

1 **Development and Validation of a Deep Learning System for the**  
2 **Diagnosis of Pediatric Diseases: A Large-Scale Real-World Data**  
3 **Study in Shanghai**

4

5 Xiaoling Ge<sup>1¶</sup>, Yi Wang<sup>2¶</sup>, Li Xie<sup>3</sup>, Yujuan Shang<sup>1</sup>, Yihui Zhai<sup>4</sup>, Zhiheng Huang<sup>5</sup>, Jianfeng

6 Huang<sup>6</sup>, Chengjie Ye<sup>7</sup>, Ao Ma<sup>1</sup>, Wanting Li<sup>7</sup>, Xiaobo Zhang<sup>6\*</sup>, Hong Xu<sup>4\*</sup>

7

8 **Affiliations**

9 <sup>1</sup>Data Management Center, Children's Hospital of Fudan University, Shanghai, China,

10 <sup>2</sup> Ministry of Education Key Laboratory of Contemporary Anthropology, Department of  
11 Anthropology and Human Genetics, School of Life Sciences, Fudan University; Human  
12 Phenome Institutes, Fudan University; International Human Phenome Institutes (Shanghai),  
13 Shanghai, China,

14 <sup>3</sup> Clinical Research Institute, Shanghai Jiao Tong University School of Medicine, Shanghai,  
15 China,

16 <sup>4</sup>Department of Nephrology, Children's Hospital of Fudan University, Shanghai, China,

17    5 Department of Gastroenterology, Children’s Hospital of Fudan University, Shanghai, China,

18    <sup>6</sup>Department of Respiratory medicine, Children’s Hospital of Fudan University, Shanghai, China,

19    <sup>7</sup>Information Technology Center, Children’s Hospital of Fudan University, Shanghai, China

20

21    \*Correspondence author

22    Email: [hxu@shmu.edu.cn](mailto:hxu@shmu.edu.cn) (H Xu); [zhangxiaobo0307@163.com](mailto:zhangxiaobo0307@163.com) (X Zhang)

23    ¶These authors contributed equally to this work.

## 24 **Abstract**

25 **Background** Artificial intelligence (AI)-assisted diagnosis is considered to be the future  
26 direction of improving the efficiency and accuracy of pediatric diseases diagnosis, while the  
27 existing research based on AI are far from sufficient because of limited data amount, inadequate  
28 coverage of disease types, or high construction costs, and have not been applied on a large scale.  
29 We aimed to develop an accurate deep learning model trained on millions of real-world data to  
30 verify the feasibility of the technology, and build the whole process of outpatient auxiliary  
31 diagnosis.

32 **Methods and findings** We applied a Chinese Natural Language Processing (NLP) and an end-  
33 to-end deep neural network classifier to the outpatient's electronic medical records (EMRs) in a  
34 single child care center in Shanghai, China, to unstructured text processing and construct an  
35 auxiliary diagnostic model, all patients were aged from 0 to 18 years. A training cohort with  
36 millions of records and an independent validation cohort with tens of thousands of records were  
37 intake separately and calculate diagnosis concordance rate (DCR) of model in each diseases  
38 group. The records with inconsistent diagnoses between human and AI were evaluated by  
39 clinical experts' group, and calculate the relative correct rate (RCR) to evaluate the diagnostic  
40 performance of the model. A total of 5,271,347 medical records were intake in model training

41 covering sixteen categories of diseases according to disease coding, reaching a DCR of 95.49%  
42 (95.48~95.51). For validation, 91,880 records were obtained from validation dataset, which  
43 reached a DCR of 93.51% (93.35~93.67) and FDCR of 72.04% (71.75~72.33). It was confirmed  
44 that the accuracy of the model was still higher than that of human with most RCR>1 in  
45 validation dataset.

46 **Conclusions** The deep learning system could support diagnosis of pediatric diseases, which has  
47 high diagnostic performance, comprehensive disease coverage, feasible technology, and can be  
48 promoted in multiple sites in the future.

49 **Funding** The Authors received no specific funding for this work.

## 50 **Introduction**

51 Diagnoses that are missed, incorrect or delayed, which are often less mentioned than problems  
52 such as drug errors and operational mistakes, are believed to affect 10%-20% of cases and result  
53 in serious consequences. The expression ability of pediatric patients is poor, and the clinical  
54 manifestations of diseases are atypical. In addition, physicians specializing in adult medicine are  
55 frequently requested to play the role of a pediatrician in grassroots hospitals, even though they  
56 tend to lack of knowledge and experience in pediatrics, and have subjective assumptions and  
57 narrow ideas in clinical thinking.<sup>1</sup> Unfortunately, most clinicians tend to avoid discussing  
58 diagnostic errors publicly, which makes it difficult to detect and correct them promptly.<sup>2</sup>

59 As we learn more about pediatric diseases, we learn more about the power of artificial  
60 intelligence (AI) tools that can be used in unprecedented ways. AI-assisted diagnosis is  
61 considered to be an important means to summarize complex diagnostic mechanisms, improve the  
62 efficiency and accuracy of pediatric clinical diagnostics,<sup>3</sup> and provide a solution to the imbalance  
63 between the supply and demand of pediatricians and patients sequentially. The deep learning  
64 method, as a subgroup of AI, is a particularly promising method that automatically learns entity  
65 representations from natural language and has been shown to match and even outperform human  
66 performance in task-specific applications. Although it requires large datasets for training, deep

67 learning has demonstrated relative robustness to noise in ground truth labels, among others. The  
68 automated capabilities of AI offer the potential to enhance the qualitative expertise of clinicians,  
69 including improving diagnostic accuracy.

70 Many clinically intelligent decision support products, such as HM Healthcare, iFLYTEK Co.Ltd  
71 intelligent medical assistant, and Baidu 01 Healthcare, have emerged and been applied in many  
72 primary hospitals in recent years in China. However, these studies mainly focus on adult diseases  
73 and lack learning at the pediatric level. In the field of pediatrics, AI technology has been applied  
74 to research on special disease auxiliary diagnosis of children's respiratory diseases,<sup>4</sup> sepsis, and  
75 noninfectious systemic inflammatory response syndrome (SIRS)<sup>5</sup> in foreign countries; however,  
76 there is no research on general pediatric auxiliary diagnosis and treatment yet. In 2019, a deep  
77 learning framework that analyzes more than 140,000 electronic medical records (EMRs) to study  
78 intelligent diagnosis, including the diagnosis of 63 pediatric diseases, using Chinese EMRs was  
79 proposed by Wu et al.<sup>6,7</sup> Moreover, Liang et al.<sup>8</sup> developed an AI intelligent diagnosis model that  
80 learned 1.36 million high-quality electronic text medical records and constructed an AI-assisted  
81 diagnosis model that has an accuracy of nearly 90% in the diagnosis of 55 common pediatric  
82 diseases. However, the existing diagnostic models based on AI are far from sufficient in terms of  
83 learning data or diseases categories, and are still limited in the diagnosis of children's diseases

84 and have not been applied on a large scale.

85 To solve the pain points of the intelligent decision-making research above, the EMRs of millions

86 of outpatients in a single-central hospital for children in Shanghai, China, Chinese Natural

87 Language Processing (NLP) and an end-to-end deep neural network classifier were used for

88 unstructured text processing and learning to imitate the deductive reasoning process of the

89 doctor's brain's assumptions. A deep learning system for pediatric diseases based on AI was

90 developed for initial diagnosis in the pediatric outpatient department in this study.

91

92

## 93 **Methods**

94

### 95 **Training Dataset of the Deep Learning System (DLS)**

96 The study began in March 2021. All outpatient visits records of patients up to eighteen years of

97 age between Jan 1, 2017, and Aug 10, 2020 in the Children's Hospital of Fudan University, were

98 included in this study for model training, regardless of sex, disease group, or disease progression.

99 Department of visit, month of visit, sex, age, chief complaint, physical examination, disease

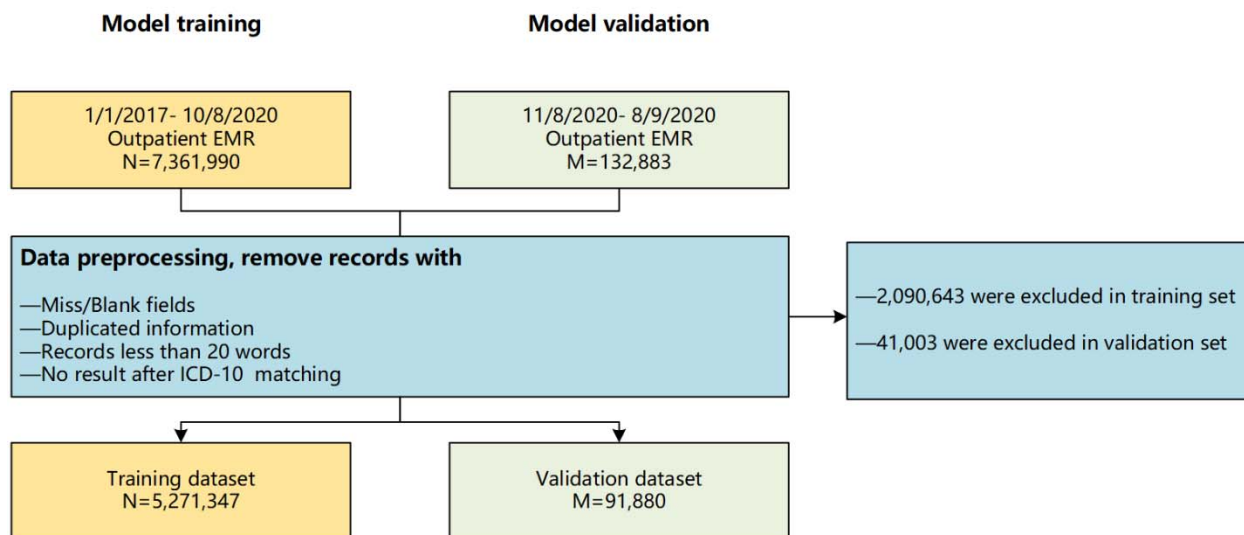
100 history, and primary diagnosis were extracted from the EMR database of the Big Data

101 Management System. This study was approved by the Research Ethics Committee of the  
102 Children’s Hospital of Fudan University. Informed consent was waived for retrospectively  
103 collected outpatient medical record data, which were anonymized. None of the researchers were  
104 able to identify individual participants during or after data collection.

105

## 106 Data Cleaning

107 We excluded records meeting the following criteria (Figure 1): 1) records with missing/blank  
108 fields; 2) records with duplicated information due to template use; and 3) records with  
109 information less than twenty words. We obtained 5,271,347 patient visit records for model  
110 training after filtering. These records cover 306 fine-grained medical majors in the hospital  
111 within the time frame of inclusion.



112

113 **Fig 1. Flow chart of study population inclusion**



## 114 **Data Labeling**

115 The codes from the International Statistical Classification of Diseases and Related Health  
116 Problems, 10<sup>th</sup> Revision (ICD-10) were assigned to each record by exact word matching of  
117 doctor diagnosis in Chinese. Multiple diagnoses were allowed in which case the first diagnosis  
118 was considered the most important main diagnosis. All cases were divided into 16 categories of  
119 diseases based on ICD-10 coding categories, and the codes derived from pregnancy, childbirth,  
120 puerperium (O00-O99), or other diseases without obvious anatomical classification (R00-Z99)  
121 were classified as “other diseases”.

122

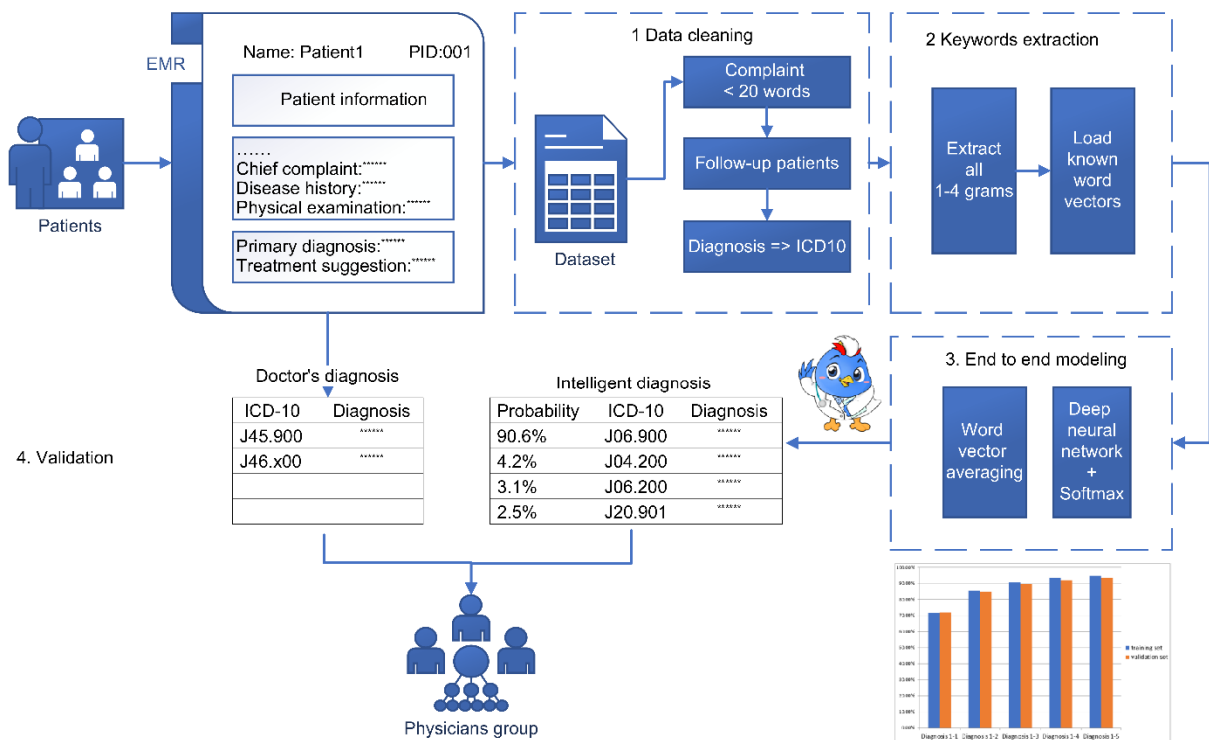
## 123 **NLP Model Construction**

124 First, department of visit, age, sex, chief complaint, physical examination, and disease history are  
125 concatenated into one sentence as input. Our NLP DLS consists of two parts: the Chinese  
126 language feature extractor and the end-to-end deep neural network classifier. The Chinese  
127 language feature extractor extracts 1–4 n-grams of Chinese characters as well as possible  
128 number/alphabet characters. For a sentence of length T, there will be 4T-6 n-grams as features.  
129 We embed this one-hot n-gram into a hidden space of 512 dimensions. The embedding matrix is  
130 initialized with public word vectors from Tencent.<sup>9</sup> For unknown n-grams, a random

131 initialization of  $N(0,0.1)$  is applied.

132 The end-to-end deep neural network consists of three fully connected hidden layers, each of  
133 which consists of 512 hidden units. A softmax layer is employed as a loss layer. A special  
134 modification is applied on the standard softmax layer to allow multiple diagnosis labels. For  $m$   
135 diagnoses, we arbitrarily weight each diagnosis as  $m, m-1, m-2 \dots 1$  because we consider the first  
136 diagnosis as the main and most important diagnosis. We then normalize the weights to 1 and feed  
137 them into softmax as “soft labels”.

138 Stochastic gradient decent (SGD) is employed to train the neural network. We perform input  
139 dropout at a rate of 50% to regularize the network. However, we turn off dropout in the hidden  
140 layers. During the training, 5% of the training dataset is held out to monitor the training process.  
141 The training is stopped when the loss on held-out samples stops improving. All the frameworks  
142 were coded by the author from scratch using the C++ programming language. When the DLS is  
143 connected to the EMR system, after the doctor fills in chief complaint, physical examination, and  
144 disease history, the model will quickly complete the calculation and provide the doctor with  
145 diagnosis opinion in the form of ‘probability value, ICD-10 code, diagnosis name’ list. The  
146 construction process of deep learning system is shown in Figure 2.



147

148 **Fig 2. Workflow of the deep learning algorithm**

149 Department of visit, age, sex, chief complaint, physical examination, and disease history are concatenated  
 150 into one sentence as input dataset. (1) Excluded records meeting the following criteria: a) records with  
 151 missing/blank fields; b) records with duplicated information due to template use; and c) records with  
 152 information less than twenty words. (2) The Chinese language feature extractor extracts 1–4 n-grams of  
 153 Chinese characters as well as possible number/alphabet characters. (3) The end-to-end deep neural  
 154 network. The embedding matrix is initialized with public word vectors from Tencent. (4) In the clinical  
 155 evaluation, an expert group was used to evaluate the parts with great differences between human and  
 156 machine diagnosis results collectively.

## 157 **Validation Dataset**

158 All outpatient visit records of patients up to eighteen years of age between August 11, 2020 and  
159 September 8, 2020 in Children's Hospital of Fudan University, were included in this study for  
160 model testing, regardless of sex, disease group, or disease progression. Department of visit,  
161 month of visit, sex, age, chief complaint, physical examination, disease history, and primary  
162 diagnosis were extracted from the Big Data Management System database (Validation dataset).  
163 After the same filtering pipeline as the training dataset, 91,880 records were obtained for testing.

164

## 165 **Comparison of the Performance of AI System with Human** 166 **Physicians**

167 Referring to the man-machine diagnosis comparison of the Dxpain clinical decision support  
168 system (CDSS), an expert group composed of one senior expert with more than five years of  
169 experience and one expert at or above the associate chief physician level with more than ten  
170 years of experience, was invited to evaluate the parts with great differences between human and  
171 machine diagnosis results collectively in each clinical department during the clinical evaluation  
172 phase. Each encountered medical record was assigned the same number of AI diagnoses based

173 on the number of diagnoses given by the doctor. Experts may choose the more accurate diagnosis  
174 by blind selection between the human diagnosis and the AI diagnosis based on medical record  
175 information consisting of department of visit, month of visit, sex, age, chief complaint, physical  
176 examination, and disease history. There are four possible outcomes in clinical evaluation: (1) AI  
177 Correct: AI's diagnosis is better than that of the doctor's; (2) Physician Correct: the doctors'  
178 diagnosis is better than that of AI; (3) Both Correct: both diagnoses are correct; and (4) Invalid:  
179 both diagnoses are incorrect.

180

## 181 **Statistical Analysis**

182 We defined the following three measurements of AI diagnostic accuracy.

183 1) First diagnostic concordant rate (FDCR): Top-1 accuracy, the percentage of records in the  
184 testing dataset in which any of the AI's top 1 diagnoses are contained in the human doctors'  
185 diagnoses.

186 2) Diagnosis concordance rate (DCR): Top-5 accuracy, the percentage of records in the testing  
187 dataset in which any of the AI's top 5 diagnoses are contained in the human doctors' diagnoses.

188 3) Relative correct rate (RCR): In a blind competition between AI and human doctors, The RCR  
189 is defined as:

$$RCR = \frac{\text{AI Correct} + \text{Both Correct}}{\text{Physician Correct} + \text{Both Correct}}$$

190 When  $RCR > 1$ , the diagnostic ability of AI is stronger than that of the doctor, and when  $RCR < 1$ ,  
191 the diagnostic ability of AI is weaker than that of the doctor.

192

## 193 **Results**

194

### 195 **Baseline Information**

196 In the training dataset, a total of 7,361,990 outpatient EMRs were collected in the outpatient  
197 EMR system from the Children's Hospital of Fudan University. After processing, 5,271,347  
198 medical records were included in model training, covering over 300 kinds of outpatient special  
199 clinics. The median age was approximately 4.58 years (0–18 years), and 56.09% of the patients  
200 were male, which was a higher percentage than female patients. Among the sixteen categories of  
201 diseases, respiratory diseases were the most common (34.53%), followed by infectious diseases  
202 (9.34%) and skin-related diseases (6.80%).

203 To validate the model, 91,880 records were obtained from EMR. There were generally more  
 204 boys than girls, which was the same as the sex distribution in the training set. As shown in Table  
 205 1, there were obvious differences in the composition of disease types between the test dataset and  
 206 the training set, but respiratory diseases were also dominant.

207

208 **Table 1. Basic Data of Outpatient Electronic Medical Records**

	Training dataset	Validation dataset
Age*	4.58±3.41	5.65±3.67
Sex	..	..
Male	2956471 (56.09%)	50065(54.49%)
Female	2314876 (43.91%)	41815 (45.51%)
Disease Groups	..	..
Infectious and parasitic diseases (A00-B99)	492496 (9.34%)	10535 (11.47%)
Tumor (C00-D48)	16682 (0.32%)	711 (0.77%)
Blood diseases (D50-D89)	73054 (1.39%)	952 (1.04%)
Endocrine, nutritional and metabolic diseases (E00-E90)	264459 (5.02%)	10403 (11.32%)
Mental and behavioral disorders (F00-F99)	102949 (1.95%)	2348 (2.56%)
Nervous system diseases (G00-G99)	95097 (1.80%)	2642 (2.88%)
Eye/ear and appendage diseases (H00-H95)	233116 (4.42%)	6479 (7.05%)
Circulatory system diseases (I00-I99)	22659 (0.43%)	574 (0.62%)
Respiratory system diseases (J00-J99)	1820065 (34.53%)	11870 (12.92%)
Digestive system diseases (K00-K93)	336072 (6.38%)	6780 (7.38%)
Skin and subcutaneous tissue diseases (L00-L99)	358481 (6.80%)	8188 (8.91%)
Musculoskeletal and connective tissue diseases (M00-M99)	99387 (1.89%)	2045 (2.23%)
Urogenital diseases (N00-N99)	108433 (2.06%)	3369 (3.67%)
Certain conditions originating in the perinatal period (P00-P96)	62462 (1.18%)	768 (0.84%)
Congenital malformations, deformations and chromosomal abnormalities (Q00-Q99)	94559 (1.79%)	2429 (2.64%)
Other diseases (O00-O99, R00-Z99)	1091376(20.70%)	21787 (23.71%)

209 \*Age: Mean±SD

210

## 211 Performance of the AI Models

212 In the experiment, we focused on the diagnostic performance of the AI model in each type of  
213 disease; that is, to explore the consistency rate of the model with the existing doctor's diagnosis  
214 in different diseases, after model training, the performance of the model in different diseases was  
215 evaluated by calculating the DCR in the two datasets (Table 2), and the corresponding 95%  
216 confidence interval was calculated. Subsequently, for cases where the AI and the doctors  
217 disagreed, a third-party judgment was made by a clinical expert group in each specialty to  
218 explore who had a more accurate diagnosis.

219

220 **Table 2. Illustration of the Diagnostic Performance of the AI Model (training dataset)**

Diseases groups	N	DCR %	95% CI
Infectious and parasitic diseases (A00-B99)	492496	96.70	(96.65~96.75)
Tumor (C00-D48)	16682	75.21	(74.56~75.87)
Blood diseases (D50-D89)	73054	92.81	(92.62~93.00)
Endocrine, nutritional and metabolic diseases (E00-E90)	264459	97.11	(97.04~97.17)
Mental and behavioral disorders (F00-F99)	102949	94.70	(94.56~94.83)
Nervous system diseases (G00-G99)	95097	92.55	(92.38~92.72)
Eye/ear and appendage diseases (H00-H95)	233116	98.26	(98.21~98.31)
Circulatory system diseases (I00-I99)	22659	80.13	(79.62~80.65)
Respiratory system diseases (J00-J99)	1820065	96.86	(96.84~96.89)
Digestive system diseases (K00-K93)	336072	94.31	(94.23~94.39)
Skin and subcutaneous tissue diseases (L00-L99)	358481	95.36	(95.29~95.43)
Musculoskeletal and connective tissue diseases (M00-M99)	99387	87.62	(87.42~87.83)
Urogenital diseases (N00-N99)	108433	93.16	(93.01~93.31)
Certain conditions originating in the perinatal period (P00-P96)	62462	93.05	(92.85~93.25)



---

Congenital malformations, deformations and chromosomal abnormalities (Q00-Q99)	94559	88.84	(88.64~89.05)
Other diseases (O00-O99, R00-Z99)	1091376	94.90	(94.86~94.94)
All diseases	5271347	95.49	(95.48~95.51)

---

221 DCR%: Diagnosis concordance rate, the percentage of records in the testing dataset in which any  
222 of the AI's top 5 diagnoses are contained in the human doctors' diagnoses.

223 95% CI: 95% confidence interval.

224 After over five million cases of training, the overall diagnostic consistency rate reached 95.49%

225 (95.48~95.51). Among the different diseases, the models performed best in eye/ear and

226 appendage diseases (H00-H95). On the other hand, the DCRs were relatively low in diagnoses of

227 categories such as tumor (C00-D48), circulatory system diseases (I00-I99), musculoskeletal and

228 connective tissue diseases (M00-M99), and any other congenital diseases (Q00-Q99), with DCRs

229 of less than 90%.

230

## 231 **Model Validation**

232 In local validation, the overall DCR reached 93.51% (93.35~93.67) while FDCR was 72.04%

233 (71.75~72.33), and the AI model performed better in eye/ear and appendage diseases(H00-H95),

234 infectious and parasitic diseases (A00-B99), endocrine, nutritional and metabolic diseases (E00-

235 E90), and respiratory system diseases (J00-J99), with DCRs exceeding 95%. In contrast, it

236 performed poorly in tumor (C00-D48) and circulatory system diseases (I00-I99), with DCRs

237 lower than 80%. In addition, the model did not perform well in the first diagnosis of Circulatory  
238 system diseases (I00-I99) and Certain conditions originating in the perinatal period (P00-P96),  
239 but the accuracy was significantly improved after expanded to the top five diagnoses. There were  
240 significant differences in the performance of the model in the training set and validation set even  
241 if the cases came from the same hospital, especially in the tumor, blood diseases, and certain  
242 conditions originating in the perinatal period (>5%). In general, the DCRs of most disease types  
243 in the test set were lower than those of the training set, but the DCRs of some disease types were  
244 slightly higher, such as in infectious and parasitic diseases (A00-B99), as well in eye/ear and  
245 appendage diseases (H00-H95).

246

247 **Table 3. Illustration of the Diagnostic Performance of the AI Model (validation dataset)**

248

Diseases groups	N	FDCR %	95% CI	DCR %	95% CI
Infectious and parasitic diseases (A00-B99)	10534	81.69	(80.95~82.43)	97.07	(96.74~97.39)
Tumor (C00-D48)	707	38.05	(34.47~41.63)	68.46	(65.03~71.89)
Blood diseases (D50-D89)	952	67.54	(64.57~70.52)	85.29	(83.04~87.54)
Endocrine, nutritional and metabolic diseases (E00-E90)	10403	71.75	(70.88~72.61)	96.68	(96.34~97.03)
Mental and behavioral disorders (F00-F99)	2346	67.22	(65.32~69.12)	90.41	(89.22~91.60)
Nervous system diseases (G00-G99)	2642	79.75	(78.22~81.28)	91.68	(90.62~92.73)
Eye/ear and appendage diseases (H00-H95)	6478	87.09	(86.28~87.91)	98.32	(98.00~98.63)
Circulatory system diseases (I00-I99)	574	44.43	(40.36~48.49)	76.13	(72.65~79.62)
Respiratory system diseases (J00-J99)	11869	69.21	(68.38~70.04)	95.69	(95.32~96.05)
Digestive system diseases (K00-K93)	6777	73.60	(72.55~74.65)	93.43	(92.84~94.02)

---

Skin and subcutaneous tissue diseases (L00-L99)	8188	68.45	(67.45~69.46)	93.61	(93.08~94.14)
Musculoskeletal and connective tissue diseases (M00-M99)	2038	63.69	(61.60~65.78)	83.37	(81.75~84.98)
Urogenital diseases (N00-N99)	3369	74.56	(73.09~76.03)	92.52	(91.63~93.41)
Certain conditions originating in the perinatal period (P00-P96)	768	55.99	(52.48~59.50)	86.20	(83.76~88.64)
Congenital malformations, deformations and chromosomal abnormalities (Q00-Q99)	2426	65.62	(63.73~67.51)	84.79	(83.36~86.22)
Other diseases (O00-O99, R00-Z99)	21781	68.73	(68.11~69.35)	92.14	(91.79~92.50)
All diseases	91852	72.04	(71.75~72.33)	93.51	(93.35~93.67)

---

249 FDCR: First diagnosis concordance rate, the percentage of records in the testing dataset in which  
250 any of the AI's top 1 diagnoses are contained in the human doctors' diagnoses.

251 DCR: Diagnosis concordance rate, the percentage of records in the testing dataset in which any  
252 of AI's top 5 diagnoses are contained in the human doctors' diagnoses.

253 95% CI: 95% confidence interval.

254

## 255 **Third-party Evaluation**

256 The medical records with inconsistent diagnoses between the AI model and the doctor were

257 evaluated by clinical experts. According to the evaluation principle, two doctors in each specialty

258 jointly evaluated the cases and provided the evaluation results of Validation set. Due to the

259 differences in disease distributions and AI diagnostic capabilities in different diseases, the

260 number of cases in the validation set was inconsistent. Based on the results of the relative

261 accuracy calculation (Table 3), experts believe that the accuracy of the model was generally

262 higher than that of doctors.

263

264 **Table 3. Validation of the Diagnostic Performance of the AI Model and Physicians**

Diseases groups	Model validation at the research site					
	N	AI Correct	Physician Correct	Both Correct	Invalid	RCR
Infectious and parasitic diseases (A00-B99)	82	25	15	18	24	1.30
Tumor (C00-D48)	31	9	5	4	13	1.44
Blood diseases (D50-D89)	21	9	3	0	9	3.00
Endocrine, nutritional and metabolic diseases (E00-E90)	57	19	16	6	16	1.14
Mental and behavioral disorders (F00-F99)	27	12	3	4	8	2.29
Nervous system diseases (G00-G99)	43	13	9	3	18	1.33
Eye/ear and appendage diseases (H00-H95)	43	12	12	6	13	1.00
Circulatory system diseases (I00-I99)	20	5	1	1	13	3.00
Respiratory system diseases (J00-J99)	86	45	11	10	20	2.62
Digestive system diseases (K00-K93)	87	28	18	14	27	1.31
Skin and subcutaneous tissue diseases (L00-L99)	110	33	25	7	45	1.25
Musculoskeletal and connective tissue diseases (M00-M99)	22	7	4	6	5	1.30
Urogenital diseases (N00-N99)	48	12	10	3	23	1.15
Certain conditions originating in the perinatal period (P00-P96)	20	8	0	4	8	3.00
Congenital malformations, deformations and chromosomal abnormalities (Q00-Q99)	73	21	11	13	28	1.42
Other diseases (O00-O99, R00-Z99)	290	94	65	37	94	1.28

265 AI: artificial intelligence, RCR: Relative correct rate = (AI Correct + Both Correct)/(Physician  
 266 Correct + Both Correct). When RCR>1, the diagnostic ability of AI is stronger than that of the  
 267 doctor, and when RCR<1, the diagnostic ability of AI is weaker than that of the doctor.

268

## 269 Discussion

270 In this study, we developed and validated a pediatric diagnosis CDSS based on a deep learning  
 271 system. Utilizing electronic health records, AI achieved high performance with satisfying DCRs

272 and RCR. In the prospective validation dataset, AI outperformed physicians from a tertiary  
273 hospital. In these two datasets, there were generally more boys than girls, which is the same as  
274 the distribution in other similar studies.<sup>8</sup> There were obvious differences in the mean age and  
275 composition of disease types among the two datasets, which may be because the validation  
276 dataset was taken from a continuous month, resulting in a bias caused by seasonal factors.

277 Although FDCR (top-1 accuracy) is only 72.04% (71.75~72.33), DCR (top-5 accuracy) has a  
278 significant improvement of 93.51% (93.35~93.67). Considering the complexity of the disease  
279 and the compatibility of diagnostic codes, there are more than one diagnostic choice in a case as  
280 a multi-label classification problem. This study mainly focused on the DCR index. As shown in  
281 the results, the AI model performed better in internal medical diseases, including infectious and  
282 parasitic diseases, endocrine/nutritional and metabolic diseases, eye/ear and appendage diseases,  
283 respiratory system diseases, etc., and best performance was in eye/ear and appendage diseases,  
284 which was because of the more concentrated distribution of common diseases in the clinic.

285 However, in surgical diseases, such as musculoskeletal and connective tissue diseases, the DCR  
286 was lower. This can be explained because surgical diseases often need to be combined with  
287 various laboratory tests and radiographic examinations for auxiliary diagnosis in the clinic.

288 Therefore, AI diagnosis was relatively correct and could provide a reference for local doctors.  
289 Pediatric medical resources are scarce and unevenly distributed in China. As of 2015, the  
290 shortage of pediatricians was as high as 200,000, and the ratio of pediatric doctors to patients  
291 was 1:3500 in China, which is far lower than the 1:1000 doctor-patient standard in developed  
292 countries.<sup>10</sup> Similar situations are common in Europe,<sup>11</sup> the US<sup>12</sup> and Japan.<sup>13</sup> In addition, the  
293 excessive concentration of high-quality medical resources leads to the exposure of primary  
294 pediatricians to fewer diseases and cases, and the lack of clinical experience results in a high  
295 incidence of misdiagnoses and diagnostic errors in primary hospitals, which delays timely access  
296 to treatment for patients. In 2013, a survey report from the American Academy of Pediatrics’  
297 Quality Improvement Innovation Networks (QuIIN) showed that 35% of pediatricians may make  
298 a diagnostic error at least once a month, and 33% have a diagnostic error at least once a year,  
299 finally resulting in adverse events.<sup>14</sup> It is supposed that the insufficient clinical knowledge of  
300 doctors, the incorrect collection of consultation information, and invalid inspection and  
301 verification are the main reasons for clinical misdiagnosis.<sup>15</sup> However, when retrieving the  
302 relevant research at home and abroad in the last ten years based on the keywords “Patient safety”  
303 and “Diagnostic Errors” in PubMed, only 6% focused on the diagnostic safety of children. A

304 previous study reported that the accuracy rate of the first diagnosis of pediatricians in the  
305 primary hospital was relatively low.<sup>16</sup> In addition, the diagnostic accuracy of primary units at  
306 lower professional levels (e.g., level 1 and level 2) was lower than that of level 3. Research on  
307 reducing the pediatric misdiagnosis rate is far behind the progress of adult and other patient  
308 safety fields.<sup>17</sup> Regarding the abovementioned status, many suggested strategies to reduce the  
309 incidence of diagnostic errors have been proposed, including improving the clinical expertise of  
310 clinicians and reducing the inherent cognitive errors of doctors to make better decisions.<sup>15</sup>  
311 Therefore, the application of AI systems could be beneficial among areas where healthcare  
312 providers are in a relative shortage. It will greatly increase the accuracy of the first diagnosis and  
313 thus significantly improve medical health care in Chinese children. In addition, this AI model  
314 could provide a new idea for the quality control of EMRs in the outpatient departments of local  
315 hospitals.

316 The AI system demonstrates high diagnostic accuracy across multiple diseases and is comparable  
317 to well-trained pediatricians in diagnosing common pediatric diseases. In the past, the CDSS was  
318 widely used to extract key clinical information based on EMRs for reasoning, which could help  
319 clinicians assess disease status, make diagnoses, choose appropriate treatments and make other

320 clinical decisions.<sup>18</sup> The clinical diagnosis and treatment guidelines and a large number of  
321 literary studies have provided a reliable source of knowledge for the system and have been  
322 universally recognized. However, there are few systems that have been put into extensive and  
323 long-term application in clinical practice worldwide. On the one hand, the construction of a  
324 knowledge base cannot meet the needs of clinicians; furthermore, most systems are not  
325 technically integrated with EMRs, resulting in disconnection from the clinical workflow, which  
326 reduces the enthusiasm of clinicians to use the CDSS. With the rapid development of medical  
327 information technology, AI technology has become a powerful support for the healthcare  
328 revolution. With patients as the center and medical institutions as the main service body, medical  
329 AI can cover the whole process from disease prevention, diagnosis, and treatment to patient  
330 rehabilitation by using AI technologies such as knowledge mapping or deep learning.<sup>19</sup> At the  
331 same time, with the widespread application of EMRs in recent years, combined with large data-  
332 level EMRs and other related information learning, AI algorithms can complete complex analysis  
333 tasks in a short period of time, feedback the best classification model results based on the input  
334 information, and assist doctors in improving the accuracy and efficiency of patient diagnosis.  
335 “Data-driven” intelligent aid decision-making based on real-world EMR data will have the



336 potential to supplement traditional rule-based decision-making methods.<sup>20</sup>

337

## 338 **Limitations**

339 First, this AI system is only suitable for identifying the diagnosis at first clinic at the present

340 stage. Second, results of laboratory medicine and imaging examination were not included in the

341 AI system. Third, this AI system cannot provide disease severity classification or treatment

342 suggestions yet. Fourth, this study lacks a comparison of multiple model algorithms.

343

## 344 **Conclusions**

345 In our study, we developed and validated an AI-based system that can provide clinical decision

346 support in the event of diagnostic uncertainty and complexity. The AI model can quickly and

347 accurately identify the diagnosis of children, which can help pediatricians make more precise

348 diagnoses while further preventing undetected cases. Our findings are of great clinical value and

349 practical significance in improving the health care of children in China and optimizing medical

350 resources.

351

## 352 **Contributors**

353 XG, YW, XZ, HX conceived of and designed the study. XG, YW, YS, AM, CY, WL, YZ, ZH, JH,

354 XZ and H X contributed to Acquisition, analysis, or interpretation of data. XG, YW, LX, YS, XZ,

355 HX contributed to draft of the manuscript, which was approved by all the authors. XG, YW, YS,

356 YZ, XZ, HX contributed to critical revision of the manuscript for important intellectual content.

357 XG, YW, LX, YS contributed to statistical analysis. XZ, HX contributed to supervision. HX and

358 XZ had full access to all the data in the study and take responsibility for the integrity of the data

359 and the accuracy of the data analysis. XG and YW contributed equally as co-first authors to this

360 work.

361

## 362 **Declaration of interests**

363 The authors have no conflicts of interest to declare

364

## 365 **Data sharing**

366 The institutional data used for training and validation are not publicly available, because they

367 contain protected patient health information. Source code of the deep neural network can be

368 made available, subject to intellectual property constraints, by contacting the co-first author  
369 (wangyi\_fudan@fudan.edu.cn).

370

## 371 **Acknowledgments**

372 For financial assisting and technical advice, we thank the project of Shanghai's Double First-  
373 Class University Construction and Development of High-Level Local Universities: Intelligent  
374 Medicine Emerging Interdisciplinary Cultivation Project, and Medical Research Data Center of  
375 Fudan University.

376

## 377 **References**

- 378 1. Cui LX. Case analysis of pediatric misdiagnosis in primary hospitals. *The Journal of*  
379 *Medical Theory and Practice*. 2020;33(6):968-969.
- 380 2. Grubenhoff JA, Ziniel SI, Cifra CL, Singhal G, McClead RE, Jr., Singh H. Pediatric  
381 clinician comfort discussing diagnostic errors for improving patient safety: a survey.  
382 *Pediatr Qual Saf*. 2020;5(2):e259.
- 383 3. Li Y, Zhang T, Yang Y, Gao Y. Artificial intelligence-aided decision support in paediatrics

- 384 clinical diagnosis: development and future prospects. *J Int Med Res.*  
385 2020;48(9):300060520945141.
- 386 4. Porter P, Abeyratne U, Swarnkar V, et al. A prospective multicentre study testing the  
387 diagnostic accuracy of an automated cough sound centred analytic system for the  
388 identification of common respiratory disorders in children. *Respir Res.* 2019;20(1):81.
- 389 5. Lamping F, Jack T, Rübsamen N, et al. Development and validation of a diagnostic model  
390 for early differentiation of sepsis and non-infectious SIRS in critically ill children - a  
391 data-driven approach using machine-learning algorithms. *BMC Pediatr.* 2018;18(1):112.
- 392 6. Li X, Wang H, He H, Du J, Chen J, Wu J. Intelligent diagnosis with Chinese electronic  
393 medical records based on convolutional neural networks. *BMC Bioinformatics.*  
394 2019;20(1):62.
- 395 7. WU JZ, Luo Z, XU S, et al. Realize the Intelligent Auxiliary Diagnosis of Pediatric  
396 Clinical Disease with the Application of Deep Learning. *China Digital Medicine.*  
397 2018;13(10):14-16.
- 398 8. Liang H, Tsui BY, Ni H, et al. Evaluation and accurate diagnoses of pediatric diseases  
399 using artificial intelligence. *Nat Med.* 2019;25(3):433-438.

- 400 9. Song Y , Shi S , Li J , et al. Directional Skip-Gram: Explicitly Distinguishing Left and  
401 Right Context for Word Embeddings. NAACL 2018.2018; (Short Paper).
- 402 10. Sun XJ, Luo CJ. An Analysis of Pediatric Treatment Dilemma form the Perspective  
403 of Medical Supply Reform. Journal of Chongqing Three Gorges University. 2021;  
404 37(1):62-70.
- 405 11. Ehrich J, Fruth J, Jansen D, Gerber-Grote A, Pettoello-Mantovani M. How to calculate  
406 the risk of shortage and surplus of pediatric workforce? J Pediatr. 2018;199:286-287.e2.
- 407 12. Eden AR, Morgan ZJ, Jetty A, Peterson LE. Proportion of family physicians caring for  
408 children is declining. J Am Board Fam Med. 2020;33(6):830-831.
- 409 13. Sasaki H, Otsubo T, Imanaka Y. Widening disparity in the geographic distribution of  
410 pediatricians in Japan. Hum Resour Health. 2013;11:59.
- 411 14. Rinke ML, Singh H, Ruberman S, et al. Primary care pediatricians' interest in diagnostic  
412 error reduction. Diagnosis (Berl). 2016;3(2):65-69.
- 413 15. Thammasitboon S, Cutrer WB. Diagnostic decision-making and strategies to improve  
414 diagnosis. Curr Probl Pediatr Adolesc Health Care. 2013;43(9):232-241.
- 415 16. Wong GW, Kwon N, Hong JG, Hsu JY, Gunasekera KD. Pediatric asthma control in Asia:

- 416 phase 2 of the Asthma Insights and Reality in Asia-Pacific (AIRIAP 2) survey. *Allergy*.  
417 2013;68(4):524-530.
- 418 17. McDonald KM, Matesic B, Contopoulos-Ioannidis DG, et al. Patient safety strategies  
419 targeted at diagnostic errors: a systematic review. *Ann Intern Med*. 2013;158(5 Pt 2):381-  
420 389.
- 421 18. Jia PL, Zhang PF, Li HD, Zhang LH, Chen Y, Zhang MM. Literature review on clinical  
422 decision support system reducing medical error. *J Evid Based Med*. 2014;7(3):219-226.
- 423 19. Ahmed Z , Mohamed K , Zeeshan S , et al. Artificial intelligence with multi-functional  
424 machine learning platform development for better healthcare and precision medicine.  
425 Database *The Journal of Biological Databases and Curation*, 2020, 2020.
- 426 20. Hoffmann M, Vander Stichele R, Bates DW, et al. Guiding principles for the use of  
427 knowledge bases and real-world data in clinical decision support systems: report by an  
428 international expert workshop at Karolinska Institutet. *Expert Rev Clin Pharmacol*.  
429 2020;13(9):925-934.

## Model training

## Model validation

1/1/2017- 10/8/2020  
Outpatient EMR  
N=7,361,990

11/8/2020- 8/9/2020  
Outpatient EMR  
M=132,883

### Data preprocessing, remove records with

- Miss/Blank fields
- Duplicated information
- Records less than 20 words
- No result after ICD-10 matching

—2,090,643 were excluded in training set  
—41,003 were excluded in validation set

Training dataset  
N=5,271,347

Validation dataset  
M=91,880

