart failure bias, then

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized Vignette-Based Multicenter Study

351 $$\beta_1 = logit(0.72 - 0.25 + 0.1) - logit(0.68) = -0.44$$
352
353 - If the participant is randomized to see the COPD bias, then
354 $$\beta_1 = logit(0.82 - 0.24 + 0.1) - logit(0.68) = -0.68$$
355

356 ## 5    General Analysis Considerations

357 ### 5.1    Timing of Analyses
358 The final analysis will be performed two weeks after the last study email invitation is sent out and
359 400 participants have completed the study.
360

361 ### 5.2    Analysis Populations
362

363 #### 5.2.1    Full Analysis Population (or Intention to Treat or Modified Intention to Treat)
364 - *All subjects who consent to taking the study and click to the first page of the study. Each*
365 *participant who does so is randomized.*

366 #### 5.2.2    Per Protocol Population
367 - *NA*

368 #### 5.2.3    Safety Population
369 - *NA*
370

371 ### 5.3    Covariates and Subgroups
372 Because all participants are randomized approximately equally across treatment groups, we do not
373 anticipant any covariates that will have an importance influence on our primary endpoints.
374 There are no *a priori* hypotheses of subgroup differences.
375

376 #### 5.3.1    Multi-center Studies
377 This is a multi-center study, where participant responses will be pooled from all centers. The rational
378 behind this is that we assume there is no meaningful center differences in treating patients with
379 ARF.
380

381 ### 5.4    Missing Data
382 The main source of missing data will be missing demographic information in participants who do not
383 complete all vignettes and the demographic questions after the survey. We assume this data will be
384 missing at random. Because this demographic information is not included as covariates in any of the
385 analysis, we do not plan to do anything to impute missing demographics data.
386

387 ### 5.5    Interim Analyses and Data Monitoring (as applicable)

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized
Vignette-Based Multicenter Study

388 **5.5.1 Purpose of Interim Analyses**

389 No interim analyses of the exposure variables will be conducted (ie the impact of AI models on
390 clinician diagnostic accuracy), however, we will measure participate baseline accuracy to confirm our
391 sample size calculations.

392 **5.5.2 Planned Schedule of Interim Analyses**

393 Participant baseline diagnostic accuracy will be measured after 300 participants are enrolled in the
394 study.

395 **5.5.3 Scope of Adaptations**

396 Not applicable.

397 **5.5.4 Stopping Rules**

398 Not applicable.

399 **5.5.5 Analysis Methods to Minimize Bias**

400 Not applicable.

401 **5.5.6 Adjustment of Confidence Intervals and p-values**

402 Not applicable.

403 **5.5.7 Interim Analysis for Sample Size Adjustment**

404 Once 300 participant responses are collected, we will measure participant baseline diagnostic
405 accuracy to confirm our baseline accuracy assumption for sample size calculations. We will not
406 measure the effects of the exposures (e.g. AI model predictions and explanations) on diagnostic
407 accuracy during the interim analysis.

408 **5.5.8 Practical Measures to Minimize Bias**

409 The study team will conduct the interim analysis to measure baseline diagnostic accuracy and will
410 not measure nor change any treatment effect assumptions in the power calculations.

411 **5.5.9 Documentation of Interim Analyses**

412

413 Data and results of the interim analysis will be stored on the HIPAA aligned compute servers of the
414 study team members.

415

416 **5.6 Multiple Testing**

417 We do not plan to perform any corrections for multiple testing in our primary endpoint of clinical
418 diagnostic accuracy across settings.
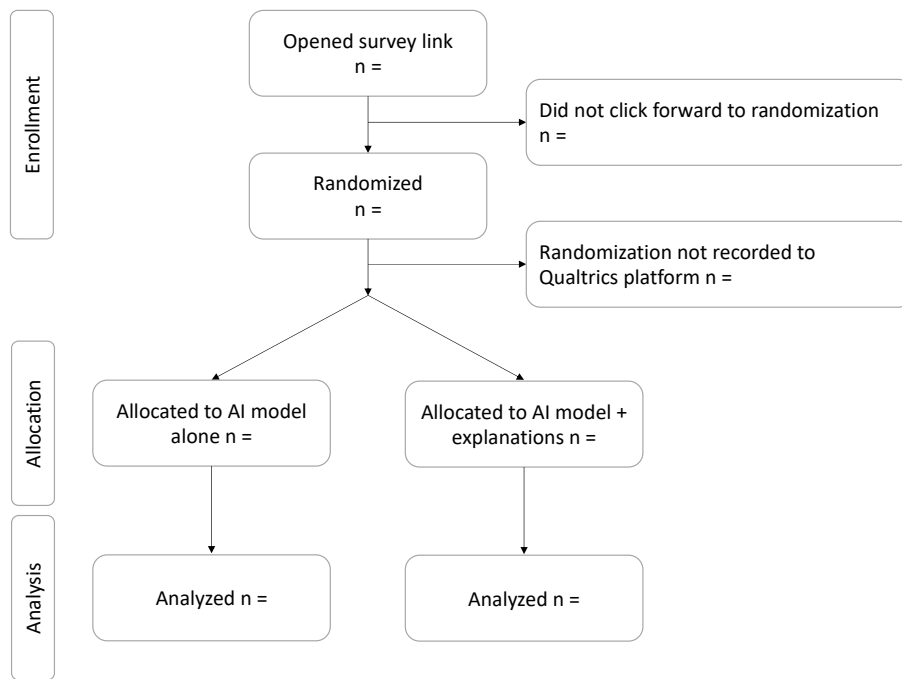
419

420 **6 Summary of Study Data**

421 The tables and figures will be based upon the full population of participants who are randomized in
422 the study and completed at least once clinical vignette. The first table will include summary statistics
423 of all study subjects, where each column represents the two treatment arms (AI Model Alone, AI
424 Model + Explanation). The primary statistical analysis results from generalized linear mixed models
425 will also be reported table format, with each row corresponding to diagnostic performance in each
426 vignette setting: Clinician Baseline, Clinician + Standard Model, Clinician + Standard Model +

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized
Vignette-Based Multicenter Study

427  Explanation, Clinician + Systematically Biased Model, Clinician + Systematically Biased Model +
428  Explanation, Clinician + Clinical Consult.

429  **6.1    Subject Disposition**

430  We will track 1) how many subjects open the survey link through an email as "Opened Survey Link,"
431  2) how many met the inclusion criteria and consented to study participation and are "randomized,"
432  3) how many randomization failures occurred because of Qualtrics platform errors, 4) how many
433  were allocated to each treatment arm as "allocated to AI model alone" or "allocated to AI model +
434  explanation," 5) how many in each arm dropped out before completing a vignette, 6) how many
435  completed at least one vignette and were "analyzed."
436  section 9."

437



438

439  **6.2    Derived variables**

440  Participant diagnostic accuracy is the primary endpoint of this study. Their responses will be
441  collected on a scale of 0-100 and responses above 50 were considered positive for each diagnosis.
442  To calculate diagnostic accuracy, this response will be compared to the reference standard labels
443  generated by a group of 5 physicians who reviewed the patients complete medical record and
444  determined the patient's diagnosis.

445  **6.3    Protocol Deviations**

446  We do not anticipate any major protocol deviations that would impact the analysis.

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized Vignette-Based Multicenter Study

447 **6.4    Demographic and Baseline Variables**

448 Collection of participant demographic information is optional and will occur after participants
449 complete all vignettes. We will collect participant age, race and ethnicity, gender, the hospital
450 setting they primarily work, their general practice area, their current role on their healthcare team,
451 and when they completed their medical training.

452 **6.5    Concurrent Illnesses and Medical Conditions**

453 Not applicable

454 **6.6    Treatment Compliance**

455 Not applicable

456 **7    Efficacy Analyses**

457

458 **7.1    Primary Efficacy Analysis**

459 The primary outcome is the participant's diagnostic accuracy after reviewing the patient vignette.
460 Participants will separately assess whether the patient in the vignette has pneumonia, heart failure,
461 and COPD, and their diagnostic accuracy for each will be analyzed as a unique response within the
462 generalized linear mixed model, with individual responses nested within study participant. To
463 determine diagnostic accuracy, participant responses will be compared to the reference standard
464 labels generated by a group of 5 physicians who reviewed the patients complete medical record.  A
465 generalized linear model with logit link will be fit for diagnostic accuracy with indicator variables for
466 each of the 5 settings evaluated (clinician baseline without AI model, standard model, standard
467 model with explanation, biased model, biased model with explanation). After fitting the model, we
468 will specifically compare diagnostic accuracy for the following settings:

469 - Baseline participant accuracy with no AI prediction model input (Clinician Baseline)
470    compared to participant accuracy with standard model predictions without explanations
471    (Clinician + Standard Model)
472 - Baseline participant accuracy with no AI prediction model input (Clinician Baseline)
473    compared to participant accuracy with standard model predictions and explanations
474    (Clinician + Standard Model + Explanations)
475 - Baseline participant accuracy with no AI prediction model input (Clinician Baseline)
476    compared to participant accuracy when systematically biased model predictions are
477    provided without explanations (Clinician + Systematically Biased Model)
478 - Baseline participant accuracy with no AI prediction model input (Clinician Baseline)
479    compared to participant accuracy when systematically biased model predictions are
480    provided with explanations (Clinician + Systematically Biased Model + Explanations)
481 - Participant accuracy when systematically biased model predictions are provided without
482    explanations (Clinician + Systematically Biased Model) compared to participant accuracy
483    when systematically biased model predictions are provided without explanations (Clinician +
484    Systematically Biased Model)

485 **7.2    Secondary Efficacy Analyses**

486 In a secondary analysis, we will examine how treatment decisions are influenced by correct or
487 incorrect model predictions. We measured the percentage of time participants made an appropriate
488 treatment decision across settings ('Clinician Baseline', 'Clinician + Model', 'Clinician + Model +

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized
Vignette-Based Multicenter Study

489 Explanation') for both standard and biased AI models. Appropriate treatment for each vignette is
490 determined based on the patients' reference diagnoses and review of the patients complete medical
491 record. We will also investigate the effects of systematically biased estimates on the distributions of
492 participant responses.

493 **7.2.1   Secondary Analyses of Primary Efficacy Endpoint**
494 Not applicable

495 **7.2.2   Analyses of Secondary Endpoints**
496 Not applicable

497 **7.3    Exploratory Efficacy Analyses**
498 Not applicable


499 **8     Safety Analyses**
500 Because this vignette survey study was deemed minimal risk, no safety analysis will be conducted

501 **8.1    Extent of Exposure**
502 Not applicable

503 **8.2    Adverse Events**
504 Not applicable

505 **8.3    Deaths, Serious Adverse Events and other Significant Adverse Events**
506 Not applicable

507 **8.4    Pregnancies (As applicable)**
508 No applicable

509 **8.5    Clinical Laboratory Evaluations**
510 Not applicable

511 **8.6    Prior and Concurrent Medications (As applicable)**
512 Not applicable

513 **8.7    Other Safety Measures**
514 Not applicable


515 **9     Pharmacokinetics (As Applicable)**
516 Not applicable
517


518 **10   Other Analyses**
519 Not applicable

Measuring the Impact of AI in the Diagnosis of Hospitalized Patients Through a Randomized Vignette-Based Multicenter Study

520 **11    Reporting Conventions**

521

522 P-values less than .001 will be reported as "p-value<.001"; P-values between .001 and .01 will be
523 reported to the nearest thousandth. P-values greater than or equal to .01 will be reported to the
524 nearest hundredth; P-values greater than .99 will be reported as "p-value>.99."


525 **12    Quality Assurance of Statistical Programming (As Applicable)**

526 All statistical analysis will be conducted in R by the first author team member. A second study team
527 member will review the R code to check for correctness, while also double checking the primary
528 analysis in Stata.


529 **13    References**

530 none

531

532