

Available online at www.sciencedirect.com

ScienceDirect

journal homepage: www.e-jmii.com

Original Article

Clinical characteristics of hospitalized children with community-acquired pneumonia and respiratory infections: Using machine learning approaches to support pathogen prediction at admission

Tu-Hsuan Chang ^{a,†}, Yun-Chung Liu ^{b,†}, Siang-Rong Lin ^c,
 Pei-Hsin Chiu ^c, Chia-Ching Chou ^{c,**}, Luan-Yin Chang ^{b,*},
 Fei-Pei Lai ^{d,e,f}



^a Department of Pediatrics, Chi Mei Medical Center, Tainan City, Taiwan

^b Department of Pediatrics, National Taiwan University Hospital, College of Medicine, National Taiwan University, Taipei City, Taiwan

^c Institute of Applied Mechanics, National Taiwan University, Taipei City, Taiwan

^d Graduate Institute of Biomedical Electronics and Bioinformatics, Taipei City, National Taiwan University, Taiwan

^e Department of Computer Science and Information Engineering, National Taiwan University, Taipei City, Taiwan

^f Department of Electrical Engineering, National Taiwan University, Taipei City, Taiwan

Received 9 September 2022; received in revised form 3 April 2023; accepted 25 April 2023

Available online 1 May 2023

KEYWORDS

Machine learning;
 Children;
 Respiratory
 infections;
 Pathogens prediction;
 Community-acquired
 pneumonia

Abstract *Background:* Acute respiratory infections (ARIs) are common in children. We developed machine learning models to predict pediatric ARI pathogens at admission.

Methods: We included hospitalized children with respiratory infections between 2010 and 2018. Clinical features were collected within 24 h of admission to construct models. The outcome of interest was the prediction of 6 common respiratory pathogens, including adenovirus, influenza virus types A and B, parainfluenza virus (PIV), respiratory syncytial virus (RSV), and *Mycoplasma pneumoniae* (MP). Model performance was estimated using area under the receiver operating characteristic curve (AUROC). Feature importance was measured using

* Corresponding author. Department of Pediatrics, National Taiwan University Hospital and College of Medicine, National Taiwan University, No. 8, Chung-Shan South Road, Taipei City, 10041, Taiwan.

** Corresponding author. Institute of Applied Mechanics, National Taiwan University, No. 1, Sec. 4, Roosevelt Rd., Da'an Dist., Taipei City, 10617, Taiwan.

E-mail addresses: ccchou@iam.ntu.edu.tw (C.-C. Chou), lychang@ntu.edu.tw (L.-Y. Chang).

† These authors have contributed equally to this work and share first authorship.

<https://doi.org/10.1016/j.jmii.2023.04.011>

1684-1182/ Copyright © 2023, Taiwan Society of Microbiology. Published by Elsevier Taiwan LLC. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

Shapley Additive exPlanation (SHAP) values.

Results: A total of 12,694 admissions were included. Models trained with 9 features (age, event pattern, fever, C-reactive protein, white blood cell count, platelet count, lymphocyte ratio, peak temperature, peak heart rate) achieved the best performance (AUROC: MP 0.87, 95% CI 0.83–0.90; RSV 0.84, 95% CI 0.82–0.86; adenovirus 0.81, 95% CI 0.77–0.84; influenza A 0.77, 95% CI 0.73–0.80; influenza B 0.70, 95% CI 0.65–0.75; PIV 0.73, 95% CI 0.69–0.77). Age was the most important feature to predict MP, RSV and PIV infections. Event patterns were useful for influenza virus prediction, and C-reactive protein had the highest SHAP value for adenovirus infections.

Conclusion: We demonstrate how artificial intelligence can assist clinicians identify potential pathogens associated with pediatric ARIs upon admission. Our models provide explainable results that could help optimize the use of diagnostic testing. Integrating our models into clinical workflows may lead to improved patient outcomes and reduce unnecessary medical costs.

Copyright © 2023, Taiwan Society of Microbiology. Published by Elsevier Taiwan LLC. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

Introduction

Acute respiratory infections (ARIs) are the most common infectious disease during childhood and are also the leading cause of pediatric hospitalization or death worldwide.¹ Pathogens of pediatric ARIs differ from community-acquired pneumonia (CAP) in adults. Infants and toddlers are more susceptible to infections, and there is a high viral pneumonia ratio.^{2,3} Influenza virus (flu), adenovirus, respiratory syncytial virus, and *Mycoplasma pneumoniae* are common etiologies among pediatric ARIs in Taiwan.^{4,5} However, the variability in clinical manifestations makes the diagnosis of pathogens a challenge. Misdiagnosis can lead to inappropriate use of antibiotics in hospitalized children with respiratory illnesses, which causes unnecessary medical care costs and undesirable outcomes for patients.^{6,7}

For pathogen diagnosis, clinicians rely on demographics, clinical manifestations, physical examinations, and recent epidemiological patterns as references. However, clinical presentations associated with pediatric ARIs commonly overlap. No universally established standards exist. In addition, reviewing different sources of health data to determine a diagnosis is time consuming, which adds to the workload of physicians in clinical settings.⁸ Tools to automatically determine the cause of ARIs can reduce the burden on clinicians and enhance overall medical care quality.

Advances in artificial intelligence (AI) have enhanced the performance of computer-aided diagnosis, prognosis prediction, and optimal intervention suggestions. Recent studies have shown that machine learning (ML) methods applied in medical decision support yield impressive results in various diseases and tasks.^{9–11} Previous works on ML for respiratory infections focused mainly on outcome predictions and image processing.^{11–14} In the area of etiology diagnosis, studies on respiratory infections are limited. Lhommet et al. tried to develop a data-driven method to differentiate viral from bacterial pneumonia at patient presentation.¹⁵ Unfortunately, neither experts nor an AI model succeeded in predicting the microbial etiology of ARIs. In another study, Mai and his colleagues used MetaMap to predict viral etiologies among 1685 pediatric admissions. However, the result was less satisfactory. The prediction models showed that the area under the receiver operating

characteristic curve ranged from 0.53 to 0.72 among the 6 common viruses.¹⁶

Timely detection of microorganisms that cause ARIs can help physicians prepare for further disease management. To date, no AI-based etiology prediction models exist for pediatric ARIs that yield clinically applicable performance. In this study, we aimed to develop an ML algorithm to predict pediatric ARI pathogens with available information at admission.

Methods

Study design

Fig. 1 illustrates the schematic diagram of pediatric ARI diagnosis and management in the study framework. The decision-making process usually begins with history taking, physical examinations and orders corresponding to laboratory or radiological studies. After obtaining adequate information, physicians' work usually includes severity assessment, differential diagnosis, pharmacological treatment, and family communication. However, the diagnosis process is dynamic. Reassessment of patients and repeated decision-making are common, especially in complicated cases. The purpose of developing a machine learning tool is to use an established dataset to simplify the work of multitasking.

Cohort and data source

Hospitalized patients admitted to National Taiwan University Hospital between 2010 and 2018 who met the following three criteria were included in the study cohort.

1. Aged between 1 month and 18 years.
2. Presented with acute respiratory infection symptoms (e.g., cough, rhinorrhea) or positive physical examination findings (e.g., crackles, stridor, wheezing) or positive radiological findings (e.g., consolidation, effusion) or tentative diagnosis of respiratory tract infection recorded within 24 h of admission (see [Supplementary Table 1](#) for full list of symptoms, physical examination findings and diagnosis list).

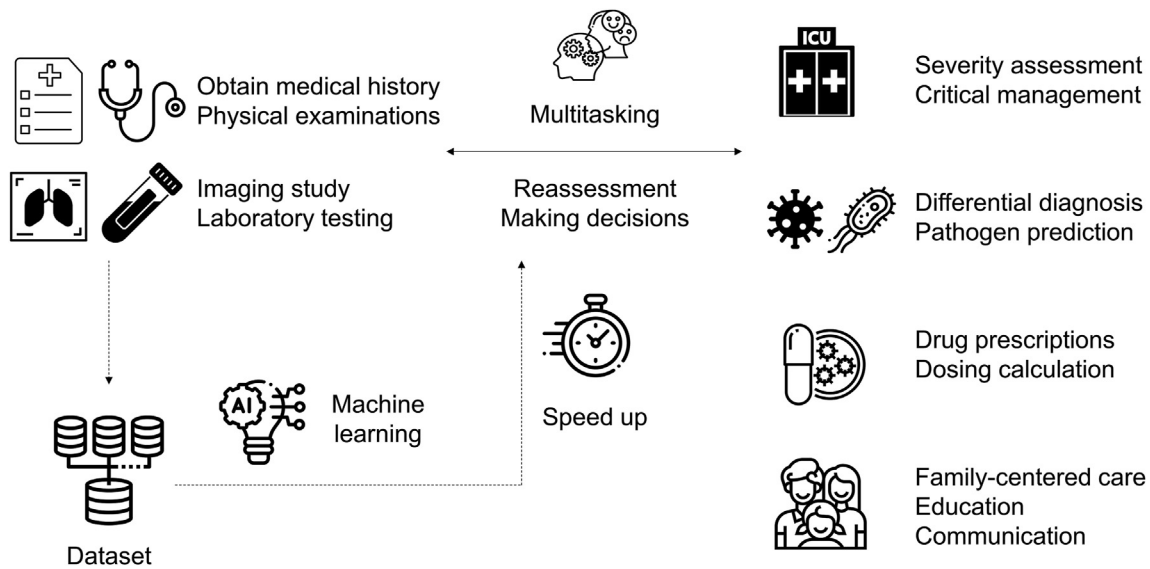


Figure 1. Schematic diagram of pediatric acute respiratory infection diagnosis and management in the study framework.

- Classified with International Classification of Disease, ninth revision (ICD-9) and tenth revision (ICD-10) codes related to respiratory infections at discharge (see [Supplementary Table 2](#) for the full list of included diagnosis codes).

Patients meeting any of the following criteria were excluded from the study.

- Patients discharged from the hospital within 14 days preceding the current admission.
- Patients with congenital pneumonia.
- Patients with immunodeficiency.
- Diagnosis with *Mycobacterium* species infections or tuberculosis.
- Opportunistic infections, such as pulmonary aspergillosis, candidiasis, or pneumocystis pneumonia.
- A diagnosis other than infections that was the likely explanation for the pulmonary infiltrates or respiratory symptoms (e.g., pneumonitis due to inhalation or ingestion of irritant substances, hemothorax).

The electronic health records of the included patients were extracted from the National Taiwan University Hospital integrative medical database and deidentified before analysis. This study was approved by the institutional review board of National Taiwan University Hospital (201912131RINB).

Outcome of interest

Our study aimed to diagnose common respiratory pathogens in pediatric patients in Taiwan, including adenovirus (ADV), influenza A (IAV) and B (IBV) virus, parainfluenza virus (PIV), respiratory syncytial virus (RSV), and *M. pneumoniae* (MP). To focus on community-acquired infections, our study was restricted to patients who were hospitalized with positive pathogen findings within a period of two days after admission and five days prior to admission, to capture the

acute infection phase. We applied viral isolation, antigen testing, and RT-PCR methods to confirm the diagnosis of each pathogen. The diagnostic tests used were determined based on clinical judgment. These tests were performed on samples taken from various sources, including nasopharyngeal swabs, throat swabs, sputum samples, pleural effusion samples, and bronchoalveolar lavage fluid samples. Although *Streptococcus pneumoniae* is an important microorganism in pediatric ARIs, the definition between infections and colonization is not clear.¹⁷ Therefore, we did not include *S. pneumoniae* in our analysis.

Between-group comparison

Clinically relevant manifestations within 24 h of admission were collected for statistical analysis among patients with different pathogen findings. Extreme values of vital signs were excluded in accordance with previous research.¹⁸ All included features are presented in [Supplementary Table 3](#).

The Kolmogorov–Smirnov test was conducted on numerical variables (i.e., laboratory results, vital signs, age) to test for normality. For nonnormally distributed variables, the medians and interquartile ranges were calculated, and the Kruskal–Wallis H test was used for between-group comparisons. For normally distributed variables, the means and standard deviations were reported, and Student's t test was used for between-group comparisons. For categorical variables (e.g., presence of underlying diseases, symptoms, and chest X-ray findings), counts and percentages were calculated, and the chi-square test was used for between-group comparisons. We applied the Benjamini–Hochberg procedure to adjust for multiple comparisons. Adjusted p values < 0.05 were considered significant.

Model training and performance evaluation

Four models were applied for pathogen prediction because of promising results yielded on classification tasks using

clinical data in various studies, namely, logistic regression (LR), random forest (RF), gradient boosting (GB), and extreme gradient boosting (XGB) models.^{9–12} Six independent binary classification models were developed to classify each targeted pathogen. Each model used the same datasets for training (60%), validation (20%), and testing (20%) in chronological order, without any interference or interaction between models.

To evaluate different algorithms against the pathogens, the area under the receiver operating characteristic curve (AUROC) was calculated. The Youden Index with an optimal combination of sensitivity and specificity was used to evaluate the separate performance of each pathogen prediction model. To further analyze the input features, we used SHapley Additive exPlanation (SHAP) values to rank feature importance.

Software

Electronic health records were extracted and preprocessed using the NumPy (version 1.16.5) and Pandas (version 0.25.1) libraries of the Python programming language version 3.7.4 (Python Software Foundation, Fredericksburg, VA, USA). Between-group statistical analyses were conducted using the SciPy package version 1.3.1. While training the models, we used Scikit-learn (the Scikit-learn Contributors, version 0.24.1) for the LR, RF, and GB algorithms. The XGBoost package (version 1.4.2) was used for the XGB algorithm.¹⁹ Evaluation was conducted using the Scikit-learn package, in addition to SHAP values (version 0.39.0). We used TreeExplainer, built based on SHAP values, to interpret and explain tree-based models in our study.²⁰

Results

Cohort selection

We depict the workflow of our cohort selection process in Fig. 2. There were 14,201 patient admissions that fulfilled the inclusion criteria in the dataset. A total of 1507 admissions that met the exclusion criteria were excluded, resulting in 12694 admissions in the study cohort. The study cohort was further separated into training, validation, and testing datasets of 7624, 2536 and 2534 admissions, respectively. We present the characteristics of the patients in each dataset in Table 1.

Characteristics of patients with each pathogen

Information regarding the demographics and clinical features of the included children with common pediatric respiratory pathogens is summarized in Table 2. Laboratory results and radiological findings can be found in Table 3. We provide detailed data in Supplementary Table 4. Among the 6 pathogens included in our models, patients infected with MP were older (median age 5.5 years, IQR: 3.5–8.1), and those infected with RSV were younger (0.8 years old, IQR: 0.3–1.5). Approximately 98.4% of patients infected with MP were admitted due to pneumonia, followed by RSV (75.8%) and PIV (70.8%). For IAV and IBV infections, patients were

classified as having upper respiratory tract infections (URI) (34.6% and 32.8%). Patients with RSV infections tended to stay in the hospital longer (5 days vs. 4 days) than patients with other pathogens. Approximately 17.3% of IAV patients received intensive care, while only 5% of MP patients were transferred to the intensive care unit (ICU).

Among all pathogens included in other models, each of them had unique characteristics. Nearly all patients infected with MP had cough (95.9%), but relatively few children reported rhinorrhea (31.3%) in our cohort. Pneumonia patches/consolidation (59.6%) and pleural effusion (6.9%) were more common in MP patients than in patients infected with other microorganisms. Patients with RSV infections had more prominent rhinorrhea (58.4%) and dyspnea (58.6%). Rhonchi (43.9%) and wheezing (50.7%) were the two most common physical examination findings in patients with RSV infections than other pathogens. Laboratory data usually revealed lymphocyte predominance and low CRP levels (0.4 mg/dL, IQR: 0.1–1.3) associated with RSV compared with other viruses. A higher proportion of children with ADV infections presented with fever (96.0%) and sore throat (27.3%), with 79.7% having a body temperature above the normal range within 24 h of admission. Laboratory results revealed higher mean WBC (12.0 k/ μ L, IQR: 8.9–16.0) and CRP (4.1 mg/dL, IQR: 1.9–7.1) values associated with ADV compared with other pathogens. Regarding IAV and IBV infections, physical examinations showed less evidence of pulmonary involvement in hospitalized children. The rates of crackles (23.3%), rhonchi (18.4%), and wheezing (13.4%) in patients infected with IAV were the lowest compared to other pathogens. However, 2.3% of IAV patients presented with cyanosis, which echoed the finding that 17.3% of patients needed intensive care. Higher proportion of patients with cardiovascular (9.6%) or respiratory (4.8%) diseases had positive results on PIV exams. Hoarseness (15.7%) and stridor (16.0%) were specific and commonly reported clinical manifestations in PIV infections.

Diagnostic performance results

On the testing set, XGB models trained with all 79 features achieved the best performance in diagnosing most of the pediatric respiratory tract infection-associated pathogens within 24 h after admission (MP: AUROC 0.87, 95% CI 0.83–0.90; RSV: AUROC 0.85, 95% CI 0.85–0.89; ADV: AUROC 0.84, 95% CI 0.75–0.84; IAV: AUROC 0.76, 95% CI 0.76–0.83; IBV: AUROC 0.73, 95% CI 0.67–0.77; PIV: AUROC 0.75, 95% CI 0.68–0.78). For both clinical interpretability and further feature selection, all features were ranked by importance by SHAP values. The top 9 features were age, fever, C-reactive protein (CRP), white blood cell (WBC) count, platelet count, lymphocyte ratio, peak temperature, peak heart rate and number of cases for each pathogen identified in our institution over the past 30 days (event pattern). With the 9 selected features, XGB models achieved better performance compared with other algorithms. Additionally, the AUROC was only slightly reduced (ranged from 0.70 to 0.87). The negative predictive values (NPVs) for each pathogen were remarkable (MP: NPV 0.99, 95% CI 0.99–0.99; RSV: NPV 0.95, 95% CI 0.95–0.95; ADV: NPV 0.98, 95% CI 0.98–0.98; IAV: NPV 0.97, 95% CI

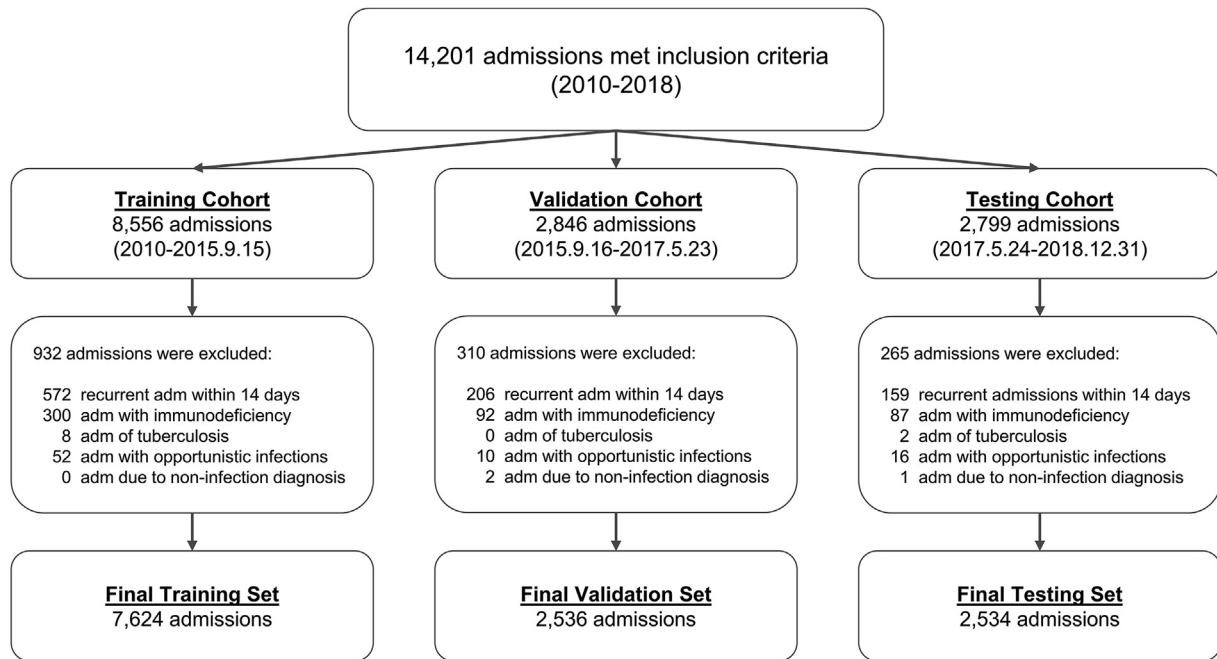


Figure 2. Workflow of cohort selection process (adm: admission).

0.97–0.97; IBV: NPV 0.98, 95% CI 0.97–0.98; PIV: NPV 0.98, 95% CI 0.98–0.98). The detailed performance matrix of each pathogen is shown in [Supplementary Table 5](#).

Fig. 3 shows the model performances with different sources of information. The basic information included demographics (DE), symptoms and physical examinations (PE). The AUROC of the MP predictive models trained with the basic information was 0.82. The addition of laboratory data (LD) increased the AUROC by 0.04. Further addition of vital sign (VS) data did not obviously increase the AUROC. Most of the models achieved their best AUROC after the addition of event pattern (EP, number of cases for each pathogen identified in our institution over the past 30 days), except the model for ADV prediction.

Explaining the rationale behind the predicted models

Fig. 4 presents the features by relative importance from the XGB algorithm based on SHAP values. The event pattern ranked in the top three for all pathogens, which had a positive influence on all model performances. Age was the most important feature in the MP, RSV, and PIV models. However, age influenced the prediction models differently for different pathogens. The XGB model tended to predict that older patients would have positive MP findings. In contrast, the models tended to predict that younger patients would have positive findings of RSV and PIV. Differential influences on prediction can be observed for WBC

Table 1 Characteristics of patients in the training, validation, and testing datasets.

	Training dataset (N = 7624)	Validation dataset (N = 2536)	Testing dataset (N = 2534)
Age	2.3 (1.1–4.6)	2.0 (0.9–4.0)	2.1 (1.0–3.9)
Sex			
Male	4337 (56.9%)	1411 (55.6%)	1485 (58.6%)
Female	3287 (43.1%)	1125 (44.4%)	1049 (41.4%)
Any chronic disease	2496 (32.7%)	865 (34.1%)	825 (32.6%)
ICU	1069 (14.0)	301 (11.9)	281 (11.1)
Death	4 (0.1)	1 (0.0)	0 (0.0)
Identified pathogens			
Respiratory syncytial virus	842 (11.0)	311 (12.3)	351 (13.9)
Adenovirus	433 (5.7)	84 (3.3)	131 (5.2)
Influenza A virus	227 (3.0)	143 (5.6)	145 (5.7)
Influenza B virus	205 (2.7)	87 (3.4)	104 (4.1)
Parainfluenza virus	224 (2.9)	70 (2.8)	100 (3.9)
<i>Mycoplasma pneumoniae</i>	122 (1.6)	129 (5.1)	68 (2.7)

Abbreviation: ICU, intensive care unit.

Table 2 Demographics and clinical features of patients with common pediatric respiratory pathogens in the dataset.

	<i>Mycoplasma pneumoniae</i> (N = 319)	RSV (N = 1504)	Adenovirus (N = 648)	Influenza A virus (N = 515)	Influenza B virus (N = 396)	Parainfluenza virus (N = 394)	P value
Demographics							
Age	5.5 (3.5–8.1)	0.8 (0.3–1.5)	3.3 (1.8–4.8)	2.5 (0.7–4.9)	3.1 (1.4–6.0)	1.2 (0.6–1.8)	<0.001
Sex (male %)	163 (51.1)	931 (61.9)	379 (58.5)	281 (54.6)	202 (51.0)	240 (60.9)	<0.001
No chronic disease	275 (86.2)	1218 (81.0)	498 (76.9)	317 (61.6)	268 (67.7)	295 (74.9)	<0.001
Event pattern (Accumulated cases in the past 30 days)	5.0 (2.5–8.0)	17.0 (11.0–24.0)	7.0 (4.0–14.0)	7.0 (4.0–15.0)	7.0 (3.0–11.0)	4.0 (2.0–7.0)	<0.001
Disease severity							
Pneumonia	314 (98.4)	1140 (75.8)	427 (65.9)	315 (61.2)	250 (63.1)	279 (70.8)	<0.001
LRI	3 (0.9)	299 (19.9)	26 (4.0)	22 (4.3)	16 (4.0)	58 (14.7)	
URI	2 (0.6)	65 (4.3)	195 (30.1)	178 (34.6)	130 (32.8)	57 (14.5)	
Outcome							
Length of stay (days)	4.0 (3.0–5.0)	5.0 (4.0–7.0)	4.0 (3.0–5.0)	4.0 (3.0–5.0)	4.0 (3.0–5.0)	4.5 (3.0–6.0)	<0.001
Intensive care	16 (5.0)	200 (13.3)	34 (5.2)	89 (17.3)	40 (10.1)	50 (12.7)	<0.001
Death	0 (0.0)	0 (0.0)	2 (0.3)	1 (0.2)	2 (0.5)	0 (0.0)	0.13
Vital signs^a							
Peak temperature (Above normal)	241 (75.5)	668 (46.2)	507 (79.7)	367 (73.7)	285 (72.7)	170 (44.7)	<0.001
Peak heart rate (Above normal)	275 (86.2)	935 (64.7)	491 (77.2)	406 (81.5)	290 (74.0)	259 (68.2)	<0.001
Peak respiratory rate (Above normal)	304 (95.3)	1124 (77.7)	497 (78.1)	361 (72.5)	313 (79.8)	280 (73.7)	<0.001
Lowest SpO ₂ (Below normal)	42 (13.5)	164 (11.6)	34 (6.0)	40 (8.7)	39 (10.5)	31 (8.3)	0.001
Clinical features							
Fever	303 (95.0)	1089 (72.4)	622 (96.0)	490 (95.1)	370 (93.4)	328 (83.2)	<0.001
Cough	306 (95.9)	1437 (95.5)	517 (79.8)	408 (79.2)	315 (79.5)	353 (89.6)	<0.001
Sputum	228 (71.5)	1026 (68.2)	326 (50.3)	243 (47.2)	194 (49.0)	224 (56.9)	<0.001
Rhinorrhea	100 (31.3)	879 (58.4)	353 (54.5)	234 (45.4)	185 (46.7)	207 (52.5)	<0.001
Dyspnea	80 (25.1)	881 (58.6)	126 (19.4)	106 (20.6)	108 (27.3)	185 (47.0)	<0.001
Sore throat	34 (10.7)	55 (3.7)	177 (27.3)	52 (10.1)	52 (13.1)	21 (5.3)	<0.001
Hoarseness	5 (1.6)	60 (4.0)	17 (2.6)	20 (3.9)	16 (4.0)	62 (15.7)	<0.001
Skin rash	11 (3.4)	66 (4.4)	31 (4.8)	17 (3.3)	25 (6.3)	17 (4.3)	0.33
Physical examination							
Rales/Crackles	193 (60.5)	615 (40.9)	168 (25.9)	120 (23.3)	116 (29.3)	124 (31.5)	<0.001
Rhonchi	97 (30.4)	661 (43.9)	166 (25.6)	95 (18.4)	78 (19.7)	148 (37.6)	<0.001
Wheezing	62 (19.4)	763 (50.7)	107 (16.5)	69 (13.4)	56 (14.1)	126 (32.0)	<0.001
Stridor	0 (0)	32 (2.1)	8 (1.2)	17 (3.3)	12 (3.0)	63 (16.0)	<0.001

Continuous variables are described as medians (IQRs) and were tested using the Kruskal–Wallis H test. Categorical variables are presented as numbers (%) and were tested using the chi-square test.^a

Abbreviations: LRI, lower respiratory tract infection; RSV, respiratory syncytial virus; URI, upper respiratory tract infection.

^a Rates of abnormal vital signs were calculated using the following formula: No. of admissions with abnormal vital signs/ No. of admissions with vital sign records

count as well. WBC count was positively correlated with the prediction of ADV but negatively correlated with the predictions on other pathogens.

Discussion

Our study proposed ML models that incorporated both subjective and objective measures collected at admission

by using a large dataset. The variables included in the model were obtained from medical histories, physical examinations, vital signs, laboratory results and radiographic findings in the order resembling a clinical approach. Our approach was novel in that we used routinely collected information to predict the presence of pathogens causing ARIs in children. Pathogen prediction results could thus be automatically produced after completing each step. The prediction performance increased as more information was

Table 3 Laboratory results and radiological findings of common pediatric respiratory pathogens in the dataset.

	<i>Mycoplasma pneumoniae</i> (N = 319)	RSV (N = 1504)	Adenovirus (N = 648)	Influenza A virus (N = 515)	Influenza B virus (N = 396)	Parainfluenza virus (N = 394)	P value
Laboratory result							
WBCs (k/ μ L)	8.1 (6.2–10.3)	9.2 (7.2–11.9)	12.0 (8.9–16.0)	8.5 (6.0–12.1)	8.6 (5.7–13.2)	10.2 (7.6–13.7)	<0.001
Hemoglobin (g/dL)	12.5 (11.9–13.3)	12.1 (11.3–12.8)	12.0 (11.3–12.7)	12.1 (11.2–13.1)	12.5 (11.6–13.3)	12.2 (11.4–12.9)	<0.001
Platelets (k/ μ L)	244.0 (201.0–306.0)	309.0 (236.2–393.0)	259.0 (208.0–316.0)	245.0 (186.0–320.0)	227.0 (174.0–297.8)	278.0 (221.0–362.0)	<0.001
Neutrophils (%)	64.8 (55.1–72.6)	34.6 (21.6–49.1)	63.0 (50.5–72.2)	56.1 (41.0–72.1)	59.0 (42.0–71.6)	42.1 (29.0–58.0)	<0.001
Lymphocytes (%)	25.1 (17.2–34.0)	52.0 (37.9–64.8)	25.0 (16.0–35.1)	29.7 (16.6–42.6)	29.0 (17.3–43.4)	43.9 (28.8–56.6)	<0.001
CRP (mg/dL)	2.4 (1.2–4.5)	0.4 (0.1–1.3)	4.1 (1.9–7.1)	1.0 (0.3–3.3)	1.3 (0.4–3.6)	0.8 (0.2–2.3)	<0.001
Radiological findings							
Infiltration/Haziness	106 (33.2)	645 (42.9)	210 (32.4)	166 (32.2)	141 (35.6)	161 (40.9)	<0.001
Patches/Opacity/ Consolidation	190 (59.6)	231 (15.4)	101 (15.6)	95 (18.4)	61 (15.4)	65 (16.5)	<0.001
Hyperinflation	2 (0.6)	80 (5.3)	14 (2.2)	11 (2.1)	8 (2.0)	10 (2.5)	<0.001
Pleural effusion	22 (6.9)	3 (0.2)	4 (0.6)	6 (1.2)	4 (1.0)	3 (0.8)	<0.001

Continuous variables are described as medians (IQRs) and were tested using the Kruskal–Wallis H test. Categorical variables are presented as numbers (%) and were tested using the chi-square test.

Abbreviations: WBCs, White blood cells; CRP, C-reactive protein.

obtained. Our approach has important implications for clinical practice, as it enables earlier and more accurate diagnosis of ARIs in children, which can shed light on targeted treatments and prevent unnecessary antibiotic use. Our approach also has the potential to reduce unnecessary diagnostic tests and medical costs.

Our models outperformed similar attempts reported in the literature and can be applied in clinical settings to facilitate decision-making and reduce physician workload. In a previous study, Chen et al. used a five-factor model (age, duration of fever, erythrocyte sedimentation rate, leukocyte count, and neutrophil proportion) to predict MP for hospitalized children with respiratory symptoms, achieving an AUROC of 0.75.²¹ Our model had a better performance on MP with an AUROC of 0.87 when different features were included in model development. Mai and his colleagues made efforts to predict common respiratory viruses in the US by combining natural language processing (NLP) tools and the ML approach.¹⁶ However, the performance of the models for 4 overlapping viruses was inferior to that of our models (AUROC: ADV, 0.53 vs. 0.81; influenza, 0.71 vs. 0.77; PIV, 0.69 vs. 0.73; RSV, 0.72 vs. 0.84). Large case numbers and independent validation datasets are the main advantages compared with previous models. The ensemble methods applied in our study outperformed models applied in previous studies (i.e., logistic regression, decision tree) due to the nature of combining multiple decision trees to minimize prediction error, which fits well with categorization tasks using clinical manifestations where the decision boundary is often nonlinear.¹⁹

The positive predictive value (PPV) of the machine learning model we trained for each pathogen, except RSV (PPV 0.34), ranged from 0.06 to 0.13, which may seem less satisfactory. However, the positive rates of pathogen exams in our dataset ranged from 0.04 to 0.29 (data not shown). In comparison, our model demonstrated a PPV ranging from 0.06 to 0.34, indicating that it can still outperform the conventional decision-making process. Clinicians usually overestimate the pretest probability compared to the available scientific literature.²² In our ML algorithm, we reached outstanding negative predictive values for the 6 pathogens (0.95–0.99). These models could be applied to optimize the use of diagnostic testing in clinical practice, avoid unnecessary exams, and achieve diagnostic stewardship. In our study, the overall performance of the IAV, IBV and PIV models (AUROC: 0.70–0.77) is less impressive than that of the ADV, RSV and MP models (AUROC: 0.81–0.87). Although different pathogens associated with ARIs vary in their clinical manifestations, sometimes the differences are subtle. It is difficult to distinguish between flu and PIV solely by clinical features and physical examination findings.²³ These results highlight the added value of vital signs in our models, which increased the overall diagnostic accuracy, especially in the flu and PIV models.

During model development, we successfully reduced the input variables from 79 to 9 without a significant loss of performance. SHAP values effectively explained the effect of each feature and made the result applicable in clinical utilization. To identify the 6 selected pathogens, age was an important feature that was present in all models. Children with MP infections tended to be older than those with RSV or PIV infections. However, in the IAV, IBV, and ADV

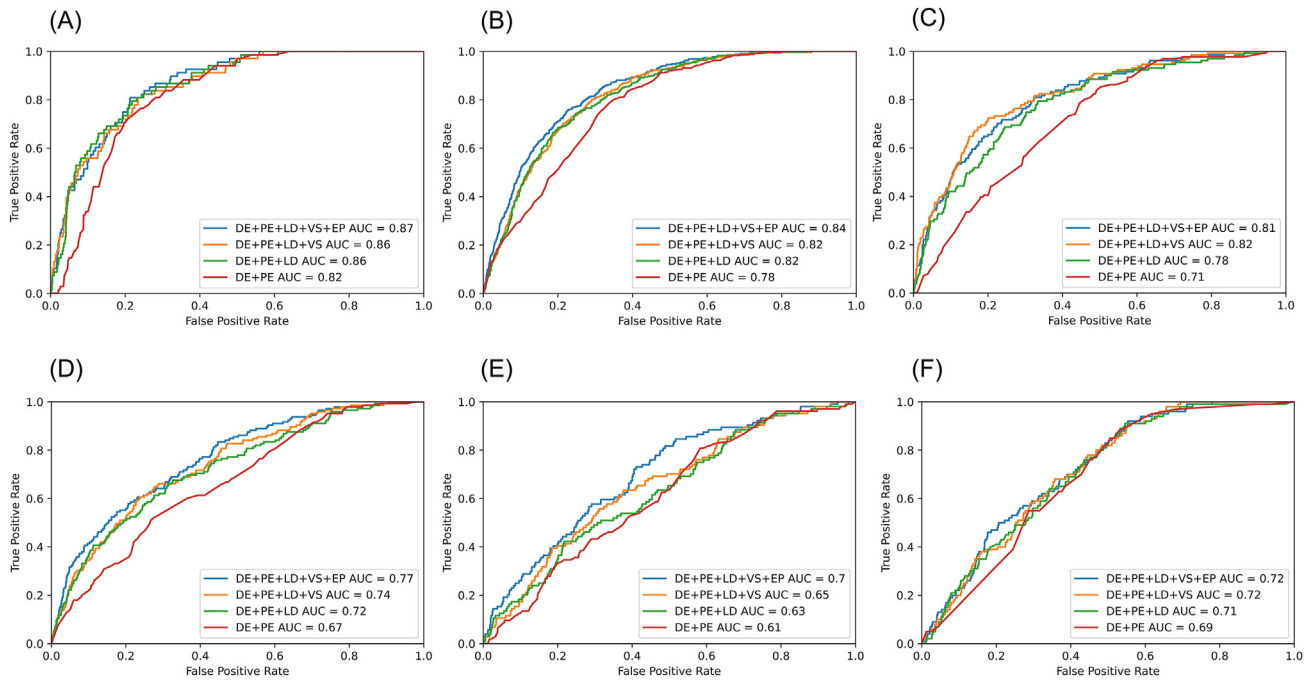


Figure 3. Performance of prediction models when inputting different information. DE: Demographics, PE: Physical examination, LD: laboratory data, VS: vital sign, EP: event pattern (accumulated cases in the past month). (A) *Mycoplasma pneumoniae* (B) Respiratory syncytial virus (C) Adenovirus (D) Influenza A virus (E) Influenza B virus (F) Parainfluenza virus.

models, the pathogens which cause epidemics in all age groups, the importance of age was lower than other pathogens in the models. Our results were compatible with previous epidemiological data.^{3,5} To better diagnose

respiratory virus infections, epidemiological patterns are very important. We lacked epidemiological trends, clusters, or contact history in the extracted features. Instead, we used the number of cases for each pathogen identified

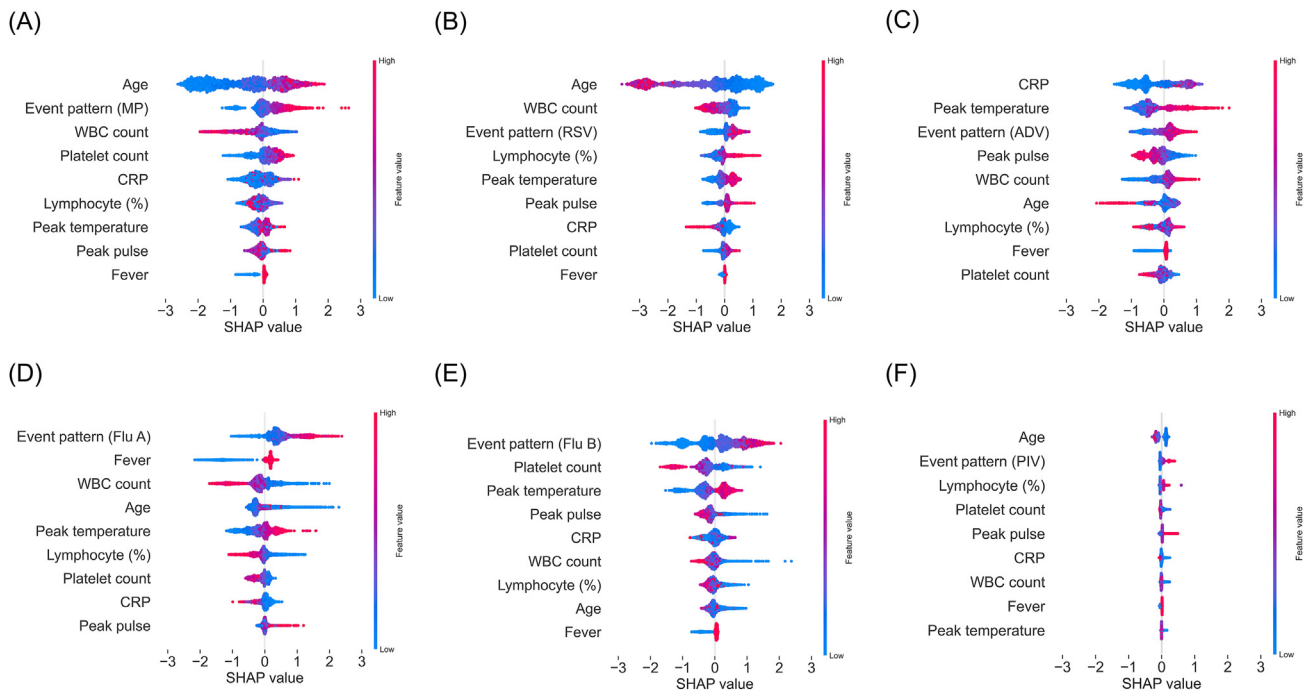


Figure 4. Feature importance analysis by SHapley Additive exPlanation (SHAP) values. All features by importance were presented. Each point represents a patient’s feature value (blue represents a low-value feature, and red stands for high-value features). A negative impact (left of the vertical line) of a feature value means prediction of the specific pathogen is less likely. In contrast, a positive impact (right of the vertical line) represents a prediction that is more likely to be the specific pathogen. (A) *Mycoplasma pneumoniae* (B) Respiratory syncytial virus (C) Adenovirus (D) Influenza A virus (E) Influenza B virus (F) Parainfluenza virus.

over the past 30 days, and this parameter significantly increased the overall performance. In the future, by incorporating data from multiple sites collectively or national real-time surveillance systems, the accuracy and sensitivity of models could be further optimized. A similar approach has been used in influenza forecasting.²⁴ ADV is unique, as it is best predicted by high fever and elevated CRP and WBC levels. The symptoms and laboratory results of adenovirus infection mimic bacterial infections, commonly associated with misuse of antibiotics.²⁵ The performance of our model in differentiating ADV from other ARIs was quite good.

Our study has some limitations. First, the models were trained on data from a single medical center, which could compromise the robustness of the models. Not all patients had been tested for all pathogens included in the dataset. This could lead to an underestimation of the true prevalence of these pathogens in our study population. Additionally, the decision to test for a particular pathogen was based on clinical judgment, which could bring some selection bias into our results. While we used multiple methods to confirm the diagnosis of each pathogen, there is always a possibility of false negative results. Therefore, the true prevalence of these pathogens may be higher than what we were able to detect in our study. Prospective validation with data from multiple centers may be helpful in improving generalizability. Second, the 6 pathogens represented only 30% of ARI etiologies in our dataset. To predict one pathogen with a relatively low prevalence, a low positive predictive value is inevitable. Third, vaccination or contact/cluster history are important information during history taking. Related items were not added to the current models. The power of prediction will be enhanced by including these features in future studies. Fourth, some common pathogens in pediatric respiratory tract infections could not be included in our analysis because of a lack of routine surveillance, such as human metapneumovirus and rhinovirus.²⁶ Additionally, the issue of bacterial coinfection is crucial in accurately predicting pathogens and prescribing appropriate antibiotics. However, there is no clear and widely accepted consensus for differentiating between coinfection and colonization. Therefore, we did not conduct further analysis on this issue, which could possibly affect the presentation of the diseases. Further investigation is needed to better understand this topic. Prospective studies with the introduction of multiplex PCR or syndromic point-of-care testing are needed to support model performance. A comparison of AI-aided decision-making and the standard of care in the real world is also necessary.

Conclusion

Our study demonstrates how artificial intelligence can assist clinicians in identifying potential pathogens in pediatric respiratory infections at admission. The results were explainable and applicable to clinical practice. Based on our model's excellent negative predictive value, clinicians can potentially avoid unnecessary medical costs and diagnostic tests, while still ensuring accurate diagnosis. Our approach has the potential to decrease the workload of clinicians and improve the overall quality of medical care.

Therefore, we recommend integrating our model into clinical workflows and using it in conjunction with clinical judgement.

Availability of data and materials

For researchers who want to access these datasets, please direct your request to the corresponding authors. Individual participant data that underlie the results reported in this article, after deidentification (text, tables, figures, and appendices), will be shared by the corresponding author upon reasonable request for academic and research purposes.

Funding

This research was funded by grants from the Ministry of Science and Technology, Taiwan (grant numbers MOST 110-2634-F-002-032, MOST 109-2314-B-002-238, MOST 110-2314-B-002-249, and MOST 111-2314-B-002-146) and Chi Mei Medical Center (CMOR11202, CMNCKU11101, CMFHR11181, and CCFHR11202). The funders had no role in the study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Acknowledgements

The authors would like to express their gratitude to the staff in the Department of Medical Research at National Taiwan University Hospital for their generous support in the extraction of the electronic health records from the National Taiwan University Hospital integrative medical database.

References

1. Williams BG, Gouws E, Boschi-Pinto C, Bryce J, Dye C. Estimates of world-wide distribution of child deaths from acute respiratory infections. *Lancet Infect Dis* 2002;2(1):25–32.
2. Lin WH, Chiu HC, Chen KF, Tsao KC, Chen YY, Li TH, et al. Molecular detection of respiratory pathogens in community-acquired pneumonia involving adults. *J Microbiol Immunol Infect* 2022;55(5):829–37.
3. Jain S, Williams DJ, Arnold SR, Ampofo K, Bramley AM, Reed C, et al. Community-acquired pneumonia requiring hospitalization among US children. *N Engl J Med* 2015;372(9):835–45.
4. Hsu HT, Huang FL, Ting PJ, Chang CC, Chen PY. The epidemiological features of pediatric viral respiratory infection during the COVID-19 pandemic in Taiwan. *J Microbiol Immunol Infect* 2022;55(6 Pt 1):1101–7.
5. Chi H, Huang YC, Liu CC, Chang KY, Huang YC, Lin HC, et al. Characteristics and etiology of hospitalized pediatric community-acquired pneumonia in Taiwan. *J Formos Med Assoc* 2020;119(10):1490–9.
6. Handy LK, Bryan M, Gerber JS, Zaoutis T, Feemster KA. Variability in antibiotic prescribing for community-acquired pneumonia. *Pediatrics* 2017;139(4).
7. Fitzpatrick T, Malcolm W, McMenamin J, Reynolds A, Guttman A, Hardelid P. Community-based antibiotic prescribing attributable to respiratory syncytial virus and other common respiratory viruses in young children: a population-

- based time-series study of scottish children. *Clin Infect Dis* 2021;**72**(12):2144–53.
8. Sinsky C, Colligan L, Li L, Prgommet M, Reynolds S, Goeders L, et al. Allocation of physician time in ambulatory practice: a time and motion study in 4 specialties. *Ann Intern Med* 2016; **165**(11):753–60.
 9. Goto T, Camargo CA, Faridi MK, Freishtat RJ, Hasegawa K. Machine learning–based prediction of clinical outcomes for children during emergency department triage. *JAMA Netw Open* 2019;**2**(1):e186937–.
 10. Liu YC, Cheng HY, Chang TH, Ho TW, Liu TC, Yen TY, et al. Evaluation of the need for intensive care in children with pneumonia: machine learning approach. *JMIR Med Inform* 2022;**10**(1):e28934.
 11. Caruana R, Lou Y, Gehrke J, Koch P, Sturm M, Elhadad N. Intelligible models for healthcare: predicting pneumonia risk and hospital 30-day readmission. In: *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining; 2015*; 2015. p. 1721–30.
 12. Hu CA, Chen CM, Fang YC, Liang SJ, Wang HC, Fang WF, et al. Using a machine learning approach to predict mortality in critically ill influenza patients: a cross-sectional retrospective multicentre study in Taiwan. *BMJ Open* 2020;**10**(2):e033898.
 13. Jaeger S, Karargyris A, Candemir S, Folio L, Siegelman J, Callaghan F, et al. Automatic tuberculosis screening using chest radiographs. *IEEE Trans Med Imag* 2013;**33**(2):233–45.
 14. Chowdhury ME, Rahman T, Khandakar A, Mazhar R, Kadir MA, Mahbub ZB, et al. Can AI help in screening viral and COVID-19 pneumonia? *IEEE Access* 2020;**8**:132665–76.
 15. Lhommet C, Garot D, Grammatico-Guillon L, Jourdainaud C, Asfar P, Faisy C, et al. Predicting the microbial cause of community-acquired pneumonia: can physicians or a data-driven method differentiate viral from bacterial pneumonia at patient presentation? *BMC Pulm Med* 2020;**20**(1):1–9.
 16. Mai MV, Krauthammer M. Controlling testing volume for respiratory viruses using machine learning and text mining. In: *AMIA annual symposium proceedings; 2016*. American Medical Informatics Association; 2016. p. 1910.
 17. Weiser JN, Ferreira DM, Paton JC. *Streptococcus pneumoniae*: transmission, colonization and invasion. *Nat Rev Microbiol* 2018;**16**(6):355–67.
 18. Kim SY, Kim S, Cho J, Kim YS, Sol IS, Sung Y, et al. A deep learning model for real-time mortality prediction in critically ill children. *Crit Care* 2019;**23**(1):1–10.
 19. Chen T, Guestrin C. Xgboost: a scalable tree boosting system. In: *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining; 2016*; 2016. p. 785–94.
 20. Lundberg SM, Erion G, Chen H, DeGrave A, Prutkin JM, Nair B, et al. From local explanations to global understanding with explainable AI for trees. *Nat Mach Intell* 2020;**2**(1):56–67.
 21. Chen J, Yin Y, Zhao L, Zhang L, Zhang J, Yuan S. *Mycoplasma pneumoniae* infection prediction model for hospitalized community-acquired pneumonia children. *Pediatr Pulmonol* 2021;**56**(12):4020–8.
 22. Morgan DJ, Pineles L, Owczarzak J, Magder L, Scherer L, Brown JP, et al. Accuracy of practitioner estimates of probability of diagnosis before and after testing. *JAMA Intern Med* 2021;**181**(6):747–55.
 23. Choi EH, Lee HJ, Kim SJ, Eun BW, Kim NH, Lee JA, et al. The association of newly identified respiratory viruses with lower respiratory tract infections in Korean children, 2000–2005. *Clin Infect Dis* 2006;**3**(5):585–92.
 24. Yang CY, Chen RJ, Chou WL, Lee YJ, Lo YS. An integrated influenza surveillance framework based on national influenza-like illness incidence and multiple hospital electronic medical records for early prediction of influenza epidemics: design and evaluation. *J Med Internet Res* 2019;**21**(2):e12341.
 25. Tabain I, Ljubin-Sternak S, Cepin-Bogovic J, Markovinovic L, Knezovic I, Mlinaric-Galinovic G. Adenovirus respiratory infections in hospitalized children: clinical findings in relation to species and serotypes. *Pediatr Infect Dis J* 2012;**31**(7):680–4.
 26. Sim JY, Chen YC, Hsu WY, Chen WY, Chou Y, Chow JC, et al. Circulating pediatric respiratory pathogens in Taiwan during 2020: dynamic change under low COVID-19 incidence. *J Microbiol Immunol Infect* 2022;**55**(6 Pt 2):1151–8.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.jmii.2023.04.011>.