

1 Clinical predictors for etiology of acute diarrhea in children in resource-limited settings

2

3 Benjamin Brintz<sup>1</sup>, Joel Howard<sup>2</sup>, Benjamin Haaland<sup>3</sup>, James A. Platts-Mills<sup>4</sup>, Tom Greene<sup>3</sup>,

4 Adam C. Levine<sup>5</sup>, Eric Nelson<sup>6</sup>, Andrew T. Pavia<sup>2</sup>, Karen L. Kotloff<sup>7</sup>, and Daniel T.

5 Leung<sup>1,8\*</sup>

6

7 <sup>1</sup>*Division of Infectious Diseases, University of Utah, Salt Lake City, USA*

8 <sup>2</sup>*Division of Pediatric Infectious Diseases, University of Utah, Salt Lake City, USA*

9 <sup>3</sup>*Division of Biostatistics, University of Utah, Salt Lake City, USA*

10 <sup>4</sup>*Division of Infectious Diseases and International Health, University of Virginia,*

11 *Charlottesville, USA*

12 <sup>5</sup>*Department of Emergency Medicine, Brown University, Providence, USA*

13 <sup>6</sup>*Departments of Pediatrics and Environmental and Global Health, University of Florida,*

14 *Gainesville, USA*

15 <sup>7</sup>*Division of Infectious Disease and Tropical Pediatrics, Center for Vaccine Development*

16 *and Global Health, University of Maryland School of Medicine, Baltimore,*

17 *USA*

18 <sup>8</sup>*Division of Microbiology and Immunology, University of Utah, Salt Lake City, USA*

19

20 \*Corresponding author

21 E-mail: [Daniel.Leung@utah.edu](mailto:Daniel.Leung@utah.edu)

22

23 Short Title: Predictors of diarrhea etiology

24

25

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

## 26 **Abstract**

27 *Background.* Diarrhea is one of the leading causes of childhood morbidity and mortality  
28 in lower- and middle-income countries. In such settings, access to laboratory diagnostics are  
29 often limited, and decisions for use of antimicrobials often empiric. Clinical predictors are a  
30 potential non-laboratory method to more accurately assess diarrheal etiology, the knowledge  
31 of which could improve management of pediatric diarrhea.

32 *Methods.* We used clinical and quantitative molecular etiologic data from the Global  
33 Enteric Multicenter Study (GEMS), a prospective, case-control study, to develop predictive  
34 models for the etiology of diarrhea. Using random forests, we screened the available  
35 variables and then assessed the performance of predictions from random forest regression  
36 models and logistic regression models using 5-fold cross-validation.

37 *Results.* We identified 1049 cases where a virus was the only etiology, and developed  
38 predictive models against 2317 cases where the etiology was known but non-viral (bacterial,  
39 protozoal, or mixed). Variables predictive of a viral etiology included age, season, height-for-  
40 age z-score (HAZ), bloody diarrhea, and vomiting. Cross-validation suggests an AUC of  
41 0.825 can be achieved with a parsimonious model of 5 variables, achieving a specificity of  
42 0.85, a sensitivity of 0.59, a NPV of 0.82 and a PPV of 0.64.

43 *Conclusion.* Predictors of the etiology of pediatric diarrhea can be used by providers in  
44 low-resources setting to inform clinical decision-making. The use of non-laboratory methods  
45 to diagnose viral causes of diarrhea could reduce inappropriate antibiotic prescription  
46 worldwide.

47

48 **Keywords.** diarrhea; clinical prediction; etiology; low- and middle- income countries; GEMS

49

50 **Author Summary:**

51 Diarrhea is one of the leading causes of death in young children worldwide. In low-resource  
52 settings, diarrhea testing is not available or too expensive, and the decision to prescribe  
53 antibiotics is often made without testing. Using clinical information to predict which cases  
54 are caused by viruses, and thus wouldn't need antibiotics, would help to improve appropriate  
55 use of antibiotics. We used data from a large study of childhood diarrhea, paired with  
56 advanced statistical methods including machine learning, to come up with the top clinical  
57 factors that could predict a viral cause of diarrhea. We compared 1049 cases where a virus  
58 was the only cause, with 2317 cases where the cause was known but not a virus. We found  
59 that age, season, nutritional status defined by height, blood diarrhea, and vomiting, were the  
60 clinical factors most predictive of whether the diarrhea was caused by a virus or not. We  
61 found that, using just those 5 factors, we were able to predict a viral cause with good  
62 accuracy. Our findings can be used by doctors to guide the appropriate use of antibiotics for  
63 diarrhea in children.

64

65

## 66 **Introduction**

67 Diarrhea is one of the leading causes of childhood morbidity and mortality in lower- and  
68 middle-income countries (LMICs) and is among the most common reasons for admission  
69 into a health facility [1]. Treatment of diarrhea is commonly empiric, with antibiotic  
70 prescription mostly based on clinical suspicion for a bacterial etiology. In resource-limited  
71 settings, laboratory etiological diagnosis is rarely made due to cost constraints or availability.  
72 A large number of patients with viral illness are prescribed antibiotics inappropriately, and  
73 the rate of use varies widely by country and setting [2]. This inappropriate use of  
74 antimicrobials can lead to toxicity, increased costs of care, and development of resistance [3].  
75 Thus, methods providing clinical decision support that accurately assesses diarrhea etiology  
76 and reduces reliance on laboratory testing are needed. Recently, tools for decision making  
77 and clinical prediction have been bolstered by the exploration of machine learning methods  
78 such as random forests, neural networks, and support vector machines [4].

79 The availability of molecular diagnostics in recent years has enabled accurate  
80 determination of etiology for pediatric diarrhea. In several large studies in LMICs, this has  
81 been used for estimating the population-based burden of various diarrheal pathogens [5-7].  
82 While etiologies of diarrhea are now better-understood, there remains a gap in knowledge  
83 regarding clinical predictors for improving clinical decision making in the setting of  
84 infectious diarrhea. In this study, we use data from the Global Enteric Multicenter Study  
85 (GEMS) [5] to examine clinical diagnostic predictors of diarrhea etiology. We provide a brief  
86 introduction to the data and data processing steps, describe our variable screening and model  
87 fitting approach, present the results of our predictive models, and discuss the implications of  
88 such models.

89

## 90 **Methods**

### 91 **Study Design and Settings**

92 GEMS is a prospective, case-control study that took place from 2007-2011 in 7 countries in  
93 Africa and South Asia (Figure S2). There were 9439 children with moderate-to-severe  
94 diarrhea (MSD) enrolled at local health care centers along with 1 to 3 matched non-diarrheal  
95 controls. An acute episode of diarrhea was defined as MSD if it fulfilled at least one of the  
96 following criteria: sunken eyes; loss of skin turgor; intravenous hydration administered or  
97 prescribed; visible blood in stool or parental report; or admission to hospital with diarrhea or  
98 dysentery or advising hospitalization. At enrollment, a stool sample was taken from each  
99 child to identify enteropathogens along with clinical information, including demographic,  
100 anthropometric, and clinical history. Methods for GEMS have been described in detail  
101 previously [5, 8, 9]. Because pathogen nucleic acids are frequently detected by PCR in  
102 children without diarrhea, we used the quantitative real-time PCR-based (qPCR) majority  
103 attribution models developed by Liu et al [6] to assign etiology of diarrhea. We derived site-  
104 and age- specific attributable fractions (AF<sub>e</sub>) for each episode, and used a cut-off of greater  
105 than 0.5 to indicate attribution of a pathogen to a particular episode. We defined viral  
106 etiology as majority attribution of the diarrhea episode by viral pathogen(s) only (i.e.  
107 excluding any co-infections with bacteria or protozoa). We defined other known etiologies as  
108 having a majority attribution of diarrhea episode by at least one other non-viral pathogen.  
109 Additionally, we defined a bacterial etiology as attribution of the diarrhea episode by any  
110 bacterial pathogen, including cases in which more than one pathogen was attributed (i.e.  
111 bacteria and virus, or bacteria and protozoa, or multiple bacteria). For patients with unknown  
112 etiologies, we presume there is an infectious cause to their diarrhea that we are not detecting,  
113 and excluded these cases from our predictive model.

114 We used the patient's clinical symptoms, epidemiologic, and anthropometric data at  
115 presentation as potential predictors of etiology. We used standard guidelines from the  
116 transparent reporting of a multivariable prediction model for individual diagnosis (TRIPOD)  
117 to develop our prediction model [10]. We focused on the prediction of a viral etiology of

118 acute diarrhea versus all other known etiologies since this would allow clinicians to  
119 comfortably withhold antibiotics. We additionally looked at the prediction of any bacterial  
120 pathogen as a way to determine if follow-up testing may be helpful in ambiguous cases.

## 121 **Data Processing**

122 We performed all data processing and analyses using R [11]. Starting with over 1000  
123 variables collected, we excluded all variables which would not be available at the time of  
124 presentation. Questions which had very few responses in certain categories (<10) were re-  
125 grouped into an “other” category as appropriate. Some variables only had 1 or 2 responses in  
126 a given category and those patients were removed from the dataset when grouping into an  
127 “other” category was not possible. For instance, only 5 patients responded they “Don’t  
128 Know” when asked if they had any blood in their stool since the illness began. We  
129 maximized the utility of the modeling process by removing highly collinear and similar  
130 variables (e.g. weight-, BMI, and BMI-for-age z-scores), while keeping variables that are  
131 clinically accessible, before observing any measurement of etiology. These steps left 156  
132 potential predictor variables for analysis.

133 In addition to the information from the GEMS survey, we developed a season variable  
134 using temperature and rain information from NOAA weather stations close to the health  
135 centers and with data during the GEMS time period<sup>12</sup>. We defined a rainy season day as a  
136 day having a center-aligned 1-month moving rain average greater than the overall rain  
137 average within the study period. We defined a hot season day as a day having a center-  
138 aligned 1-month moving temperature average greater than the overall temperature average  
139 within the study period.

## 140 **Statistical Modeling and Assessment**

141 We used random forests as a screening step to obtain an order of variable importance toward  
142 the goal of building a parsimonious model. The random forest method uses an ensemble  
143 approach by generating multiple decision trees (1000 trees, throughout) and approaching  
144 variable importance by determining a reduction in mean squared prediction error for each

145 variable on the “out-of-bag” samples (or testing samples) created while bootstrapping the  
146 data. During this step, categorical variables are treated as a single variable rather than a  
147 variable for each level.

148 We used 5-fold cross-validation to attain an estimate of generalizable model  
149 performance. For each cross-validation iteration, we took the order of the importance  
150 measure from the screening step to determine which variables we used to fit separate logistic  
151 regression model and random forest models with various predictor subset sizes. Subsets  
152 examined were sizes 1 through 10, 15, 20, 30, 40, and 50. In each iteration of cross-  
153 validation we made predictions on the test set and obtained measures of performance: the  
154 receiver operating characteristic (ROC) curve, and area under the ROC curve (AUC), also  
155 known as the C-statistic, along with AUC confidence intervals [13]. For a diagnostic  
156 threshold balancing the relative costs of false positives and false negatives, we calculated the  
157 positive predictive value (PPV) and the negative predictive value (NPV) as functions of the  
158 derived sensitivity and specificity of the prediction, using the prevalence of the  
159 corresponding etiology in GEMS.

## 160 **Ethics approval**

161 The GEMS study protocol was approved by ethics committees at the University of Maryland,  
162 Baltimore and at each field site. Parents or caregivers of participants provided written  
163 informed consent, and a witnessed consent was obtained for illiterate parents or caretakers.

## 164 **Results**

165 Of the 3366 patients in the GEMS study, 9439 patients had MSD and are included in this  
166 analysis (Figure S3), 1049 had a viral etiology and 2069 had a bacterial etiology (Table 1).  
167 Using random forest screening, we found that age, season, bloody diarrhea, height-for-age z-  
168 score (HAZ), and vomiting were the five variables most predictive of a viral etiology (Table  
169 2), and that top predictive variables for bacterial etiology were similar (Supplemental Table  
170 S1).

171 When we performed 5-fold cross-validated logistic regression and random forest models,  
172 the average AUC across 100 random iterations of cross-validation ranged from 0.71 (1  
173 variable) to 0.84 (7 or more variables) for prediction of viral etiology (Figure 1) with similar  
174 results for bacterial etiology (Figure S4). To demonstrate the direction and magnitude of the  
175 effect of the top 10 variables from variable importance screening by fitting a logistic  
176 regression on the entire data set (Table 3). Lower age, a higher HAZ, more vomiting, no  
177 blood in the stool, and a dry/cold season, were associated with viral etiology. As expected,  
178 the opposite associations were found for bacterial etiology (Supplemental Table S2). We  
179 found similar results in a sensitivity analysis with rotavirus removed, though some effect  
180 magnitudes were reduced. To estimate the achievable sensitivity and specificity by each  
181 model at various predictor sizes, we generated ROC curves from cross-validation, and found  
182 that using a parsimonious model of 5 variables, we achieved a specificity of 0.85 and a  
183 sensitivity of close to 0.60 for prediction of viral etiology (Figure 2). For predicting a  
184 bacterial cause, our models achieved a sensitivity of 0.85 and a specificity of 0.63 (Figure  
185 S5). Using the prevalence of viral etiology in GEMS, our prediction model had a NPV of  
186 0.82 and a PPV of 0.64.

187  
188 Figure 1: Average AUC and 95% CIs from 100 iterations of cross-validation for both a  
189 logistic regression (LR) and random forest (RF) as the number of variables in the model  
190 increases and inset shows zoomed in graphs of 1 through 10 variables.



191

192 Figure 2: Interpolated estimates of ROC curves from the cross-validation for logistic  
193 regression and random forest models with variable sizes of 5, 10, and 20. The faded dashed  
194 lines represent examples of how we could achieve a sensitivity of 0.6 and a specificity of  
195 0.85 for prediction of viral etiology.

196

197 When we examined the predictors associated with viral etiology for each of the 7 sites in  
198 GEMS by filtering the entire dataset by site, we found all had a similar order of variable  
199 importance with some minor differences (Table 4). We then looked at the performance of the  
200 prediction model filtered for specific sites and specific continents within each cross-  
201 validation iteration's test set, and found that at Asian sites the predictions had an AUC almost  
202 0.07 better than African sites on average. Looking at individual sites, in Kenya the model  
203 predictions had the worst average AUC while Bangladesh had the best average AUC. Across  
204 all sites, the AUC of a 5-variable model was similar to a 10-variable model with less than  
205 0.02 lower average AUC.

206 We then performed an external validation by testing the logistic regression on each site  
207 individually following training on the other sites in the same continent, and found  
208 performance metrics similar to the cross-validation results, with AUC ranging from 0.65 to  
209 0.92 across the seven sites. As with the internal cross-validation, we found 5-variable models  
210 to have similar performance to 10-variable models. We found similar results for the bacterial  
211 etiology prediction (Supplemental Table S3).

212

## 213 Discussion

214 Our use of data from GEMS, which involved 3366 diarrheal episodes with known etiology in  
215 7 countries and with over 150 clinically-relevant parameters collected for each episode,  
216 allowed for a robust analysis that revealed the ability of clinical variables alone to predict  
217 diarrheal etiology with a high degree of accuracy. Using machine learning algorithms, we  
218 found that a model with just 5 variables (age, season, HAZ, bloody diarrhea, and vomiting),  
219 could accurately predict viral etiology, with a cross-validated AUC of 0.825. Translation of  
220 these findings towards clinical decision making has the potential to improve management,  
221 including appropriate antibiotic use, in LMICs.

222

223 Previous studies predicting etiology of diarrheal illness<sup>14-17</sup>, have been limited by the low  
224 number of participants, amount of clinical data collected, pathogen variety, number of  
225 pathogens detected, method of detection, lack of controls without diarrhea, single center  
226 design, and the need for stool testing. Etiological prediction is particularly challenging in  
227 LMIC settings, where multi-pathogen detection is common in children with diarrhea, and  
228 presumed pathogens can be isolated from asymptomatic individuals in up to 50% of study  
229 controls<sup>18</sup>. New molecular diagnostic methods used on the GEMS samples involved a  
230 quantitative assessment of 32 potential pathogens, with matched case-control pairs, to ascribe  
231 an etiological attributable fraction (AF<sub>e</sub>) for each episode. This quantitative method, in  
232 context of a case-control study, is thus able to account for the high rate of asymptomatic  
233 detection of pathogens by molecular testing in children in LMICs, which can confound the  
234 attribution of etiology. Using these data, we built several models to evaluate the effect of  
235 clinical indicators on whether children presenting with acute diarrhea had a viral etiology (or  
236 bacterial etiology). We showed that AUCs improved for the first 7 variables but thereafter  
237 the addition of more variables did not improve the model. Notably, we found that an AUC of  
238 0.825 could be achieved with 5 variables, enabling the translation of this predictive model to  
239 a parsimonious rule which could be used in clinical decision-support.

240

241 When considering sensitivity and specificity in the context of diarrheal etiology, we  
242 assumed a high specificity target for prediction of “viral only” etiology (Figure 2), and  
243 similarly, a high sensitivity target for bacterial etiology (Figure S4), both of which would  
244 minimize the risk of not giving antibiotics to a child with a bacterial infection. While current  
245 WHO guidelines recommend antibiotics only for children with dysentery and for children  
246 with acute water diarrhea (AWD) with severe dehydration in cholera endemic regions, there  
247 is evidence suggesting treatment of non-dysenteric *Shigella* infections may be beneficial [19,  
248 20]. Our prediction model showed that for predicting a viral etiology, for a desired specificity  
249 of 0.85, we achieved a sensitivity of 0.59. We found that the most significant predictors for  
250 differentiating viral from other etiologies were: age, HAZ, season, bloody diarrhea, and  
251 vomiting. Vomiting, a higher HAZ, and dry/cold season were evidence towards a viral  
252 etiology, while an older age and bloody diarrhea were evidence against a viral etiology.

253 The predictors we identified are consistent with those of previous studies. Bloody  
254 diarrhea as a predictor of a bacterial cause of diarrhea, especially for shigellosis, has been  
255 well established<sup>14-17, 21-23</sup> and informs the IMCI guidelines that dysentery be treated with  
256 antibiotics. Vomiting as a predictor of a viral process has similarly been shown in previous  
257 studies<sup>14, 16</sup>. It is well established that younger children have a higher incidence of diarrhea<sup>24</sup>  
258 and some studies have suggested that younger age is also more indicative of a viral process<sup>16,</sup>  
259 <sup>22, 24-26</sup>. We showed that age was the most important predictor with mean age of viral case  
260 being 13.0 months, and 22.1 months for bacterial cases.

261  
262 Using data gathered from NOAA weather stations proximal to our study sites during the  
263 study period, we were able to develop seasonal variables based on temperature and rainfall.  
264 We show that a viral etiology of diarrhea is associated with a drier, colder climate, consistent  
265 with observation from previous studies from the USA<sup>16</sup> and India<sup>26</sup>. The positive association  
266 of anthropometrics (higher HAZ and MUAC) with viral etiology may suggest that improved  
267 nutrition is more protective of a bacterial than a viral process. Symptoms found in earlier  
268 studies to be predictive of etiology, but which did not improve predictive performance in our  
269 analysis, include fever, number of stools per day, duration of diarrhea, and presence of

270 mucous<sup>14-17, 23</sup>. Similarly, variables related to hygiene and sanitation did not help with  
271 prediction of etiology.

272

273 Given that GEMS was conducted in 7 countries across Africa and Asia, we examined the  
274 model performance across sites. We found that the model attained an average AUC of about  
275 0.86 in Asian sites and about 0.79 in African sites, likely due to poor performance of the  
276 model in Kenya and good performance in Bangladesh. This suggests that additional external  
277 validation will be necessary to assess both performance and generalizability. Indeed, even  
278 within continent, countries had varying AUCs. We also found that, when externally validated  
279 against other sites from the same continent, use of five variables achieve similar AUC as use  
280 of 10 variables. Future studies should aim to capture country- or continent-specific trends so  
281 that outbreaks or volatility can be accounted for in the predictions.

282

283 Our study has a number of limitations. First, our predictive model does not distinguish  
284 between different bacterial etiologies, which may require different therapy. Additionally, it  
285 does not predict for parasitic infections. In GEMS [6], a number of bacterial pathogens had  
286 few to no cases detected using  $A_{Fe} > 0.5$ , including EHEC, Yersinia, LT ETEC, EAEC,  
287 atypical EPEC, and Clostridium difficile. This was due to these organisms' presence in  
288 control children without diarrhea, making attribution difficult. While it is possible that these  
289 could have co-occurred with a viral pathogen, there is limited evidence that antibiotic  
290 treatment of these etiologies would be beneficial in this setting. External validation is  
291 essential for this and all clinical prediction models, as demonstrated by our heterogenous  
292 result by continent. GEMS was conducted before the widespread use of rotavirus vaccine and  
293 rotavirus was the dominant viral pathogen; thus, the model will need to be validated in  
294 settings where rotavirus vaccination campaigns have had substantial impact. Finally, we note  
295 that because variable selection was used before fitting the logistic regression model, the role  
296 of the variables in terms of p-values and confidence intervals may be over-stated.

297

298 In conclusion, utilizing a large number of cases and quantitative molecular methods of  
299 pathogen detection with etiologic attribution based on a case-control study, we showed that  
300 etiology prediction could be attained for episodes of acute diarrhea with as few as 5  
301 variables. Our findings confirm previously considered predictors of viral etiology including  
302 lack of bloody diarrhea, vomiting, younger age, and a dry and cool climate, and reveal  
303 additional predictors of viral etiology associated with anthropometric measures. These  
304 findings have the potential to provide clinicians in lower-resource settings with better  
305 informed clinical decision making, including identification of the subset of children from  
306 whom antibiotics may be safely withheld and a group who may benefit from antimicrobials  
307 and/or adjunctive microbiologic testing.

## References

- [1] Walker CLF, Rudan I, Liu L, Nair H, Theodoratou E, Bhutta ZA, et al. Global burden of childhood pneumonia and diarrhoea. *The Lancet*. 2013;381(9875):1405– 1416.
- [2] Rogawski ET, Platts-Mills JA, Seidman JC, John S, Mahfuz M, Ulak M, et al. Use of antibiotics in children younger than two years in eight countries: a prospective cohort study. *Bulletin of the World Health Organization*. 2017;95(1):49.
- [3] World Health Organization. Antimicrobial resistance: global report on surveillance. World Health Organization; 2014.
- [4] Eom JH, Kim SC, Zhang BT. AptaCDSS-E: A classifier ensemble-based clinical decision support system for cardiovascular disease level prediction. *Expert Systems with Applications*. 2008;34(4):2465–2479.
- [5] Kotloff KL, Nataro JP, Blackwelder WC. Burden and aetiology of diarrhoeal disease in infants and young children in developing countries (the Global Enteric Multicenter Study, GEMS): a prospective, case-control study. *The Lancet*. 2013;382(9888):209– 222.
- [6] Liu J, Platts-Mills JA, Juma J, Kabir F, Nkeze J, Okoi C, et al. Use of quantitative molecular diagnostic methods to identify causes of diarrhoea in children: a reanalysis of the GEMS case-control study. *The Lancet*. 2016;388(10051):1291–1301.
- [7] Platts-Mills JA, Liu J, Rogawski ET, Kabir F, Lertsethtakarn P, Sigua M, et al. Use of quantitative molecular diagnostic methods to assess the aetiology, burden, and clinical characteristics of diarrhoea in children in low-resource settings: a reanalysis of the MAL-ED cohort study. *The Lancet Global Health*. 2018;6(12):e1309–e1318.
- [8] Kotloff KL, Blackwelder WC, Nasrin D, Nataro JP, Farag TH, van Eijk A, et al. The Global Enteric Multicenter Study (GEMS) of diarrheal disease in infants and young children in developing countries: epidemiologic and clinical methods of the case/control study. *Clinical infectious diseases*. 2012;55(suppl\_4):S232–S245.

- [9] Panchalingam S, Antonio M, Hossain A, Mandomando I, Ochieng B, Oundo J, et al. Diagnostic microbiologic methods in the GEMS-1 case/control study. *Clinical infectious diseases*. 2012;55(suppl 4):S294–S302.
- [10] Moons KG, Altman DG, Reitsma JB, Ioannidis JP, Macaskill P, Steyerberg EW, et al. Transparent Reporting of a multivariable prediction model for Individual Prognosis or Diagnosis (TRIPOD): explanation and elaboration. *Annals of internal medicine*. 2015;162(1):W1–W73.
- [11] R Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria; 2018. Available from: <https://www.R-project.org/>.
- [12] Chao DL, Roose A, Roh M, Kotloff KL, Proctor JL. The seasonality of diarrheal pathogens: A retrospective study of seven sites over three years. *BioRxiv*. 2019;p. 541581.
- [13] LeDell E, Petersen M, van der Laan M. Computationally efficient confidence intervals for cross-validated area under the ROC curve estimates. *Electronic journal of statistics*. 2015;9(1):1583.
- [14] DeWitt TG, Humphrey KF, McCarthy P. Clinical predictors of acute bacterial diarrhea in young children. *Pediatrics*. 1985;76(4):551–556.
- [15] Fontana M, Zuin G, Paccagnini S, Ceriani R, Quaranta S, Villa M, et al. Simple clinical score and laboratory-based method to predict bacterial etiology of acute diarrhea in childhood. *The Pediatric infectious disease journal*. 1987;6(12):1088– 1091.
- [16] Klein EJ, Boster DR, Stapp JR, Wells JG, Qin X, Clausen CR, et al. Diarrhea etiology in a children’s hospital emergency department: a prospective cohort study. *Clinical Infectious Diseases*. 2006;43(7):807–813.
- [17] Velasco AC, de Agüero Barrio MG. Clinical and laboratory indicators of etiology of diarrhea. *Anales espanoles de pediatria*. 1992;36(6):423–427.

- [18] van Coppentraet LB, Dullaert-de Boer M, Ruijs G, Van der Reijden W, van der Zanden A, Weel J, et al. Case–control comparison of bacterial and protozoan microorganisms associated with gastroenteritis: application of molecular detection. *Clinical Microbiology and Infection*. 2015;21(6):592–e9.
- [19] Tickell KD, Brander RL, Atlas HE, Pernica JM, Walson JL, Pavlinac PB. Identification and management of Shigella infection in children with diarrhoea: a systematic review and meta-analysis. *The Lancet Global Health*. 2017;5(12):e1235–e1248.
- [20] Rogawski ET, Liu J, Platts-Mills JA, Kabir F, Lertsethtakarn P, Sigua M, et al. Use of quantitative molecular diagnostic methods to investigate the effect of enteropathogen infections on linear growth in children in low-resource settings: longitudinal analysis of results from the MAL-ED cohort study. *The Lancet Global Health*. 2018;6(12):e1319–e1328.
- [21] Singh T, Verma M, Chhatwal J, Chacko B, Kaur H, Prabhakar H. Predictive utility of clinical and stool parameters in bacterial diarrhoea in children. *Indian journal of medical sciences*. 1995;49(12):285–290.
- [22] Suwatano O. Acute diarrhea in under five-year-old children admitted to King Mongkut Prachomkiao Hospital, Phetchaburi province. *Journal of the Medical Association of Thailand= Chotmaihet Thangphaet*. 1997;80(1):26–33.
- [23] Denno DM, Stapp JR, Boster DR, Qin X, Clausen CR, Del Beccaro KH, et al. Etiology of diarrhea in pediatric outpatient settings. *The Pediatric infectious disease journal*. 2005;24(2):142–148.
- [24] Saidi SM, Lijima Y, Sang WK, Mwangudza AK, Oundo JO, Taga K, et al. Epidemiological study on infectious diarrheal diseases in children in a coastal rural area of Kenya. *Microbiology and immunology*. 1997;41(10):773–778.



- [25] Baselga CA, Alonso MG, Bernal MS, Bueno GL, Bueno ML, Gracia MC, et al. Bacterial diarrhea in infancy: epidemiologic study of 256 cases. *Anales espanoles de pediatria*. 1991;34(3):203–206.
- [26] Niyogi S, Saha M, De S. Enteropathogens associated with acute diarrhoeal diseases. *Indian journal of public health*. 1994;38(2):29.
- [27] Santos FS, Santos FCS, Santos LHd, Leite AM, Mello DFd. Breastfeeding and protection against diarrhea: an integrative review of literature. *Einstein (São Paulo)*. 2015;13(3):435–440.
- [28] Quigley MA, Kelly YJ, Sacker A. Breastfeeding and hospitalization for diarrheal and respiratory infection in the United Kingdom Millennium Cohort Study. *Pediatrics*. 2007;119(4):e837–e842.
- [29] Arifeen S, Black RE, Antelman G, Baqui A, Caulfield L, Becker S, et al. Exclusive breastfeeding reduces acute respiratory infection and diarrhea deaths among infants in Dhaka slums. *Pediatrics*. 2001;108(4):E67.
- [30] Yoon PW, Black RE, Moulton LH, Becker S. Effect of not breastfeeding on the risk of diarrheal and respiratory mortality in children under 2 years of age in Metro Cebu, The Philippines. *American Journal of Epidemiology*. 1996;143(11):1142–1148.

## Supplementary Figure Legends

### S1 Checklist: STROBE Checklist

Figure S2: The left map shows the locations of the 4 study sites in Africa. Right map shows the locations of 3 study sites in South Asia. The map was generated using the `get_map` and `ggmap` functions in R version 3.6.1.

Figure S3: Average AUC and 95% CIs from 100 iterations of cross-validation for both a logistic regression (LR) and random forest (RF) as the number of variables in the model increases and inset shows zoomed in graphs of 1 through 10 variables.

Figure S4: Consort diagram of the reduction of patients from 22567 in the GEMS dataset to the 3366 cases in our study. Note that we only filtered out non-responses for response variables that were in the top 50 of our screening step.

Figure S5: Interpolated estimates of ROC curves from the cross-validation for logistic regression and random forest models with variable sizes of 5, 10, and 20. The faded dashed lines represent examples of how we could achieve a sensitivity of 0.85 and a specificity of 0.60 for any bacteria.

## Tables

Table 1: Number of cases attributed to each pathogen with an attributable fraction above 0.5.

Pathogen	Cases
<i>Adenovirus 40/41</i>	222
<i>Aeromonas</i>	59
<i>Astrovirus</i>	111
<i>C. jejuni/C. coli</i>	85
<i>Cryptosporidium</i>	301
<i>Cyclospora cayetanensis</i>	16
<i>Entamoeba histolytica</i>	29
<i>Helicobacter pylori</i>	131
<i>Isospora</i>	3
<i>Norovirus GII</i>	70
<i>Rotavirus</i>	967
<i>Salmonella</i>	67
<i>Sapovirus</i>	75
<i>Shigella/EIEC</i>	1376
<i>Vibrio cholerae</i>	152
<i>EAEC</i>	1
<i>ST-EPEC (STh)</i>	407
<i>Typical EPEC (bfpA)</i>	43
Occurrences	Cases
<i>Protozoal</i>	218
<i>Viral</i>	1049
<i>Viral-Protozoal</i>	30
<i>Bacterial</i>	1664
<i>Bacterial-Protozoal</i>	92
<i>Bacterial-Viral</i>	307
<i>Bacterial-Viral-Protozoal</i>	6

*Table 2: Rank of variable importance by reduction in residual sum of squares (RSS) using random forest regression.*

<b>Viral Etiology</b>	
<b>Variable Name</b>	<b>RSS Reduction</b>
Age	51.6
Season	29.0
Blood in stool	26.1
HAZ	24.7
Vomiting	23.0
Breastfed	22.0
MUAC	20.9
Resp. Rate	18.5
Wealth Index	18.3
Temperature	16.7

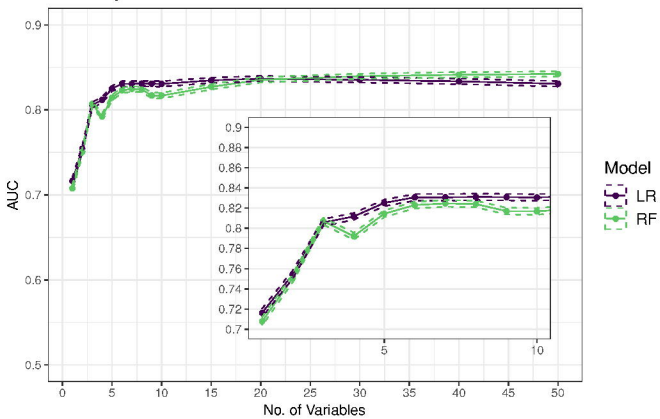
*Table 3: The odds ratios, 95% confidence interval, and p-value from a logistic regression model.*

Variable Name	Viral Only	
	Odds Ratios (95% CI)	P-value
Intercept	1.975 (0.053 – 72.894)	0.7117
Age (mo.)	0.956 (0.944 – 0.967)	<0.0001
Season		
Dry/Cold	Reference	
Rainy/Cold	0.197 (0.145 – 0.268)	<0.0001
Dry/Hot	0.304 (0.244 – 0.379)	<0.0001
Rainy/Hot	0.338 (0.268 – 0.426)	<0.0001
Blood in stool	0.129 (0.096 – 0.173)	<0.0001
HAZ	1.168 (1.081 – 1.262)	0.0001
Vomiting	2.383 (1.995 – 2.847)	<0.0001
Breastfed		
None	Reference	
Partially	2.359 (1.827 – 3.046)	<0.0001
Exclusively	2.400 (1.554 – 3.705)	0.0001
MUAC	1.031 (0.963 – 1.105)	0.3773
Resp. Rate (per min.)	0.990 (0.979 – 1.000)	0.0541
Wealth Index	1.066 (0.976 – 1.164)	0.1559
Temperature (°C)	0.988 (0.897 – 1.088)	0.8022

Table 4: The table contains both site-specific variable importance ordering and a cross-validated average overall AUC, AUC by country, and AUC by continent and confidence intervals from a 5 (bold) and 10 (ital.) variable logistic regression model for predicting a viral etiology with variables based on the overall variable importance. Lastly, it shows the AUC and a 95% confidence interval resulting from testing the logistic regression with variables based on the overall variable importance on each site individually following its training on the other countries in the same continent

	Africa				Asia		
Country Variable	The Gambia	Mali	Mozambique	Kenya	India	Bangladesh	Pakistan
1	Age	Age	Age	Age	Age	Age	Age
2	Season	Season	Season	HAZ	MUAC	Blood in stool	Breastfed
3	HAZ	Vomiting	Breastfed	MUAC	HAZ	Season	HAZ
4	Blood in stool	MUAC	HAZ	Resp. Rate	Season	Sunken Eyes	Resp. Rate
5	MUAC	HAZ	Temp.	Breastfed	Resp. Rate	Vomiting	MUAC
6	Temp.	Resp. Rate	MUAC	Temp.	Blood in stool	MUAC	Temp.
7	Resp. Rate	Breastfed	Resp. Rate	Wealth Index	Wealth Index	Rectal Straining	Wealth Index
8	Wealth Index	Wealth Index	Wealth Index	# Share Facility	# Share Facility	Temp.	Vomiting
9	People in House	Temp.	Vomiting	People in House	Temp.	HAZ	People in House
10	Vomiting	People in House	People in House	Days of Episode	People in House	Wealth Index	Blood in stool
Cntry AUCs	<b>0.850</b> (0.841-0.858) <i>0.847</i> (0.838-0.855)	<b>0.792</b> (0.780-0.803) <i>0.796</i> (0.785-0.807)	<b>0.833</b> (0.823-0.843) <i>0.839</i> (0.828-0.848)	<b>0.686</b> (0.674-0.698) <i>0.693</i> (0.681-0.705)	<b>0.812</b> (0.805-0.820) <i>0.813</i> (0.806-0.821)	<b>0.927</b> (0.922-0.933) <i>0.923</i> (0.918-0.929)	<b>0.788</b> (0.778-0.798) <i>0.801</i> (0.791-0.811)
Cont. AUCs	<b>0.791 (0.786-0.796)</b> <i>0.793 (0.788-0.798)</i>				<b>0.856 (0.852-0.860)</b> <i>0.862 (0.858-0.866)</i>		
Overall AUC	<b>0.825 (0.822-0.828)</b> <i>0.831 (0.827-0.834)</i>						
Cont. Ext. Val.	<b>0.809</b> (0.766-0.852) <i>0.803</i> (0.760-0.846)	<b>0.789</b> (0.737-0.841) <i>0.796</i> (0.745-0.846)	<b>0.830</b> (0.786-0.874) <i>0.826</i> (0.781-0.870)	<b>0.671</b> (0.617-0.724) <i>0.670</i> (0.616-0.724)	<b>0.811</b> (0.776-0.846) <i>0.813</i> (0.778-0.847)	<b>0.924</b> (0.899-0.949) <i>0.922</i> (0.896-0.948)	<b>0.790</b> (0.747-0.834) <i>0.795</i> (0.751-0.838)

Viral Only Cross-validated AUCs



Viral Only: ROC Curve

