

Combining Rapid Antigen Testing and Syndromic Data Improves Sensitivity and Specificity in Real-World COVID-19 Detection

Fergus J Chadwick^{a,b}, Yacob Haddou^{a,b}, Tasnuva Chowdhury^a, David Pascall^c,
Shayan Chowdhury^e, Jessica Clark^{a,b}, Joanna Andrecka^f, Mikolaj
Kundergorski^{d,b}, Craig Wilkie^{d,b}, Eric Brum^f, Tahmina Shirin^g, A S M
Alamgir^g, Mahbubur Rahman^g, Ahmed Nawsher Alam^g, Farzana Khan^g, Janine
Illian^{d,b}, Ben Swallow^{d,b}, Davina L Hill^{a,b}, Dirk Husmeier^d, Jason
Matthiopoulos^{a,b}, Katie Hampson^{a,b}, Ayesha Sania^h

^a*Institute of Biodiversity, Animal Health and Comparative Medicine, University of Glasgow*

^b*COVID-19 in LMICs Research Group, University of Glasgow*

^c*MRC Biostatistics Unit, University of Cambridge*

^d*School of Mathematics and Statistics, University of Glasgow*

^e*a2i Programme, ICT Ministry/UNDP Bangladesh*

^f*UN FAO in support of the UN Interagency Support Team, Bangladesh*

^g*Institute of Epidemiology, Disease Control and Research, Ministry of Health, Bangladesh*

^h*Division of Developmental Neuroscience, Department of Psychiatry, Columbia University*

1. Abstract (Max 250 Words - Currently over)

Background

The majority of the world's population live in low- and middle-income countries (LMICs) where access to gold-standard diagnostics like RT-PCR is often limited. Rapid antigen testing (RAT) and syndromic diagnosis are two alternative, inexpensive and easy-to-deploy surveillance methods but there are concerns that they lack the sensitivity and specificity to effectively guide practice.

Methods

Bangladesh's Institute of Epidemiology Disease Control And Research (IEDCR) identified potential COVID-19 patients in Dhaka using syndromic surveillance.

*Corresponding Author

Email addresses: f.chadwick.1@research.gla.ac.uk (Fergus J Chadwick), yacob.haddou@glasgow.ac.uk (Yacob Haddou), tasnuvachowdhury2004@gmail.com (Tasnuva Chowdhury), david.pascall@mrc-bsu.cam.ac.uk (David Pascall), shayan.chowdhury@a2i.gov.bd (Shayan Chowdhury), Jessica.Clark@glasgow.ac.uk (Jessica Clark), aandrecka@gmail.com (Joanna Andrecka), mikolaj.kundergorski@gmail.com (Mikolaj Kundergorski), craig.wilkie@glasgow.ac.uk (Craig Wilkie), eric.brum@fao.org (Eric Brum), tahmina.shirin14@gmail.com (Tahmina Shirin), aalamgir@gmail.com (A S M Alamgir), dr_mahbub@yahoo.com (Mahbubur Rahman), anawsher@yahoo.com (Ahmed Nawsher Alam), farzanakhan_25@yahoo.com (Farzana Khan), janine.illian@glasgow.ac.uk (Janine Illian), ben.swallow@glasgow.ac.uk (Ben Swallow), davina.hill@glasgow.ac.uk (Davina L Hill), dirk.husmeier@glasgow.ac.uk (Dirk Husmeier), jason.matthiopoulos@glasgow.ac.uk (Jason Matthiopoulos), katie.hampson@glasgow.ac.uk (Katie Hampson), ays328@mail.harvard.edu (Ayesha Sania)

28 A sample ($n = 1172$) of these patients was tested using RAT and syndromic data
29 were collected. Models were fit to predict RT-PCR status using the RAT data,
30 the syndromic data, and the two combined. Model performance was measured
31 using predictive power and classification performance under three epidemiological
32 scenarios: “Agnostic,” “Rising Cases” and “Low-Level Cases.”

33 *Findings*

34 Combined data models yielded equal or improved performance over syndromic-
35 and RAT-only models across all three epidemiological scenarios and when com-
36 pared as more generic prediction and classification engines. In the “Rising Cases”
37 scenario, which most closely represents the current situation in many LMICs,
38 the combined data model false negative rate is 26 percentage points lower than that of
39 the RAT only model. Although the syndromic only model matches the combined
40 models false negative rate, its false positive rate is 31 percentage points higher.

41 *Interpretation*

42 A few accurate tests may be less useful at the population level than many
43 more imperfect ones. Small, scalable improvements in the accuracy of mass-
44 deployed but imperfect tests can then make a very big difference for pandemic
45 control.

46 We demonstrate that such improvements can be achieved by statistically utilising
47 complementary strengths and weaknesses across two imperfect diagnostics, we
48 can greatly improve the detection of COVID-19.

49 *Funding*

50 The Bill and Melinda Gates Foundation and the Wellcome Trust.

51 **2. Introduction (~1107 Words)**

52 Identification and isolation of COVID-19 cases remains key to the pandemic
53 response across the globe. The faster and more accurately we can identify cases,
54 the more effectively we can provide clinical care, reduce transmission of infection
55 and develop population-level interventions. RT-PCR testing has rapidly become
56 the default, gold-standard test for COVID-19 in applied settings due to its high
57 sensitivity and specificity for COVID-19 [2]. Most of the world’s population,
58 however, live in low- and middle-income countries (LMICs) where the laboratory
59 facilities needed to carry out RT-PCR tests are often scarce and hard to reach
60 [4]. COVID-19 diagnosis worldwide, therefore, must be made accessible using
61 inexpensive methods that can be carried out locally [6].

62 An increasingly popular alternative to RT-PCR is rapid antigen testing
63 (RAT) [7]. Like RT-PCR, these tests have high specificity for COVID-19 while
64 being less expensive, easier to implement, and faster to produce results [8].
65 RATs also require less commitment and discomfort for patients. For RT-PCR
66 testing, patients must travel to a designated site (such as a hospital or testing
67 booth) or have highly visible PPE-clad officials visit their home. Then, invasive
68 nasopharyngeal swabs must be taken and there is a delay in receiving the result
69 (between one day and a week in Bangladesh). In contrast, RAT can be conducted
70 on nasal or saliva samples, completed in the home with minimal PPE and results

71 are available in 30 minutes. RATs can be taken by persons with limited training,
 72 thus decreasing the time and expense associated with identifying cases. Together,
 73 these traits make RATs an appealing alternative to RT-PCR. However, several
 74 concerns have been raised about the sensitivity of RAT [9] leading to more false
 75 negative diagnoses.

76 Another alternative to RT-PCR, one that has been used since the start of
 77 the pandemic, is identifying cases through symptom-thresholding [10]. In this
 78 approach, a patient presenting with a fever and one or more viral pneumonia
 79 symptoms is treated as a COVID-19 positive patient. The main advantage
 80 of this approach is the ease of implementation. As with RAT the process is
 81 faster, cheaper and less invasive than RT-PCR, but unlike RAT the process
 82 relies on minimal equipment and thus can be scaled quickly and easily. For
 83 example, in Bangladesh, an LMIC, much of the initial support and reporting of
 84 infections locally is provided by community support teams (CSTs) composed of
 85 local volunteers with basic training. The CSTs can easily collect symptomatic
 86 data in the community and provide care where the thresholds are met. However,
 87 these thresholds were developed early in the outbreak, and thus were necessarily
 88 drawn from clinical intuition, rather than data, and for different variants and
 89 populations than they are now applied to. Consequently, the relationship between
 90 the thresholds and the true COVID-19 status is often weak, with low specificity
 91 leading to a very large number of false positive diagnoses.

92 A natural extension to these symptom-threshold approaches is syndromic
 93 modelling. Here, a patient presenting with a fever and one or more viral
 94 pneumonia symptoms is treated as a potential COVID-19 patient. However,
 95 rather than using a set of pre-determined criteria, a range of symptomatic and
 96 risk factor data are collected and then a sub-sample of patients is tested using
 97 RT-PCR for COVID-19 [11]. These data are used to fit a model that allows
 98 more accurate prediction of how likely a patient is to have COVID-19 through
 99 the identification of COVID-19 syndromes [13]. It is worth highlighting at this
 100 point that in resource-limited settings, there is very limited provision for testing
 101 of asymptomatic cases, despite their important role in disease transmission [14].
 102 Even while focusing solely on symptomatic patients, syndromic modelling is a
 103 complex and nuanced task. The strength of relationships between symptoms
 104 and diseases is not stable through time or across sampling strategies since the
 105 relative importance of each symptom for disease diagnosis, in part, depends on
 106 the prevalence of other diseases causing similar symptoms in the community [15].
 107 For example, if another disease for which loss of smell is a symptom becomes
 108 common, that symptom becomes a worse predictor for COVID-19. Similarly,
 109 if everyone who presents has a cough and thus is included in the sample, then
 110 coughing will likely have a very low correlation with COVID-19 (even if the
 111 two are strongly related in the general population). Symptoms are also inter-
 112 related, meaning that they cannot be interpreted independently. The majority of
 113 methods used currently do not account for these changes through time, symptom-
 114 to-symptom correlations or the relationship between the population sampled
 115 and the target population. Even then, the many types of common respiratory
 116 disease generally means that even then these models tend to have relatively high

JMC: Not
an official
term

JMC: Need
to highlight
that we are
only dealing
with symp-
tomatic pa-
tients, this
felt a good
place to do
it

JMC: query
cross-
sectional
sampling

117 false positive rates (low specificity) for COVID-19 [15], although much lower
118 than the symptom-threshold approach.

119 Poor sensitivity and specificity are problematic in diagnostics but higher
120 error rates than gold-standard methods may be tolerable depending on their
121 scale and impact given the local situation. Low specificity means a large number
122 of false positive classifications, where the patient is told they have COVID-19
123 but they actually do not. This might lead to patients unnecessarily self-isolating
124 and receiving support which can be expensive to the individuals and local public
125 health bodies, as well as reducing available resources for those who need them
126 [16]. Similarly, low sensitivity means more false negative classifications, where
127 the patient is told they do not have COVID-19 but they actually do, which can
128 lead to a health-risk for the individual and to the disease spreading further [17].
129 Although the default approach is generally to minimise both misclassification
130 rates (our “Agnostic” scenario below), the true costs of these misclassifications
131 will depend on local context. When the prevalence of the disease is low, false
132 positives may create local scepticism about the value of testing, or when there
133 are strong population-level mitigations already in place (such as a nationwide
134 lockdown), then false positives might be more costly than false negatives [16],
135 corresponding to our “Low-Level Cases” scenario. If the disease is abundant
136 or increasing rapidly then false negatives are likely to be more costly, as in our
137 “Rising Cases” scenario. Often the situation will be even more nuanced and a
138 different balance will need to be struck [4].

139 The “best” diagnostic, therefore, is not a single universal test. The two
140 dominant testing methods available in LMICs when not optimised for the
141 local situation are highly flawed. Relying solely on symptomatic diagnosis
142 will likely overestimate the number of individuals with COVID-19 due to its
143 lack of specificity. Conversely, RATs will give a false impression of control
144 due to the number of positive cases that will be missed. In this paper, we
145 demonstrate that by combining these two testing methods we can utilise their
146 complementary strengths, ameliorate their respective weaknesses, and optimise
147 them for different epidemiological scenarios. We aim to compare the performance
148 of these two testing methods and the combined approach both in terms of general
149 prediction and as diagnostics under three epidemiological scenarios with different
150 misclassification requirements. We show that the optimised combined data
151 models achieve equal-to-much-lower error rates than the next best method in all
152 metrics. We then discuss the role of statistically integrating data from multiple
153 imperfect testing methods in resource limited settings to improve the diagnosis
154 of diseases, particularly COVID-19.

155 **3. Methods (~1019 Words)**

156 Participants included in this study were identified for COVID-19 testing after
157 self-reporting symptoms to the Bangladesh government’s national hotlines for
158 COVID-19 support. Recruitment took place across Dhaka (the capital city of
159 Bangladesh) between 2nd April 2021 and 5th May 2021.

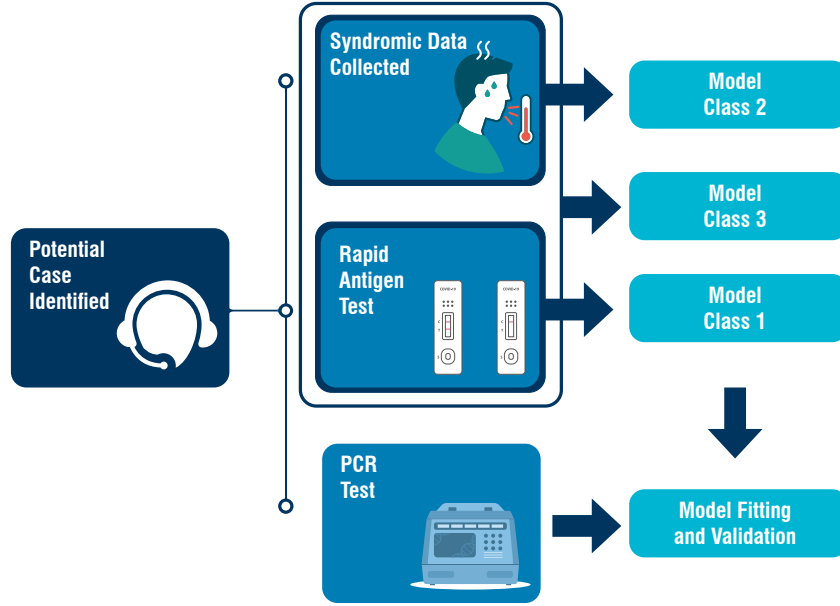


Figure 1: Schematic description of identification of likely COVID-19 patients by community support teams (CSTs), swab collection and model definitions. The teams collected syndromic data (age, gender and presence/absence of 14 predetermined symptoms), and two sets of naso-pharyngeal swabs (one each for Rapid Antigen Testing and RT-PCR). We then used rapid antigen testing (RAT) and syndromic data, two imperfect but inexpensive diagnostics, to generate three model classes: RAT result only in Model Class 1, syndromic data only in Model Class 2, and both RAT result and syndromic data in Model Class 3. The RT-PCR test result is used to train and test each model using temporal cross-validation.

160 Patients were selected for further testing conditional on the presence of a
 161 fever ($>38^{\circ}\text{C}$) at the point of testing and one or more of 14 additional symptoms
 162 associated with COVID-19 (breathing problems, coughing, diarrhoea, fever
 163 (ongoing), a headache, loss of taste, loss of smell, muscle pain, red eyes, a runny
 164 nose, a sore throat, tiredness, vomiting or a wet cough). The patient's age
 165 and gender were also recorded, but these data were not included in the patient
 166 selection criteria.

167 Nasal swabs and syndromic data were collected from the patient by medical
 168 technologists. One swab each was used for rapid antigen testing (RAT) and
 169 RT-PCR (gold-standard for COVID-19 status). The full questionnaire and
 170 testing protocols are provided in Appendix XX. Participants provided written
 171 informed consent to sample collection and for their test results to be analyzed in
 172 the study.

173 We examined the ability of the two imperfect identification methods, the
 174 syndromic profile and RAT result, to predict the patient's COVID-19 status

175 when used separately and together. The different data combinations define three
176 model classes (Figure 1).

177 Model Class 1 uses only the RAT result and is the simplest of the three.
178 It equates a positive RAT result with the patient being PCR positive, and a
179 negative RAT result with PCR negativity. Model Class 2 uses only the syndromic
180 data and Model Class 3 combines the RAT result with the syndromic data.

181 For Model Class 2, we used a Bayesian multivariate probit model [18]. The
182 multivariate probit structure allows the model to account for the correlations
183 between, and binary nature of, the symptoms (e.g. loss of taste is often correlated
184 with loss of smell). By using a Bayesian formulation, we are able to quantify
185 and propagate the uncertainty in the parameter estimates. Structurally, the
186 multivariate probit model allows the symptoms and COVID-19 status to be
187 treated as correlated binary outcomes with an intrinsic rate (the intercept for each
188 variable) and the patient’s age and gender, while propagating and quantifying
189 uncertainty.

190 In Model Class 3, we utilise the specificity of RAT by treating RAT positive
191 patients as PCR positive patients. The RAT negative patients are then modelled
192 using the sensitive syndromic approach using Model Class 2 to capture additional
193 PCR-positive patients that are missed by the RAT. This approach leverages
194 the fact that RAT-negative-PCR-positive patients may have different syndromic
195 profiles than RAT-positive-PCR-positive patients and allows the model to adapt
196 more specifically to that group. The models were fitted to the data using Bayesian
197 inference techniques based on Hamiltonian Monte Carlo in the Stan programming
198 language [19].

199 We conducted backwards model selection (starting with the most complex
200 model feasible, with all 14 symptoms and both covariates) to identify a subset of
201 models with the highest predictive power under temporal cross-validation (Figure
202 2). Reducing the number of possible models to a small number of the most
203 predictive models was necessary to reduce computational demand and reduce the
204 risk of overfitting models to the test scenarios. The large number of symptoms
205 means that there is a high number of potential model configurations (>131 000
206 for 14 symptoms and two covariates) which might, by chance, perform well on
207 the test sets (even under the challenging conditions of temporal cross-validation)
208 but lack transferability. By first using general predictive power to narrow down
209 the number of candidate models and then testing those models under more
210 specific scenarios, we are more likely to choose models which generalise well
211 to new data. The number of candidate models used was not pre-determined.
212 In fitting the models it became clear that there were “jumps” in performance
213 (as defined below) between models containing five and four symptoms, so the
214 models with zero to four symptoms were used as the candidate models.

215 We scored the models’ predictive power using cross-entropy. Cross-entropy
216 measures the accuracy of models that generate probabilities of binary outcomes,
217 rather than make binary classifications, similar in concept to a mean square error
218 for normally-distributed data, but adapted for binary data [20]. A cross-entropy
219 value close to zero corresponds to high levels of accuracy, with larger values
220 indicating lower accuracy. As the score only uses the predicted probability and

Find more
ways to cite
Figure
refig:data-
flowchart in
this text.

221 true values, it is possible to directly compare the predictions of any model for
222 the same test set. More details on the model structure and selection process,
223 including code, are available in Appendix XX.

224 We then compared models as classifiers using their false positive and false
225 negative rates in three epidemiological scenarios. In applied settings, models
226 must often be evaluated on their performance as classifiers rather than just as
227 prediction engines (i.e. their ability to say a patient is COVID-19 positive or
228 negative, not simply the probability the patient might be COVID-19 positive or
229 negative). To generate a classification, a probability threshold must be chosen
230 over which patients are classified as COVID-19 positive.

231 Classifier performance was compared using receiver operating characteristic
232 (ROC) curves and error rates under three epidemiological scenarios. ROC
233 curves show the true and false positive rates that each model can achieve. To
234 extract the error rate under the epidemiological scenarios (described below
235 and in Table 1), we use the ROC calculations, to identify the probability
236 threshold which most closely meets the scenario requirement (see Table 1 for
237 requirements and Appendix XX for calculation details). Comparing specific
238 scenarios allows classifier performance to be demonstrated in relevant scenarios.
239 Whether measuring classifier performance in specific scenarios or more generally,
240 decisions need to be made about the relative cost and acceptable levels of the two
241 types of misclassification (false positives and negatives). We strongly emphasise
242 that local context should be the guide in applying these methods.

243 In Scenario 1, we do not consider epidemiological context but simply costing
244 false negative and false positive rates equally. We do this by maximising the two
245 correct classification rates both individually and in total, as measured by the
246 harmonic mean (as opposed to the arithmetic mean which would only maximise
247 the rates in total). Scenario 2 corresponds to the current situation in Bangladesh
248 at time of writing (July 2021), with COVID-19 cases beginning to rapidly increase
249 again. Under these circumstances, false negatives are extremely costly relative
250 to false positives due to the exponential growth of the disease. In Scenario 3,
251 the pandemic is not declining but maintaining a steady rate of cases. In this
252 situation, policy-makers may be keen to keep false positive diagnoses low to
253 prevent lockdown fatigue and to keep the workforce active.

Add ROC
calculations
to appendix

254 4. Results (~353 words)

255 A total of 1172 subjects had data available for the current analyses. The
256 mean age of women participants (47% of the sample) was 37 (SD = 14), and for
257 men (53% of the sample) was 36 (SD = 14). Participants were identified by the
258 community support teams (CSTs) and drawn from across Dhaka.

259 Model selection for Model Class 2 (syndromic data only) and 3 (syndromic
260 and RAT data), each retained age as an explanatory variable and showed a
261 marked decline in predictive power at more than 4 symptoms. The final four
262 symptoms in order of importance (i.e. the most important symptom was retained
263 in all of the final 4 models, the least important symptom was only retained in
264 the 4 symptom model) were wet cough, runny nose, loss of smell and breathing

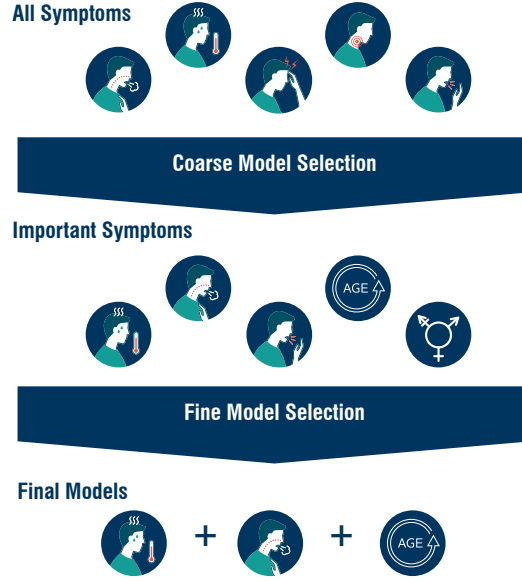


Figure 2: Schematic for rounds of model selection in the multivariate probit component of Model Classes 2 and 3. With 14 symptoms (only 5 shown here for demonstration purposes) and two covariates there are over 131000 possible model combinations. To make exploring these possible models computationally feasible and to reduce the risk of overfitting, we carried out two rounds of model selection. First, the data are divided into temporal cross-validation sets. The multivariate probit connects symptoms to the RT-PCR result through a correlation matrix. In the coarse model selection, the most complex feasible model (all symptoms and covariates) is fit to the training data. The estimated correlations between each symptom and the RT-PCR result are compared for each cross-validation set. The symptoms that have non-zero correlations in a systematic direction (i.e. all positively or all negatively correlated with RT-PCR result) are retained. The process is then repeated on each retained set of symptoms until the four symptoms in each model class with the strongest correlation to RT-PCR result. We then conduct a more exhaustive model selection on all the possible permutations of the four symptoms and two covariates. In this round, each model is fit to training data and used to predict for the test set, and the quality of those predictions is measured using cross-entropy scoring. The cross-entropy score is then used to select the best predictive model for each level of model complexity. Only these final models are then used for classification. This reduces the set of models tested as classifiers from $>131\,000$ to just four per model class.

Table 1: For each epidemiological scenario there is a requirement and a performance criterion. The requirement refers to a base level of performance the model must achieve; in general this will be a maximum acceptable error rate of some kind. These requirements were determined in discussion with members of the Institute of Epidemiology, Disease Control and Research, Ministry of Health, Bangladesh (IEDCR). The requirement determines a probability threshold for each model which most closely meets that requirement. The performance criterion is then used to determine which model performs the 'best' given that the requirement has been met.

Scenario Name	Requirement	Performance Criterion (Error)
1 Agnostic	Maximise correct classification rates	Sum of error rates
2 Rising Cases	Max. 20% false negative rate	False negative rate
3 Low-Level Cases	Max. 20% false positive rate	False positive rate

problems for Model Class 2, and fever, wet cough, tiredness and diarrhoea for Model Class 3. For both Model Class 2 and Model Class 3 model selection retained age but not gender as a covariate.

In the comparison of model predictive performance, Model Class 1 (RAT only) performed worst with an out-of-sample cross-entropy of 3.24 (cross-entropy values further from zero correspond to worse predictive performance). The median cross-entropy values were between 2.53 and 2.59 for models in Class 2 (syndromic data only). Models in Class 3 (combined data model) performed best with cross-entropy values between 1.44 and 1.47 (see Figure 3).

General model classification performance is shown by the full ROC curves for each model (Figure 4).

Scenario specific classification performance is shown in Figure 5. In Scenario 1 ("Agnostic," see Table 1), the median error was 0.47 for models in Class 1 and Class 3 and between 0.87 and 0.9 for models in Class 2 (Figure 5). In Scenario 2 ("Rising Cases"), Model Class 1 was unable to meet the required false negative rate. The median errors were between 0.74 and 0.76 for models in Class 2, and 0.43 and 0.51 for models in Class 3 (Figure 5). In Scenario 3 ("Low-Level Cases"), the error in Class 1 was 0.02 and the median errors ranged from 0.19 to 0.2 for Class 2, and 0.15 to 0.2 for Class 3 (Figure 5).

DHC: out-of sample vs test set cross-entropy

5. Discussion (~1314 Words)

We have demonstrated that combining rapid antigen tests (RATs) with syndromic modelling yields better prediction of COVID-19 status and greater flexibility than each diagnostic individually. These improvements are non-trivial in real-world settings. In Bangladesh, there are currently 15 000 new cases being identified every day, using only the limited supply of RT-PCR, the pandemic growth is accelerating and every missed case has a compounding effect. Scenario 2 ("Rising Cases") was developed with the need to keep false negative rates

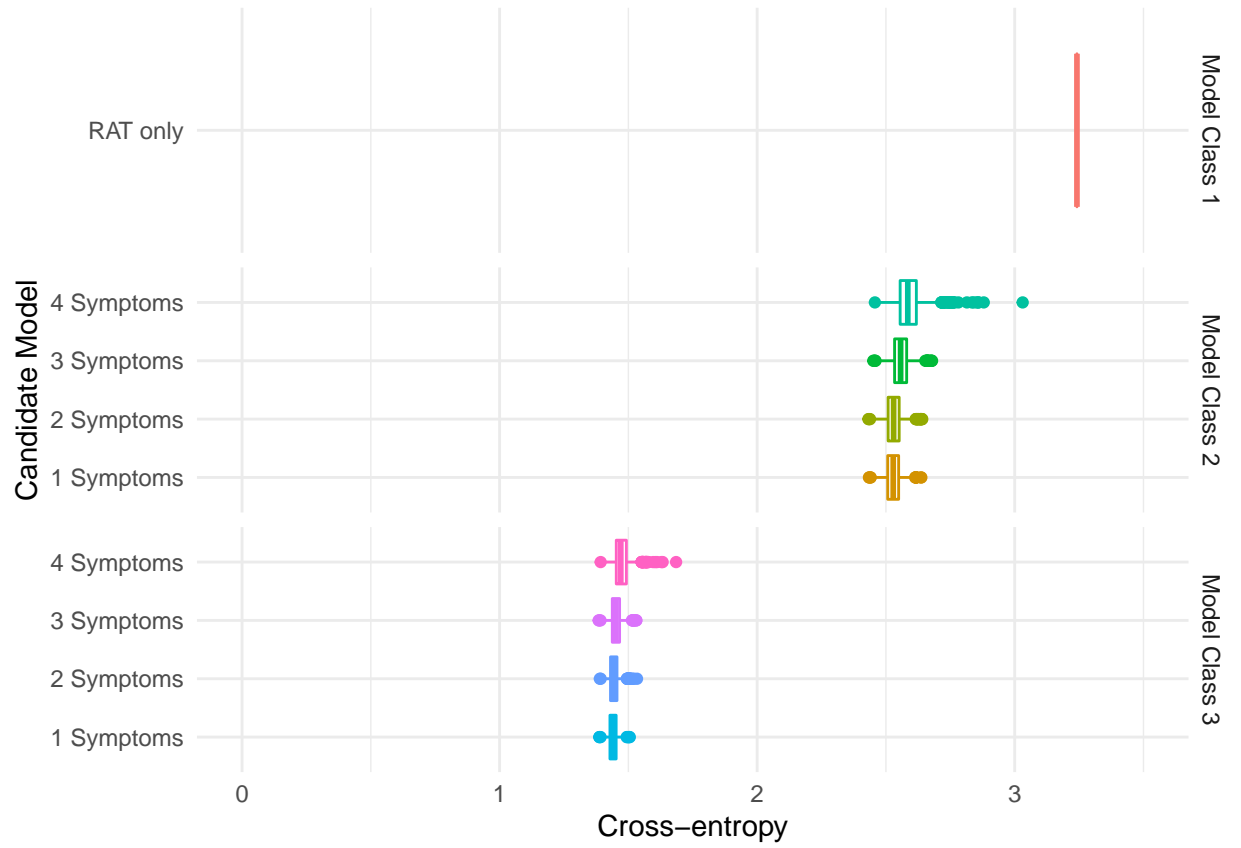


Figure 3: Predictive performance of candidate models. Interquartile ranges for the posterior cross-entropy of the best candidate models at each level of model complexity tested under temporal cross-validation. cross-entropy is a measure of distance from the truth, so values closer to zero indicate better models. The intermediate complexity models perform best at prediction, although performance is similar across all the models within each model class (1: rapid antigen testing (RAT) only; 2: syndromic data only; and 3: combined RAT and syndromic data).

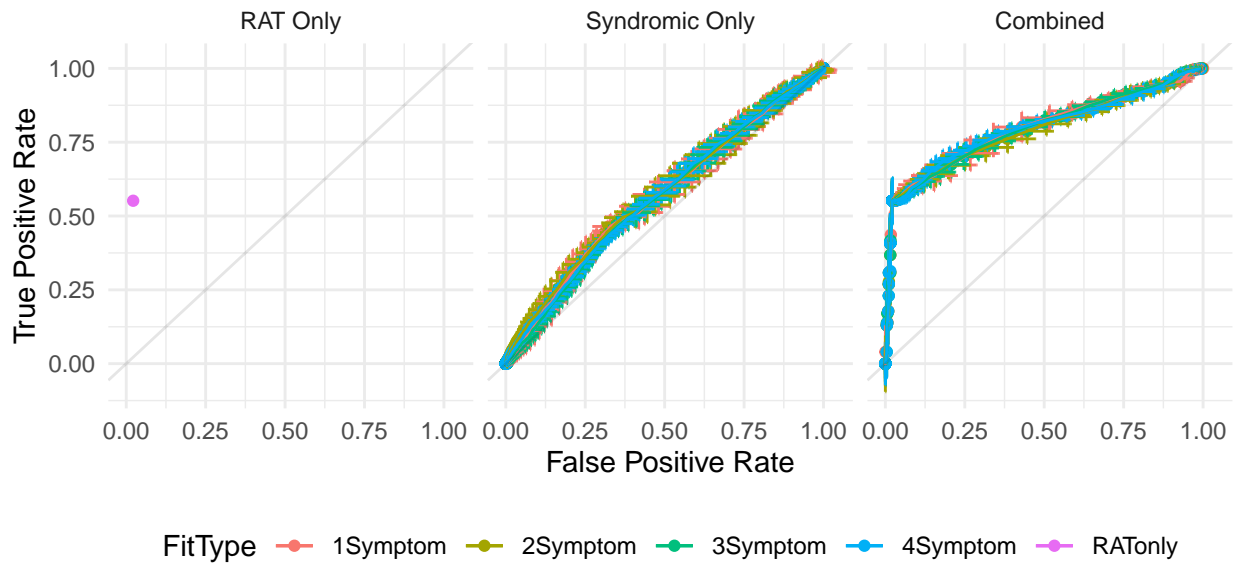


Figure 4: Receiver operating characteristics for rapid antigen testing (RAT) only approach (Model Class 1) and posterior mean (\pm posterior standard deviation) receiver operating characteristics for Class 2 (syndromic data only) and 3 (syndromic and RAT data) models. These curves demonstrate the performance of the model for any hypothetical scenario as defined by the axes (as opposed to Figure 5 which demonstrates model performance in specific epidemiological scenarios which are realisations of a single point in this space).

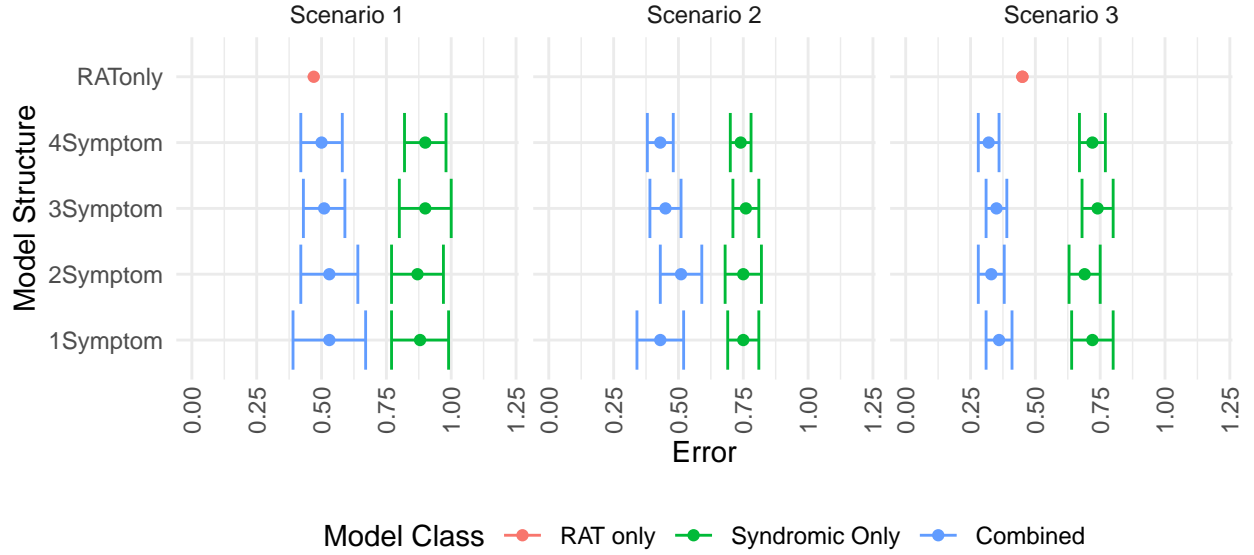


Figure 5: Performance of models under each scenario measured by errors defined in Table 2. Low errors correspond to better model performance. There is no error rate defined for the Model Class 1 (RAT only model) in Scenario 2 as the model failed to meet the requirement for that scenario (making the error functionally infinite).

low and maps well onto the situation in Bangladesh (see Table 1). In this scenario, the combined data model (Model Class 3) false negative rate is 26 percentage points lower than that of the RAT only model (Model Class 1). Although the syndromic only model (Model Class 2) matches the combined models false negative rate, its false positive rate is 31 percentage points higher. These are large performance gains for any diagnostic but when deployed at the scale of Bangladesh and similar countries, these improvements represent catching tens of thousands of cases that would otherwise be missed. Furthermore, this boost is achieved with data that are already being collected in Bangladesh and other low- and middle- income countries (LMICs). Outwith developing and rerunning the models presented in this paper, these improvements are essentially cost-free and eminently scalable.

The pattern is similar in epidemiological Scenarios 1 (“Agnostic”) and 3 (“Low-Level Cases”), with the combined model class performing performing equally well or better than the other two classes (Figure 5). These three scenarios only offer snapshots of performance, however, and we strongly advocate defining model performance in terms of false negative and false positive rates with reference to local conditions. An indication of how these models will perform under any condition can be obtained by comparing the more generic model performance metrics for prediction and classification (Figures 3 and 4, respectively). These figures demonstrate both the added flexibility of the more

313 complex model classes that allow them to be tailored to specific needs and
314 the need to combine the high-quality but inflexible RAT results with the more
315 flexible but lower quality syndromic data. Interestingly, the most of the Class
316 2 models performs approximately as well as chance except the simplest which
317 performs worse than chance. A model that performs worse than random can still
318 be useful if one takes the inverse decision. Even a flexible model which performs
319 as well as random classification can be useful if those error rates reflect those
320 needed in a given local situation. Fortunately, Model Class 3 is both flexible and
321 performs better than random.

322 We have deliberately not emphasised the final symptoms chosen through
323 model selection in this paper as we are focusing on prediction and classification for
324 a unique sub-population: self-referring, symptomatic patients. We do, however,
325 highlight that while fever and loss of smell were the two most important symptoms
326 in the two classes of syndromic models, the other symptoms retained were different
327 (with cough and wet cough retained in the combined syndromic and RAT model,
328 Class 3, and loss of taste and vomiting in the syndromic only model, Class 2).
329 Further research is needed to understand the mechanisms by which symptoms
330 predict COVID-19 and by which RAT misses COVID-19. Of particular interest
331 is whether individuals that are missed by RAT are less infectious, which could
332 be explored by using Threshold Cycle (Ct) values from the RT-PCR to compare
333 viral load with respect to prediction by the different methods [21]. We note
334 also that, as expected, age was retained in model selection. We were, however,
335 surprised that gender was removed during model selection. Gender is thought to
336 play a major role in infection risk [23]. As we are looking to predict symptomatic
337 COVID-19 in symptomatic individuals, generalised risk of infection is perhaps
338 less predictive than expected, potentially due to the balancing of risk and burden
339 [24].

340 Using a large sample collected under field-realistic conditions, we have rigor-
341 ously tested our approach. By taking a statistical modelling approach to case
342 identification, we are able to update our diagnostic process in real time, allowing
343 this method to readily adapt to new variants (or even new diseases) or new
344 priorities for resource allocation. The modelling frameworks we have used are
345 also sufficiently flexible to accommodate new data sources. Of particular interest
346 are extensions to include the “pandemic context” in the model using space-time
347 data. Furthermore, by using more sophisticated modelling structures that work
348 at the scale of probabilities, rather than binary tests, it is possible to tune error
349 rates to better reflect the local relative costs of false positives and false negatives.
350 Naturally, these strengths have complementary limitations. Our models require
351 updating in real-time and can only achieve good performance if the validation
352 data are of high quality. Similarly, targeting error rates is only sensible if those
353 rates properly reflect local conditions which is hard to do in practice. These
354 limitations should be seriously considered but the alternatives for imperfect
355 testing methods are diagnostics that cannot be tailored to local conditions at all
356 (and, as such may perform worse than a method which is sub-optimally tailored
357 to local conditions) or diagnostics which make these decisions implicitly and not
358 explicitly. We choose to make these decisions explicitly to allow them to be more

@Dirk - can
you please
clarify this
or suggest
reasons the
model is per-
forming so
badly. The
red line here
is a univari-
ate probit
regression
with one
continuous
covariate
so I don't
understand
why it's per-
forming so
poorly un-
less the tem-
poral cross-
validation
sets are
wildly dif-
ferent from
each other?

readily challenged, researched and improved upon. We also emphasise the need for rigorous experimental design to ensure findings from the sample population are applicable to the target population and the need for further research into understanding error rate trade-offs in applied settings.

We believe that the combined syndromic and rapid antigen testing approach represents the most promising approach to large-scale testing in LMICs for COVID-19 at present. By using the small amount of RT-PCR testing possible and formally integrating multiple imperfect, non-gold-standard methods, we can tune these diagnostics to our local conditions. Where data collection is being coordinated electronically (e.g. through mobile applications), the models can be used for diagnosis in the field and updated in real-time. We have demonstrated that these improvements can be impressive in real-world scenarios, and will have a large impact when scaled to the population sizes in LMICs. The methodology we have outlined here is applicable to a wide range of diseases and settings across LMICs. One of the biggest challenges in diagnosing and tracking many diseases in resource-limited settings is the low availability of access to gold-standard testing (such as RT-PCR in the case of COVID-19) and high error rates of alternative testing methods. In this paper, we have outlined the process for coupling a small number of gold-standard tests with formal statistical integration of alternative testing methods, to generate high quality diagnostic models. This process readily maps onto many other case identification problems, including the diagnosis of several neglected tropical diseases. For example, malaria (gold standard (GS) is also RT-PCR, imperfect methods (IM) include antigen tests, syndromic diagnosis and blood smears), schistosomiasis (GS: RT-PCR or autopsy; IM: Kato Katz egg counts, antibody detection) and rabies (GS: fluorescent antibody test; IM: light microscopy, differential diagnosis).

The management of global pandemics can only be done with global testing. While the quest to achieve this using only gold-standard diagnostic methods is laudable, it is also often impractical. Imperfect diagnostics are frequently imperfect in different ways, and these differences are ripe for statistical treatment. What is more, these approaches are often more agile than gold-standard diagnostics in situations of flux, for example, in the early stages of new pandemics or disease strains, when fast responses are essential.

By investing in understanding how to utilise the complementary strengths of imperfect testing and deploy the limited gold-standard testing available for validation, we can provide good quality testing at the scale needed to fight infectious diseases.

6. Funding (~26 words)

The Bill and Melinda Gates Foundation funded work by FAO (INV-022851), and University of Glasgow reports funding from Wellcome (207569/Z/17/Z). The authors declare no competing interests.

400 7. Acknowledgements (~67 words)

401 We would like to thank members of the community support teams in
 402 Bangladesh who have provided essential services throughout the pandemic.
 403 Earlier drafts of this manuscript benefited from the input of Daniel Haydon,
 404 Anne-Sophie Bonnet-Lebrun, Luca Nelli, Crinan Jarrett, Rita Claudia Cardoso
 405 Ribeiro, Halfan Ngowo, Heather McDevitt and Gina Bertolacci. The University
 406 of Glasgow COVID-19 in LMICs Group provided the environment in which to
 407 develop this work.

408 References (Max 30)

- 409 [1] Corman VM, Landt O, Kaiser M, Molenkamp R, Meijer A, Chu DK, et al.
 410 Detection of 2019 novel coronavirus (2019-nCoV) by real-time RT-PCR.
 Eurosurveillance 2020;25:2000045.
- 411 [2] Tahamtan A, Ardebili A. Real-time RT-PCR in COVID-19 detection:
 412 Issues affecting the results. Expert Review of Molecular Diagnostics
 2020;20:453–4.
- 413 [3] Chowdhury R, Luhar S, Khan N, Choudhury SR, Matin I, Franco OH.
 414 Long-term strategies to control COVID-19 in low and middle-income
 countries: An options overview of community-based, non-pharmacological
 interventions. European Journal of Epidemiology 2020;35:743–8.
- 415 [4] Vandenberg O, Martiny D, Rochas O, Belkum A van, Kozlakidis Z. Con-
 siderations for diagnostic COVID-19 tests. Nature Reviews Microbiology
 2021;19:171–83.
- 416 [5] Cash R, Patel V. Has COVID-19 subverted global health? The Lancet
 417 2020;395:1687–8.
- 418 [6] Olalekan A, Iwalokun B, Akinloye OM, Popoola O, Samuel TA, Akinloye
 419 O. COVID-19 rapid diagnostic test could contain transmission in low-
 and middle-income countries. African Journal of Laboratory Medicine
 2020;9:1–8.
- 420 [7] Linares M, Pérez-Tanoira R, Carrero A, Romanyk J, Pérez-García F,
 421 Gómez-Herruz P, et al. Panbio antigen rapid test is reliable to diagnose
 SARS-CoV-2 infection in the first 7 days after the onset of symptoms.
 422 Journal of Clinical Virology 2020;133:104659.
- 423 [8] Boum Y, Fai KN, Nikolay B, Mboringong AB, Bebell LM, Ndifon M,
 et al. Performance and operational feasibility of antigen and antibody
 rapid diagnostic tests for COVID-19 in symptomatic and asymptomatic
 patients in cameroon: A clinical, prospective, diagnostic accuracy study.
 424 The Lancet Infectious Diseases 2021.
- 425 [9] Mak GC, Cheng PK, Lau SS, Wong KK, Lau C, Lam ET, et al. Evaluation
 of rapid antigen test for detection of SARS-CoV-2 virus. Journal of
 426 Clinical Virology 2020;129:104500.

427 [10] Jin Y-H, Cai L, Cheng Z-S, Cheng H, Deng T, Fan Y-P, et al. A rapid
advice guideline for the diagnosis and treatment of 2019 novel coronavirus
(2019-nCoV) infected pneumonia (standard version). *Military Medical*
428 *Research* 2020;7:1–23.

429 [11] Sim JXY, Conceicao EP, Wee LE, Aung MK, Seow SYW, Teo RCY, et
al. Utilizing the electronic health records to create a syndromic staff
surveillance system during the COVID-19 outbreak. *American Journal of*
430 *Infection Control* 2021;49:685–9.

431 [12] Undurraga EA, Chowell G, Mizumoto K. COVID-19 case fatality risk
by age and gender in a high testing setting in latin america: Chile,
432 march–august 2020. *Infectious Diseases of Poverty* 2021;10:1–1.

433 [13] Wenham C, Smith J, Morgan R. COVID-19: The gendered impacts of
the outbreak. *The Lancet* 2020;395:846–8.

434 [14] Mayorga L, Samartino CG, Flores G, Masuelli S, Sánchez MV, Mayorga
435 LS, et al. A modelling study highlights the power of detecting and
isolating asymptomatic or very mildly affected individuals for COVID-19
epidemic management. *BMC Public Health* 2020;20:1–1.

436 [15] Garry S, Abdelmagid N, Baxter L, Roberts N, Waroux O le P de, Ismail
437 S, et al. Considerations for planning COVID-19 treatment services in
humanitarian responses. *Conflict and Health* 2020;14:1–1.

438 [16] Surkova E, Nikolayevskyy V, Drobniewski F. False-positive COVID-19
439 results: Hidden problems and costs. *The Lancet Respiratory Medicine*
2020;8:1167–8.

440 [17] West CP, Montori VM, Sampathkumar P. COVID-19 testing: The threat
441 of false-negative results. *Mayo clinic proceedings*, vol. 95, Elsevier; 2020,
p. 1127–9.

442 [18] Albert JH, Chib S. Bayesian analysis of binary and polychotomous re-
443 sponse data. *Journal of the American Statistical Association* 1993;88:669–
444 79.

445 [19] Carpenter B, Gelman A, Hoffman MD, Lee D, Goodrich B, Betancourt M,
et al. Stan: A probabilistic programming language. *Journal of Statistical*
446 *Software* 2017;76:1–32.

447 [20] Gneiting T, Raftery AE. Strictly proper scoring rules, prediction, and
448 estimation. *Journal of the American Statistical Association* 2007;102:359–
78.

449 [21] Albert E, Torres I, Bueno F, Huntley D, Molla E, Fernández-Fuentes MÁ,
et al. Field evaluation of a rapid antigen test (panbio™ COVID-19 ag
rapid test device) for COVID-19 diagnosis in primary healthcare centres.
450 *Clinical Microbiology and Infection* 2021;27:472–e7.

451 [22] Menni C, Valdes AM, Freidin MB, Sudre CH, Nguyen LH, Drew DA,
et al. Real-time tracking of self-reported symptoms to predict potential
452 COVID-19. *Nature Medicine* 2020;26:1037–40.

- 453 [23] Wolff D, Nee S, Hickey NS, Marschollek M. Risk factors for covid-19
454 severity and fatality: A structured literature review. *Infection* 2020;1–4.
- 455 [24] Joe W, Kumar A, Rajpal S, Mishra U, Subramanian S. Equal risk, unequal
burden? Gender differentials in COVID-19 mortality in india. *Journal of*
456 *Global Health Science* 2020;2.