

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, United Nations Development Program, ICT Ministry, 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 (288 words)

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

Community support teams in Dhaka, Bangladesh identified potential COVID-19 patients in Dhaka using syndromic surveillance. A sample (n = 1172) of

*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.kundegorski@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)

these patients was tested using RAT and syndromic data were collected. Models were fit to predict RT-PCR status using the RAT data, the syndromic data, and the two combined. Model performance was measured using predictive power and classification performance under three epidemiological scenarios: “Agnostic,” “Rising Cases” and “Low-Level Cases.”

Findings

Combined data models yielded equal or improved performance over syndromic- and RAT-only models across all three epidemiological scenarios and when compared as more generic prediction and classification engines. In the “Rising Cases” scenario, which most closely represents the current situation in many LMICs, the combined data model false negative rate is 26 (IQR: 24, 29) percentage points lower than that of the RAT only model. Although the syndromic only model matches the combined models false negative rate, its false positive rate is 31 (IQR: 29, 34) percentage points higher.

Interpretation

Small, scalable improvements in the accuracy of mass-deployed but imperfect diagnostic methods can then make a very big difference for pandemic control. We demonstrate that such improvements can be achieved by statistically utilising complementary strengths and weaknesses across two imperfect diagnostics, we can greatly improve the detection of COVID-19.

Funding

The Bill and Melinda Gates Foundation and the Wellcome Trust.

2. Introduction (1080 Words)

Identification and isolation of COVID-19 cases remains key to the pandemic response across the globe. The faster and more accurately we can identify cases, the more effectively we can provide clinical care, reduce transmission of infection and develop population-level interventions. RT-PCR testing has rapidly become the default, gold-standard test for COVID-19 in applied settings (although see [1]) due to its high sensitivity and specificity for COVID-19 [3]. Most of the world’s population, however, live in low- and middle-income countries (LMICs) where the laboratory facilities needed to carry out RT-PCR tests are often scarce and hard to reach [5], and patient diagnosis and support comes from telemedicine or community support teams (CSTs) composed of local volunteers with basic training. COVID-19 diagnosis worldwide, therefore, must be made accessible using inexpensive methods that can be carried out locally [7].

An increasingly popular alternative to RT-PCR is rapid antigen testing (RAT) [8]. Like RT-PCR, these tests have high specificity for COVID-19 while being less expensive, easier to implement, and faster but with lower sensitivity [9]. For RT-PCR testing, patients must travel to a designated site or have officials visit their home in enhanced personal protective equipment. In contrast, RATs can be conducted on nasal swabs, completed in the home with minimal PPE, and results are available in 30 minutes. RATs can be taken by persons with limited training, thus decreasing the time and expense associated with identifying cases.

71 Together, these traits make RATs an appealing alternative to RT-PCR, however,
72 concerns have been raised that the lower sensitivity of RAT [10] leads to more
73 false negative diagnoses.

74 Another diagnostic that has been used since the start of the pandemic
75 is symptom-thresholding [11]. Here, a patient presenting with a fever and
76 one or more symptoms is treated as a COVID-19 positive patient. The main
77 advantage of this approach is the ease of implementation. As with RAT, symptom-
78 thresholding is faster, cheaper and less invasive than RT-PCR. Unlike RAT,
79 symptom-thresholding can be scaled immediately at the onset of a pandemic,
80 however, it is also reliant on thresholds developed then. These thresholds were
81 necessarily drawn from clinical intuition, rather than data, often for different
82 variants and populations than they are now applied to. Consequently, the
83 relationship between the thresholds and the true COVID-19 status is often weak,
84 with low specificity leading to a very large number of false positive diagnoses. A
85 natural extension, therefore, is syndromic modelling. In this approach, rather
86 than using a set of pre-determined thresholds, a range of symptomatic and risk
87 factor data (such as age and gender) are collected and then a sub-sample of
88 patients is tested using RT-PCR for validation [12]. These data are used to fit a
89 model that allows more accurate prediction of how likely a patient is to have
90 COVID-19 through the identification of COVID-19 syndromes [14].

91 It is worth highlighting at this point that in resource-limited settings there is
92 very limited provision for testing of asymptomatic cases, despite their important
93 role in disease transmission [15] . Even while focusing solely on symptomatic
94 patients, syndromic modelling is a complex and nuanced task. Disease syndromes
95 can change between populations, when new variants emerge, and as other diseases
96 become more or less common [16]. These changes can make syndromic models
97 generalise poorly. For example, if another disease for which loss of smell is
98 a symptom becomes common, loss of smell is no longer strongly indicative of
99 COVID-19. Similarly, if everyone who presents has a cough, regardless of their
100 COVID-19 status, then coughing will show no relationship with COVID-19
101 (even if the two are strongly related in the general population). Furthermore,
102 symptoms do not always occur in isolation, some, like loss of smell and loss of
103 taste, are strongly related. Unfortunately, the majority of syndromic modelling
104 methods currently used do not account for these complexities. Even where they
105 can, at least partially, be accounted for, the many types of common respiratory
106 disease generally means that syndromic modelling still tends to have quite low
107 specificity [16].

108 Moderate to poor sensitivity and specificity are problematic in diagnostics but
109 may be tolerable depending on their scale and impact given the local situation.
110 Low specificity means a patient is likely to be told they have COVID-19 when
111 they do not (a high false positive rate), leading to patients unnecessarily self-
112 isolating and receiving support. This is expensive to the individual and to
113 local public health bodies, reducing available resources for those who need them
114 [17]. Similarly, low sensitivity means more patients being told they do not have
115 COVID-19 when they actually do (a high false negative rate), leading to the
116 individual not getting appropriate support or taking action to prevent the disease

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 spreading further [18]. Although the default approach is generally to minimise
 118 both misclassification rates (our “Agnostic” scenario below), the true costs of
 119 these misclassifications will depend on local context. When the prevalence of the
 120 disease is low, false negatives will be correspondingly low and false positives may
 121 create local scepticism leading to poor adherence longer term. In this situation
 122 (our “Low-Level Cases” scenario), false positives might be more costly than false
 123 negatives [17]. If the disease is abundant or increasing rapidly then changes in
 124 the false negative rate might have an outsized impact on the pandemic trajectory
 125 and thus be more costly, as in our “Rising Cases” scenario. Often the situation
 126 will be even more nuanced and a different balance will need to be struck [5].

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

143 **3. Methods (965 words)**

144 *3.1. Data Collection*

145 Participants included in this study were identified for COVID-19 testing by
 146 community support teams (CSTs). Recruitment took place across Dhaka (the
 147 capital city of Bangladesh) between 19th May 2021 and 11th July 2021.

148 Patients were selected for further testing if they had a fever ($>38^{\circ}\text{C}$) at the
 149 point of testing and one or more of 14 symptoms associated with COVID-19
 150 (breathing problems, coughing, diarrhoea, fever (ongoing), a headache, loss of
 151 taste, loss of smell, muscle pain, red eyes, a runny nose, a sore throat, tiredness,
 152 vomiting or a wet cough). If selected, the CSTs collected the patient’s age and
 153 gender, and took two nasal swabs.

154 One swab each was used for rapid antigen testing (RAT) and RT-PCR.
 155 The full questionnaire and testing protocols are provided in Supplementary 1.
 156 Participants provided written informed consent to sample collection and for their
 157 results to be analyzed in the study.

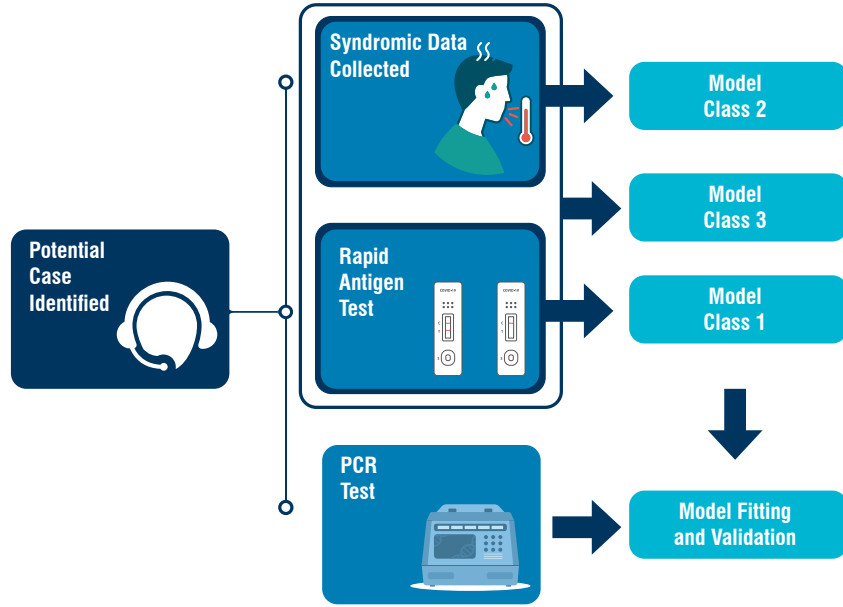


Figure 1: Schematic description of identification of likely COVID-19 patients by community support teams (CSTs), swab collection and model definitions. The teams collected syndromic data (age, gender and presence/absence of 14 predetermined symptoms), and two sets of naso-pharyngeal swabs (one each for 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.

158 *3.2. Modelling*

159 *3.2.1. Structure*

160 We examined the ability of the two imperfect identification methods, syn-
161 dromic modelling and RAT, to predict the patient’s COVID-19 status when used
162 separately and together. These combinations define three model classes (Figure
163 1).

164 Model Class 1 uses only the RAT result. It equates being RAT-positive with
165 the patient being PCR-positive for COVID-19 (hereafter, PCR-positive), and
166 being RAT-negative with PCR-negativity.

167 Model Class 2 uses only the syndromic data. For this model, we used a
168 Bayesian multivariate probit model [19]. The multivariate probit structure allows
169 the model to account for the binary and correlated nature of the symptoms
170 while conditioning on the risk factors of age and gender. By using a Bayesian
171 formulation, we are able to quantify uncertainty in the parameter estimates.

172 Model Class 3 combines the two data sources. We utilise the specificity of RAT
173 by treating RAT-positive patients as PCR-positive patients. The RAT-negative
174 patients are modelled using the sensitive syndromic approach using Model Class
175 2 to capture PCR-positive patients that are missed by the RAT. This approach
176 leverages the potential different syndromic profiles of PCR-positive patients who
177 are RAT-positive and -negative, allowing the model to adapt solely to the latter.
178 The models were fitted to the data using Bayesian inference techniques based on
179 Hamiltonian Monte Carlo in the Stan programming language [20].

180 *3.2.2. Model Selection*

181 We conducted backwards model selection (starting with the most complex,
182 biologically plausible model) to identify a subset of models with the highest
183 predictive power under temporal cross-validation (Figure 2). Reducing the
184 number of possible models was necessary to reduce computational demand and
185 reduce the risk of overfitting models to the test scenarios. The large number
186 of symptoms corresponds to a high number of potential model configurations
187 ($>131\,000$ for 14 symptoms and two covariates) which might perform well on the
188 test sets (even under the challenging conditions of temporal cross-validation) but
189 lack transferability. By first using general predictive power to narrow down the
190 number of candidate models and then testing those models, we are more likely
191 to choose models which generalise well to new data. The number of candidate
192 models used was not pre-determined but it was clear when fitting the models
193 that there were “jumps” in performance (as defined below) between models
194 containing five and four symptoms, so the models with one to four symptoms
195 were used as the candidate models. Zero symptom models were not included
196 in the analysis as they do not correspond to a feasible policy (with covariates
197 they would require governments to ask individuals of a given gender and age
198 as COVID-19 positive, and without covariates they would involve randomly
199 assigning individuals as COVID-19 positive).

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

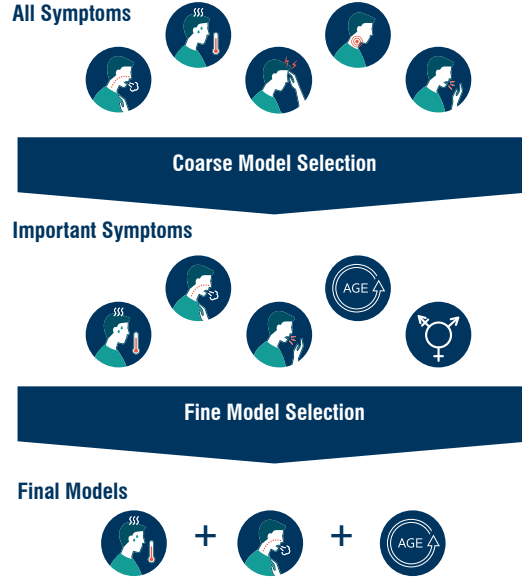


Figure 2: Schematic for rounds of model selection in the multivariate probit component of Model Classes 2 and 3. With 14 symptoms (only 5 shown here for demonstration purposes) and two covariates there are over 131000 possible model combinations. To make exploring these possible models computationally 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

3.2.3. Predictive Performance

We scored the models' predictive power using cross-entropy. Cross-entropy measures the accuracy of models that generate probabilities of binary outcomes, rather than make binary classifications, similar in concept to a mean square error for normally-distributed data, but adapted for binary data [21]. A cross-entropy value close to zero corresponds to high levels of accuracy, with larger values indicating lower accuracy. More details on the model structure and selection process, including code, are available in Supplementary 2.

3.2.4. Classification Performance

In applied settings, models must often be evaluated on their performance as classifiers rather than just as prediction engines (i.e. their ability to say a patient is COVID-19 positive or negative, not simply the probability the patient might be COVID-19 positive or negative). To generate a classification, a probability threshold must be chosen over which patients are classified as COVID-19 positive.

Classifier performance was compared both generically (using receiver operating characteristic (ROC) curves [22]) and under three epidemiological scenarios (using error terms described in Table 1). We strongly emphasise that generic performance here is only used to show the flexibility of the model classes; the best model for a local situation can only be determined if the relative cost of false positives and false negatives is known.

In Scenario 1, we do not consider epidemiological context but simply minimise false negative and false positive rates equally. We do this by maximising the two correct classification rates both individually and in total, as measured by the harmonic mean (as opposed to the arithmetic mean which would only maximise the rates in total). Scenario 2 corresponds to the current situation in Bangladesh at time of writing (July 2021), with COVID-19 cases beginning to rapidly increase again. Under these circumstances, false negatives are extremely costly relative to false positives due to the exponential growth of the disease. In Scenario 3, the pandemic is not declining but maintaining a steady rate of cases. In this

229 situation, policy-makers may be keen to keep false positive diagnoses low to
230 prevent lockdown fatigue and to keep the workforce active. The requirements in
231 Scenario 2 and 3 were developed in discussion with the Institute of Epidemiology,
232 Disease Control and Research (IEDCR), Bangladesh, for illustrative purposes.

233 4. Results (476 words)

234 Of 1241 subjects surveyed, a total of 1172 subjects had complete data
235 available for the current analyses with the remainder removed due to duplication
236 of barcodes or missing data. The mean age of women participants (47% of the
237 sample) was 37 (SD = 14) years, and for men (53% of the sample) was 36 (SD =
238 14) years. Participants were identified by the community support teams (CSTs)
239 and drawn from across Dhaka.

240 Model selection for both Model Class 2 (syndromic data only) and 3 (syn-
241 dromic and RAT data) showed a marked decline in predictive power at more than
242 4 symptoms. The covariate gender was dropped for both model classes while
243 age was dropped in Class 2 but retained in Class 3. The final four symptoms
244 in order of importance (i.e. the most important symptom was retained in all
245 of the final 4 models, the least important symptom was only retained in the 4
246 symptom model) were loss of taste, diarrhoea, vomit and fever for Model Class
247 2, and fever, wet cough, cough and loss of taste for Model Class 3.

248 In the comparison of model predictive performance, Model Class 1 (RAT
249 only) performed worst with an out-of-sample cross-entropy of 3.24 (cross-entropy
250 values further from zero correspond to worse predictive performance). The
251 median cross-entropy values were between 2.53 and 2.59 for models in Class 2.
252 Models in Class 3 performed best with cross-entropy values between 1.44 and
253 1.47 (see Figure 3).

254 Generic model classification performance is shown by their ROC curves
255 (Figure ??).

DHC out-
of sam-
ple vs test
set cross-
entropy

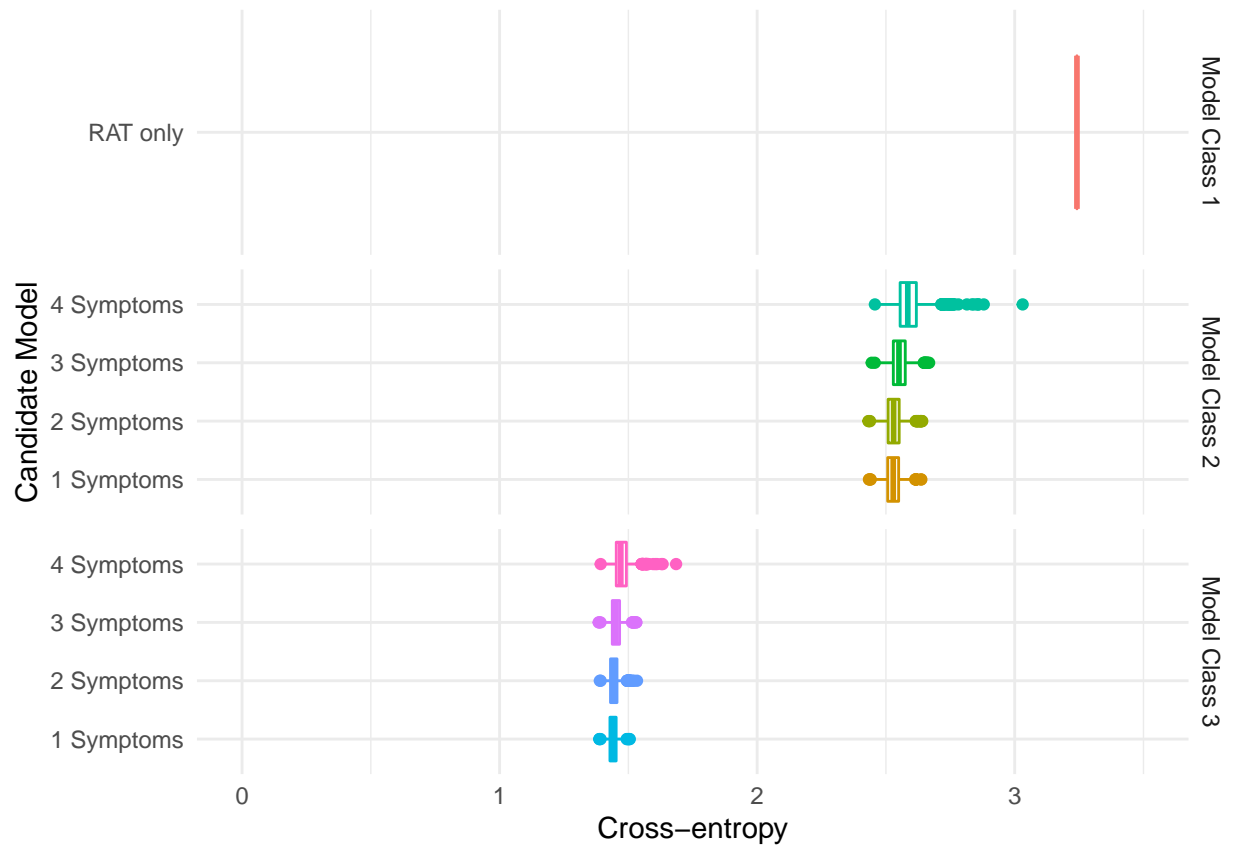
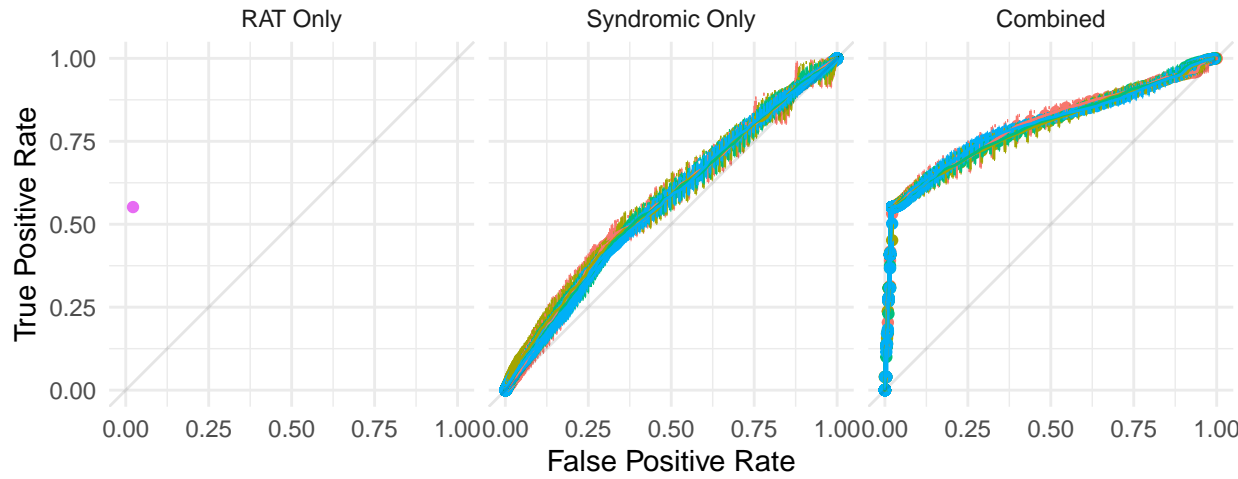
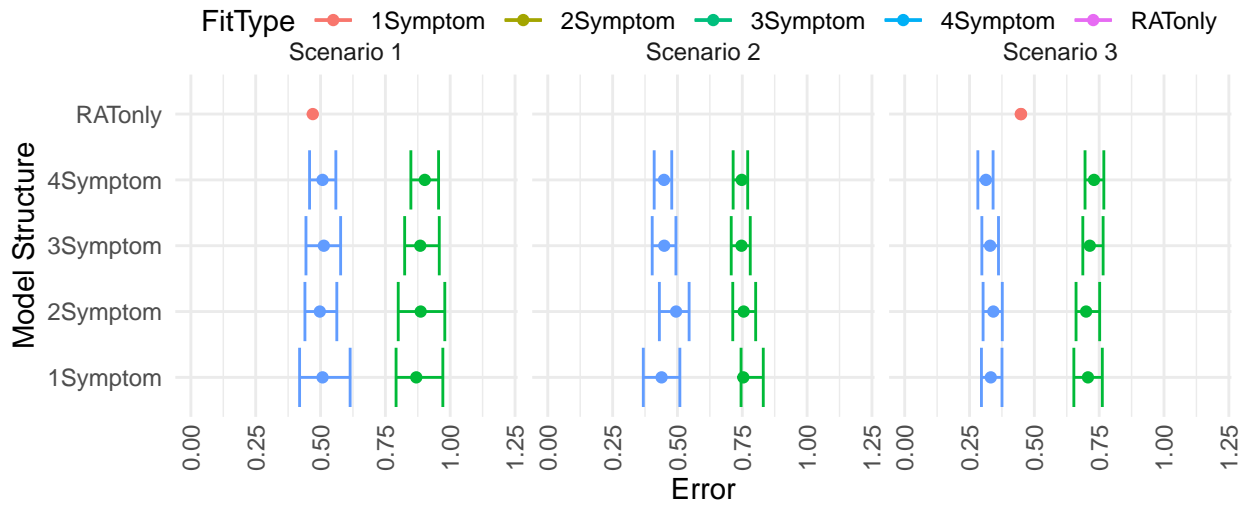


Figure 3: Predictive performance of candidate models. Interquartile ranges for the posterior cross-entropy of the best candidate models at each level of model complexity tested under temporal cross-validation. cross-entropy is a measure of distance from the truth, so values closer to zero indicate better 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).



256



257

258 Scenario specific classification performance is shown in Figure ???. Across
 259 all three scenarios (defined in Table 1), the best models in Class 3 performed
 260 equally well or better than the other two model classes. In Scenario 1 (“Agnos-
 261 tic”), models in Classes 1 and 3 performed equally well (overlapping posterior
 262 interquartile ranges) and distinctly better (no overlap in posterior interquartile
 263 range) than models in Class 2. In Scenario 2 (“Rising Cases”), Model Class 1
 264 failed to meet the requirement and so was excluded, and Model Class 3 once
 265 again outperformed Class 2. In Scenario 3 (“Low-Level Cases”), Model Class 2
 266 once again performed worst, and Model Class 3 achieved the lowest error, with
 267 Model Class 1 falling in between the two (closer to Class 3 than 2). For Classes
 268 2 and 3 across all the scenarios the number of symptoms made relatively little

269 difference within the final four candidate models in terms of median performance,
 270 although the more complex models have higher precision. It should be noted that
 271 the candidate models are chosen as a result of a selection process and performed
 272 much better than more complex models (i.e. those with 5 or more symptoms) or
 273 simpler models (with no symptoms but an intercept and covariates) in terms of
 274 cross-entropy and ROC, indicating they would likely also perform worse in these
 275 scenarios.

276 5. Discussion (815 words)

277 We have demonstrated that combining rapid antigen tests (RATs) with
 278 syndromic modelling yields better identification of COVID-19 cases than either
 279 diagnostic in isolation. These gains in performance are mirrored across metrics of
 280 prediction, generic classification and scenario-specific classification. The biggest
 281 improvement is seen in Scenario 2 (“Rising Cases”) which was developed around
 282 the current situation in Bangladesh (see Table 1 where the pandemic is once
 283 again accelerating. In this scenario, the combined data model (Model Class 3)
 284 false negative rate is 26 (IQR: 24, 29) percentage points lower than that of the RAT
 285 only model (Model Class 1). Although the syndromic only model (Model Class
 286 2) matches the combined models false negative rate, its false positive rate is 31
 287 (IQR: 29, 34) percentage points higher.

288 In a country where there are currently 15 000 new cases being identified
 289 every day, these improvements are non-trivial, representing tens of thousands of
 290 daily cases that would otherwise be missed. Furthermore, this boost in diagnostic
 291 performance is achieved with data that are already being collected in Bangladesh
 292 and other low- and middle- income countries (LMICs). Outwith developing and
 293 rerunning the models presented in this paper, these improvements are essentially
 294 cost-free and eminently scalable.

295 The pattern is similar in epidemiological Scenarios 1 (“Agnostic”) and 3
 296 (“Low-Level Cases”), with the combined model class performing performing
 297 equally well or better than the other two classes (Figure ??). These three
 298 scenarios only offer snapshots of performance. An indication of how these models
 299 will perform under any condition can be obtained by comparing the more generic
 300 model performance metrics for prediction and classification (Figures 3 and ??,
 301 respectively). These figures demonstrate both the added flexibility of the more
 302 complex model classes that allow them to be tailored to specific needs and
 303 the need to combine the high-quality but inflexible RAT results with the more
 304 flexible but lower quality syndromic data.

305 The final symptoms chosen through model selection should be interpreted
 306 cautiously. These models were developed for prediction and classification in a
 307 unique sub-population: CST-identified, symptomatic patients. The symptoms
 308 and risk factors retained in the model classes differed, despite these data being
 309 collected over a short time period from the same population. These differences
 310 may point to mechanisms by which CST-identified and RAT-positive patients
 311 differ from other groups. Of particular interest is whether individuals that are

missed by RAT are less infectious, which could be explored by using viral load measured as Threshold Cycle (Ct) values from the RT-PCR [23].

Our methodology has been developed using a large sample size drawn under field-realistic conditions and has thus developed with the practicalities of mass deployment in mind. Improving case identification using statistical methods allows us to update our diagnostic process in real-time, allowing rapid adaptation to new variants or even new diseases. The modelling frameworks we have used are also sufficiently flexible to accommodate new data sources (such as background case numbers) or changes in the local relative costs of false positives and false negatives.

Naturally, these strengths have complementary limitations. Our models require updating in real-time and can only achieve good performance if the validation data are of good quality. Similarly, targeting misdiagnosis rates is only sensible if those rates properly reflect local conditions which can be challenging. While these limitations should be seriously considered, we believe that the alternatives simply hide these problems. We choose to make these decisions explicitly to allow them to be more readily challenged, researched and improved upon. These challenges represent promising new avenues for impactful research that improve our understanding of estimating misdiagnosis rate trade offs and how to translate sample population findings to target populations.

We believe that combined syndromic and rapid antigen testing approach is the most promising method for large-scale testing in LMICs for COVID-19 at present. 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 framework we outline above is adaptable for other diagnostic problems. Malaria, schistosomiasis, rabies and many other diseases are all currently monitored either sparsely with gold-standard methods (such as RT-PCR, autopsies, fluorescent antibody testing) or at a large scale with more error-prone methods (RATs, blood smears, egg counts, differential diagnosis).

The management of global pandemics can only be done with testing at scale. 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.

356 7. Acknowledgements (69 words)

357 We would like to thank members of the community support teams in
 358 Bangladesh who have provided essential services throughout the pandemic.
 359 Earlier drafts of this manuscript benefited from the input of Paul Johnson,
 360 Daniel Haydon, Anne-Sophie Bonnet-Lebrun, Luca Nelli, Crinan Jarrett, Rita
 361 Claudia Cardoso Ribeiro, Halfan Ngowo, Heather McDevitt and Gina Bertolacci.
 362 The University of Glasgow COVID-19 in LMICs Group provided the environment
 363 in which to develop this work.

364 References (Max 30)

- 365 [1] Dramé M, Teguo MT, Proye E, Hequet F, Hentzien M, Kanagaratnam L,
 et al. Should RT-PCR be considered a gold standard in the diagnosis of
 366 covid-19? *Journal of Medical Virology* 2020.
- 367 [2] Corman VM, Landt O, Kaiser M, Molenkamp R, Meijer A, Chu DK, et al.
 Detection of 2019 novel coronavirus (2019-nCoV) by real-time RT-PCR.
 368 *Eurosurveillance* 2020;25:2000045.
- 369 [3] Tahamtan A, Ardebili A. Real-time RT-PCR in COVID-19 detection:
 Issues affecting the results. *Expert Review of Molecular Diagnostics*
 370 2020;20:453–4.
- 371 [4] Chowdhury R, Luhar S, Khan N, Choudhury SR, Matin I, Franco OH.
 Long-term strategies to control COVID-19 in low and middle-income
 countries: An options overview of community-based, non-pharmacological
 372 interventions. *European Journal of Epidemiology* 2020;35:743–8.
- 373 [5] Vandenberg O, Martiny D, Rochas O, Belkum A van, Kozlakidis Z. Con-
 siderations for diagnostic COVID-19 tests. *Nature Reviews Microbiology*
 374 2021;19:171–83.
- 375 [6] Cash R, Patel V. Has COVID-19 subverted global health? *The Lancet*
 376 2020;395:1687–8.
- 377 [7] Olalekan A, Iwalokun B, Akinloye OM, Popoola O, Samuel TA, Akinloye
 O. COVID-19 rapid diagnostic test could contain transmission in low-
 and middle-income countries. *African Journal of Laboratory Medicine*
 378 2020;9:1–8.
- 379 [8] Linares M, Pérez-Tanoira R, Carrero A, Romanyk J, Pérez-García F,
 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.
 380 *Journal of Clinical Virology* 2020;133:104659.
- 381 [9] 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.
 382 *The Lancet Infectious Diseases* 2021.

- 383 [10] 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*
384 *Clinical Virology* 2020;129:104500.
- 385 [11] 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*
386 *Research* 2020;7:1–23.
- 387 [12] 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*
388 *Infection Control* 2021;49:685–9.
- 389 [13] Undurraga EA, Chowell G, Mizumoto K. COVID-19 case fatality risk
by age and gender in a high testing setting in latin america: Chile,
390 march–august 2020. *Infectious Diseases of Poverty* 2021;10:1–1.
- 391 [14] Wenham C, Smith J, Morgan R. COVID-19: The gendered impacts of
the outbreak. *The Lancet* 2020;395:846–8.
- 392 [15] Mayorga L, Samartino CG, Flores G, Masuelli S, Sánchez MV, Mayorga
LS, et al. A modelling study highlights the power of detecting and
isolating asymptomatic or very mildly affected individuals for COVID-19
393 epidemic management. *BMC Public Health* 2020;20:1–1.
- 394 [16] Garry S, Abdelmagid N, Baxter L, Roberts N, Waroux O le P de, Ismail
S, et al. Considerations for planning COVID-19 treatment services in
395 humanitarian responses. *Conflict and Health* 2020;14:1–1.
- 396 [17] Surkova E, Nikolayevskyy V, Drobniewski F. False-positive COVID-19
results: Hidden problems and costs. *The Lancet Respiratory Medicine*
397 2020;8:1167–8.
- 398 [18] West CP, Montori VM, Sampathkumar P. COVID-19 testing: The threat
of false-negative results. *Mayo clinic proceedings*, vol. 95, Elsevier; 2020,
399 p. 1127–9.
- 400 [19] Albert JH, Chib S. Bayesian analysis of binary and polychotomous re-
sponse data. *Journal of the American Statistical Association* 1993;88:669–
401 79.
- 402 [20] Carpenter B, Gelman A, Hoffman MD, Lee D, Goodrich B, Betancourt M,
et al. Stan: A probabilistic programming language. *Journal of Statistical*
403 *Software* 2017;76:1–32.
- 404 [21] Gneiting T, Raftery AE. Strictly proper scoring rules, prediction, and
estimation. *Journal of the American Statistical Association* 2007;102:359–
405 78.
- 406 [22] Hoo ZH, Candlish J, Teare D. What is an ROC curve? 2017.
- 407
- 408

- 409 [23] 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.
410 Clinical Microbiology and Infection 2021;27:472–e7.