

# Can natural language processing help differentiate inflammatory intestinal diseases in China? Models applying random forest and convolutional neural network approaches

**Yuanren Tong**

Chinese Academy of Medical Sciences & Peking Union Medical College Institute of Biomedical Engineering

**Keming Lu**

Tsinghua University

**Yingyun Yang**

Peking Union Medical College Hospital

**Ji Li**

Peking Union Medical College Hospital

**Yucong Lin**

Tsinghua University

**Dong Wu**

Peking Union Medical College Hospital

**Aiming Yang**

Peking Union Medical College Hospital

**yue li (✉ [yuelipumch@126.com](mailto:yuelipumch@126.com))**

Peking Union Medical College Hospital <https://orcid.org/0000-0002-1133-7317>

**Sheng Yu**

Tsinghua University

**Jiaming Qian**

Peking Union Medical College Hospital

---

## Research article

**Keywords:** inflammatory bowel disease, intestinal tuberculosis, natural language processing

**Posted Date:** September 30th, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-39653/v3>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published on September 29th, 2020. See the published version at <https://doi.org/10.1186/s12911-020-01277-w>.

1 **Can natural language processing help differentiate inflammatory intestinal diseases in**  
2 **China? Models applying random forest and convolutional neural network approaches**

3 Author list: Yuanren Tong<sup>1†</sup>; Keming Lu<sup>2†</sup>; Yingyun Yang<sup>1</sup>; Ji Li<sup>1</sup>; Yucong Lin<sup>3,4</sup>; Dong Wu<sup>1</sup>;  
4 Aiming Yang<sup>1</sup>; Yue Li<sup>1\*</sup>; Sheng Yu<sup>3,4,5\*</sup>; Jiaming Qian<sup>1</sup>

5 1. Department of gastroenterology, Peking Union Medical College Hospital, Chinese Academy  
6 of Medical Sciences and Peking Union Medical College, Beijing, China, 100730

7 2. Department of Automation, Tsinghua University, Beijing, China, 100084

8 3. Center for Statistical Science, Tsinghua University, Beijing, China, Beijing, China, 100084

9 4. Department of Industrial Engineering, Tsinghua University, Beijing, China, 100084

10 5. Institute for Data Science, Tsinghua University, Beijing, China. 100084

11 † These authors contributed equally to this article.

12

13 Mr. Yuanren Tong: tongyr14@mails.tsinghua.edu.cn

14 Mr. Keming Lu: lkm16@mails.tsinghua.edu.cn

15 Dr. Yingyun Yang: yingyunyang@126.com

16 Dr. Ji Li: lij0235@pumch.cn

17 Dr. Yucong Lin: linc16@mails.tsinghua.edu.cn

18 Dr. Dong Wu: wudong@pumch.cn

19 Prof. Aiming Yang: yangam@pumch.cn

20 Prof. Jiaming Qian: qianjiaming1957@126.com

21 Dr. Yue Li: yuelee76@gmail.com

22 Dr. Sheng Yu: syu@tsinghua.edu.cn

23 **\*Corresponding author:**

24 **Yue Li, M.D.**

25 [yuelee76@gmail.com](mailto:yuelee76@gmail.com)

26 Tel/Fax: +86-10-69155751

27 Department of Gastroenterology, Peking Union Medical College Hospital, Chinese Academy

28 of Medical Sciences and Peking Union Medical College, Beijing, China, 100730

29 **Sheng Yu, Ph.D.**

30 [syu@tsinghua.edu.cn](mailto:syu@tsinghua.edu.cn)

31 Tel/Fax: +86-10-62783842

32 Center for Statistical Science& Department of Industrial Engineering& Institute for Data

33 Science, Tsinghua University, Beijing, China, 100084

34 **Conference:**

35 The abstract of this article has won the **first prize of the Young Investigator Award** during

36 the **Asian Pacific Digestive Week (APDW) 2019** held in Kolkata, India

37

38 **Abstract**

39 **Background:** Differentiating between ulcerative colitis (UC), Crohn's disease (CD) and  
40 intestinal tuberculosis (ITB) using endoscopy is challenging. We aimed to realize automatic  
41 differential diagnosis among these diseases through machine learning algorithms.

42 **Methods:** A total of 6399 consecutive patients (5128 UC, 875 CD and 396 ITB) who had  
43 undergone colonoscopy examinations in the Peking Union Medical College Hospital from  
44 January 2008 to November 2018 were enrolled. The input was the description of the  
45 endoscopic image in the form of free text. Word segmentation and key word filtering were  
46 conducted as data preprocessing. Random forest (RF) and convolutional neural network  
47 (CNN) approaches were applied to different disease entities. Three two-class classifiers (UC  
48 and CD, UC and ITB, and CD and ITB) and a three-class classifier (UC, CD and ITB) were  
49 built.

50 **Results:** The classifiers built in this research performed well, and the CNN had better  
51 performance in general. The RF sensitivities/specificities of UC-CD, UC-ITB, and CD-ITB  
52 were 0.89/0.84, 0.83/0.82, and 0.72/0.77, respectively, while the values for the CNN of  
53 CD-ITB were 0.90/0.77. The precisions/recalls of UC-CD-ITB when employing RF were  
54 0.97/0.97, 0.65/0.53, and 0.68/0.76, respectively, and when employing the CNN were  
55 0.99/0.97, 0.87/0.83, and 0.52/0.81, respectively.

56 **Conclusions:** Classifiers built by RF and CNN approaches had excellent performance when  
57 classifying UC with CD or ITB. For the differentiation of CD and ITB, high specificity and  
58 sensitivity were achieved as well. Artificial intelligence through machine learning is very

59 promising in helping unexperienced endoscopists differentiate inflammatory intestinal  
60 diseases.

61 **Key words:** inflammatory bowel disease; intestinal tuberculosis; natural language  
62 processing

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77 **Key Summary:**

78 **Summarize the established knowledge on this subject:**

- 79 ● Differential diagnosis among UC, CD and ITB is an important clinical problem.
- 80 ● Endoscopy is very important in differential diagnosis among UC, CD and ITB.
- 81 ● Standard terminology system has been built for describing endoscopic images and  
82 many endoscopic features are summarized.
- 83 ● The accuracy of differential diagnosis relies on the experience of clinical doctors.

84 **What are the significant and/or new findings of this study?**

- 85 ● This study built classifiers with two algorithms that had high accuracy and recall rate,  
86 using the descriptions of endoscopic images in the form of free-text.
- 87 ● The classifiers could allow clinical doctors to make accurate differential diagnosis  
88 barely based on objective descriptions of the endoscopic images.
- 89 ● The classifiers could help find new features or new combination of features that could  
90 contribute to differential diagnosis
- 91 ● This research provided a new point of view of combining natural language processing  
92 and clinical endoscopy.

93

94

95

96

97 **Introduction**

98 Inflammatory bowel disease (IBD), including ulcerative colitis (UC) and Crohn's disease (CD),  
99 are idiopathic and chronic digestive tract inflammatory diseases with repeated remission and  
100 relapses. Intestinal tuberculosis (ITB) is intestinal inflammation due to *Mycobacterium*  
101 *tuberculosis* infection. The ultimate courses and prognoses of IBD and ITB can be different.  
102 For example, ITB can be cured with early diagnosis and proper antituberculosis treatment,  
103 while recurrences of CD and UC are common and require life-long follow-up. Misdiagnosis  
104 and inappropriate treatment can cause either a prolonged disease course or severe adverse  
105 effects [1]. Thus, timely and accurate diagnosis and differentiation of IBD and ITB are very  
106 important, especially in China, which has the third-highest tuberculosis incidence rate in the  
107 world according to the World Health Organization [2].

108

109 Endoscopic examination is of vital importance in the diagnosis of IBDs and is always  
110 conducted first. The endoscopic features of IBDs and ITB have been well described [3-5].  
111 However, the three inflammatory intestinal diseases, especially ITB and CD, can be difficult to  
112 distinguish because they have very similar manifestations in terms of both clinical symptoms  
113 and endoscopic appearance. This fact often causes incorrect endoscopic diagnosis and results  
114 in delayed treatment.

115

116 There have been some studies focused on performing differential diagnosis between IBD and  
117 ITB. Yue Li et al. pointed out that interferon  $\gamma$ -release assays could help distinguish CD and

118 ITB [6]. Ren Mao et al. proposed that computed tomographic enterography assisted in  
119 differentiating CD and ITB [7]. **These studies conducted differentiation from a traditional**  
120 **point of view, which could not provide immediate feedback to clinicians.**

121

122 **Some recent studies used statistical methods to make differentiations. Yao He et al. built**  
123 **a comprehensive diagnostic nomogram to differentiate between CD and ITB [8]. Y. J. Lee**  
124 **et al. formulated a scoring system to evaluate the weight of each endoscopic feature for**  
125 **identifications of ITB and CD [9]. These studies provided novel insights from multiple**  
126 **aspects into differential diagnosis between IBDs and ITB. However, the sample sizes in**  
127 **these studies were relatively small. In addition, the scoring system developed by Y. J. Lee**  
128 **has not been validated in a subsequent group; thus, the result might be overly optimistic.**  
129 **Another point is that previously proposed diagnosis methods require much manual work**  
130 **and can be time consuming; thus, they are not able to provide timely insights for clinical**  
131 **doctors.**

132

133 The aim of this study was to perform more accurate endoscopic diagnosis and differentiation  
134 of IBD and ITB with the help of natural language processing and machine learning in order to  
135 assist physicians, especially those with limited experience. In addition, the applications of  
136 these novel algorithms to clinical studies would add to the literature, which could promote  
137 the application of artificial intelligence to clinical problems.

138

139 **Patients/Material and methods**

140 **Overview of the study**

141 With standard terminologies for clinical doctors to describe endoscopic images already  
142 established [10-12], we developed several classifiers based on different algorithms to  
143 automatically classify UC, CD and ITB with free-text endoscopic descriptions as the input.

144

145 **Study population**

146 Electronic health records (EHRs) of a total of 6399 consecutive patients who had undergone  
147 colonoscopies in Peking Union Medical College Hospital (PUMCH) and were clinically  
148 diagnosed as having UC (n=5128), CD (n=875), or ITB (n=396) from January 2008 to  
149 November 2018 were collected successively. This research was approved by the Ethics  
150 Committee of Peking Union Medical College Hospital on September 31<sup>st</sup>, 2018 (IRB S-K894).

151 The clinical diagnoses of UC and CD were made via a combination of medical history,  
152 endoscopic features, pathological features, and treatment follow-up based on the Chinese  
153 consensus of IBD (2018) by IBD specialists in this hospital. Diagnosis of ITB was obtained by  
154 the presence of any of the following: (1) positive acid-fast bacilli on histological examination  
155 or positive *M. tuberculosis* culture; (2) radiological, colonoscopic, and/or other proven TB; or  
156 (3) full response to anti-TB therapy. Colonoscopies were performed by well-trained  
157 endoscopists at PUMCH using Olympus CF-Q260 or H260 colonoscopes.

158

159 Based on the well-established terminology used by endoscopists to describe colonoscopic  
160 images, we extracted descriptions of colonoscopic images of the patients' index colonoscopy  
161 in the form of free text. Clinically confirmed diagnoses extracted from the hospital information  
162 system (HIS) were used as labels.

163

#### 164 **Data processing**

165 Figure 1 shows the flow path of data processing. **An example of input data could be found**  
166 **in the supplementary material.**

167

168 The image descriptions were preprocessed with natural language processing (NLP)  
169 techniques to extract linguistic features before being input into the classifiers. First, Chinese  
170 word segmentation was applied to the description to tokenize the input text, using the Python  
171 package 'jieba' and enhanced by the Xiangya Professional Medical Dictionary. Punctuation  
172 and words without actual clinical meanings such as 'patients' and 'prepare for the examination'  
173 were deleted.

174

175 The second step of NLP was keyword filtering, which aimed to identify informative keywords  
176 in the description. Term frequency-inverse document frequency (TF-IDF) was applied to filter  
177 keywords. The TF-IDF value was defined as follows:

178

$$\text{TF}(X) = \frac{\text{Number of times term X appears in the document}}{\text{Total number of terms in the document}},$$

179 
$$\text{IDF}(X) = \ln \frac{\text{Total number of documents}}{\text{Number of documents containing term } X}$$

180 
$$\text{TF-IDF} = \text{TF} \times \text{IDF}.$$

181 Terms whose TF-IDF values for a document were out of the range of [0.3, 0.7] were removed  
182 from the input for that document.

183

184 The last step of NLP was dimension reduction with non-negative matrix factorization (NMF)  
185 [13], which could improve the interpretability of the extracted features and allow clinical  
186 doctors to understand the results better.

187

188 In addition to single words, N-grams and L1 regularization (also known as the least absolute  
189 shrinkage and selection operator, or LASSO) were applied in the above process as well. Details  
190 were provided in Supplementary Materials for conciseness.

191

## 192 **Development of classifiers**

193 Random forest (RF) is appreciated for the advantage in weighting the importance of features.

194 Convolutional neural networks (CNNs) were selected for their capability to extract features  
195 automatically. In addition, all algorithms applied are able to analyze free text directly, thus  
196 requiring little manual work.

197

198 RF was applied to two-class classifications (UC and CD, UC and ITB, and CD and ITB) and the

199 three-class classification of UC, CD and ITB, while the CNN was applied to the three-class  
200 classification and the CD/ITB classification. The reason for not applying the CNN to the two-  
201 class classifications including UC was an unbalanced sample number.

202

203 **The labeled dataset was randomly split into a training group (70%) and testing group**  
204 **(30%) for both RF and CNN.**

205 *(1) Random forest*

206 RF was applied to data processed by NMF. The RF parameters can be found in Supplementary  
207 Materials Table 1. Due to the unbalanced data, cost-sensitive learning was employed in order  
208 to assign different weights to different diseases. This approach could improve the  
209 performance of the model of CD and ITB, which had a small number of samples.

210

211 **For choosing hyperparameters on RF, ten-fold cross-validation was applied on the**  
212 **training set (70% of total data). The training set were split into ten equal subsets**  
213 **randomly. Nine subsets were used to train for hyperparameters, while the remaining**  
214 **subset was used as the validation set. The chosen hyperparameters were then applied**  
215 **for further train and test.**

216

217 Regarding feature extraction, we first extracted 50 features **sorted by variable importance**  
218 **of RF**, which comprised phrases produced by segmentation. All the features **were then**

219 reviewed by two experienced clinical doctors to combine similar features and omit  
220 meaningless or duplicated features.

221

## 222 *(2) Convolutional neural network*

223 To train CNN models, we extracted approximately 110,000 descriptions from the endoscopic  
224 center of PUMCH without labels. The data were employed to train a GloVe[14] model for  
225 word embedding, which was used to initiate the CNN embedding layer.

226

227 The CNN model applied a structure similar to the Text-CNN model proposed by Yoon Kim  
228 [15]. A word list was built, and an integer was allocated to each word. The input sentences  
229 were segmented into words and represented by a corresponding integer sequence. The  
230 integer sequence was then embedded into a 100-dimension vector. The vectors were used  
231 as the input for the CNN. The parameters of the CNN can be found in Supplementary  
232 Materials Table 2.

233

234 The model was optimized by the Adam algorithm [16]. FocalLoss was used as the loss function  
235 to handle the imbalanced sample sizes and accelerate convergence [17]. A total of 100  
236 iterations were employed to allow the model to converge.

237

## 238 **Result visualization**

239 T-distributed stochastic neighbor embedding (t-SNE) [18] was applied to reduce the  
240 dimensionality of the features to visualize the results.

241

## 242 **Statistical analysis**

243 The receiver operator characteristic (ROC) [19] curve was applied to evaluate the performance  
244 of the two-class classifiers. The sensitivity, specificity and area under the curve (AUC) were  
245 calculated. **AUC could evaluate the classifiers globally and was not sensitive to the ratio**  
246 **of positive and negative samples.** The precision (also known as the positive predictive value),  
247 recall (also known as sensitivity), and F1 score were used to evaluate the performance of the  
248 three-class classifier. The statistics mentioned above were defined as follows:

$$249 \quad \text{Recall} = \frac{\text{True Positives}}{\text{Positives}} = \frac{\text{True Positives}}{\text{Ture Positives+False Negatives'}}$$

$$250 \quad \text{Specificity} = \frac{\text{True Negatives}}{\text{Negatives}},$$

$$251 \quad \text{Precision} = \frac{\text{True Positives}}{\text{True Positives + False Positives'}}$$

$$252 \quad \text{F1-score} = \frac{2\text{Precision}\times\text{Recall}}{\text{Precision} + \text{Recall}}$$

253 All statistical analyses were performed by R 3.5.1 software and Python 3.7.

254

## 255 **Results**

### 256 **Basic characteristics of the study population**

257 For the 6399 patients enrolled in this study, the male-to-female ratios were 1.12:1, 2.59:1,

258 and 0.79:1 for patients with UC, CD and ITB, respectively. The mean (standard deviation, SD)  
259 ages of patients with UC, CD and ITB were  $42.94 \pm 13.66$ ,  $36.04 \pm 14.15$ , and  $42.21 \pm 16.10$ ,  
260 respectively. The basic characteristics of the 6399 patients are summarized in Table 1.

261

### 262 **Convergence of the CNN**

263 After approximately 20 iterations, the model reached convergence, which indicated that the  
264 trained CNN model had stabilized. Supplementary Figure 1 shows the convergence of the  
265 loss and the accuracy of the CNN

266

### 267 **Differential diagnosis between UC and CD**

268 The performance of the classifier is presented in Table 2. The sensitivity, specificity and AUC  
269 were 0.890, 0.837 and 0.936, respectively, when using RF. The CNN was not applied. The  
270 features extracted by RF to distinguish UC and CD are summarized in Table 3. Mucosal  
271 damage characteristics including diffuse congestion, ulcers with purulent secretion, loss of  
272 vascular texture, bleeding tendency after touch, and location of lesions for example rectum  
273 involved, without lumen stenosis, without involvement of ileocecal valve are extracted features  
274 indicating UC.

275

### 276 **Differential diagnosis between UC and ITB**

277 The performance of the classifier is reported in Table 2. The sensitivity, specificity and AUC

278 were 0.833, 0.818 and 0.892, respectively, when using RF. The CNN was not applied. The  
279 features extracted by RF are summarized in Table 4. Undamaged terminal ileum mucosa,  
280 pseudopolyps or neoplasm, undamaged haustral pattern, and involvement of ileocecal valve  
281 are extracted major features indicating ITB.

282

### 283 **Differential diagnosis between CD and ITB**

284 The performance of the classifier is reported in Table 2. The sensitivity, specificity and AUC  
285 were 0.723, 0.772 and 0.816 when using RF and 0.900, 0.771 and 0.910 when using the CNN,  
286 respectively. The features extracted by RF are summarized in Table 5. Involvement of ileocecal  
287 valve and undamaged ileum are extracted as top features for differentiating ITB from CD.  
288 Features including lumen stenosis, mucosal congestion, loss of vascular texture and ulcers  
289 covered with white exudates indicate CD.

290

### 291 **Visualization of the classification**

292 To present the classification results more intuitively, the results are visualized in Figure 2 by  
293 the t-SNE algorithm. Note that the X and Y axes are the result of dimension reduction of  
294 features and bear no real world meaning. Each dot represents an endoscopic description of  
295 a patient. The distance between two dots is proportional to the dissimilarity of the extracted  
296 features of two samples. The CNN almost grouped all same diseases together, while RF  
297 formed many small clusters and many of which had mixed colors. This shows that the CNN  
298 could automatically extract features that were highly capable at characterizing the three

299 differential diagnoses. Further research on interpreting the feature representations from CNN  
300 is warranted.

301

## 302 **Discussion**

303 To the best of our knowledge, this is the first study to apply RF and the CNN to realize  
304 automatic classification of IBD and ITB based on endoscopic results in the form of free text.

305 Classifiers built by both RF and the CNN had extremely high sensitivity and specificity for two-  
306 class classification problems UC and CD and UC and ITB and the three-class classification  
307 problem (UC, CD and ITB), which indicates that UC can be potentially automatically diagnosed  
308 based on a simple objective description, while differential diagnosis of CD and ITB can be  
309 strongly indicated by the classifier.

310

311 Diffuse mucosal damage, bleeding tendency and ulcer with purulent secretion were major  
312 features that led to a classification of UC. Location of the lesion was also an indicator. However,  
313 continuous lesions, which were always mentioned in diagnoses in other studies, were not  
314 clearly proposed. The reason could be that mucosal diffuse damage covers the meaning of  
315 continuous lesions; thus, continuous lesions were not considered an independent feature. In  
316 general, the features automatically extracted by RF were similar to those reported in previous  
317 studies and guidelines [3-5].

318

319 Differential diagnosis between CD and ITB has always been clinically challenging. A study  
320 showed that approximately 65% of patients with CD had been incorrectly diagnosed as having  
321 ITB at least once in China [9]. Another study found that approximately 40% of CD patients had  
322 received a trial anti-TB treatment because of misdiagnosis [20]. Although histologic or  
323 pathologic features including caseating granuloma, isolation of *M. tuberculosis* or a positive  
324 acid-fast staining result can provide extremely strong evidence for TB diagnosis, previous  
325 studies pointed out that these examinations were positive in fewer than 50% patients [9].

326

327 Studies focused on endoscopic differential diagnosis have raised several differential diagnosis  
328 points, including the number of involved colon segments, ileocecal valves, direction of ulcers,  
329 special appearance (cobblestone, scars and pseudopolyps) and location of lesions [1].  
330 However, the weights and potential combinations of these points have not been clarified. In  
331 addition, these clinical features for differential diagnosis often exist in an individual patient  
332 simultaneously. The two reasons above can make doctors, especially those with limited  
333 experience, have great difficulty in making an accurate endoscopic diagnosis. The classifiers  
334 built in this research address this gap by assigning weights to each feature, allowing doctors  
335 to simply record the endoscopic image faithfully and objectively and could result in high  
336 precision and a high recall rate, which would benefit clinical work.

337

338 In addition, the classifiers built by RF found some features that were new or rarely noted  
339 before. Ulcers covered with white exudate were considered a sign of CD instead of ITB. The

340 white exudate could be due to an inflammatory response. Although we cannot claim the  
341 diagnostic value of this point immediately, we hope that this finding can lead to further  
342 research. Mucosal congestion and loss of vascular textures used to be considered features of  
343 UC. However, the classifier considered these features signs of CD when performing differential  
344 diagnosis between CD and ITB, which might provide some new insights. Other features  
345 extracted by RF are reported in Table 5; these were mentioned in previous studies or  
346 guidelines [3-5].

347

348 The detailed features of ulcers (longitudinal, transverse, cobblestone appearance, etc.) were  
349 not considered important features for differential diagnosis of CD and ITB by the classifier in  
350 this study. This is different from the findings of previous studies. We suppose that the reason  
351 may be that the patients with typical features (e.g., cobblestone appearance of CD) comprised  
352 only a very small portion of all patients enrolled in the research. These signs might be  
353 important, but the statistical power is not strong enough to detect this difference. In other  
354 words, examination of these features has high specificity but low sensitivity. Thus, these  
355 features are not proper for differential diagnosis of all patients.

356

357 From the points above, features extracted by RF could also provide a new view for doctors  
358 to make differential diagnosis. Classifiers built by RF could find and assign different weights  
359 to each feature. Besides, these classifiers could also combine the features together to  
360 provide potential extra information. The two points above might help doctors to make

361 better differential diagnosis. **In addition, Interpreting the RF results in comprehensive**  
362 **manner may help encourage the doctors and clinicians to use machine learning tools**  
363 **in healthcare field.**

364

365 The CNN could not be applied to the two-class classifiers including UC in this study because  
366 of extremely unbalanced sample sizes (UC: CD=5128:875; UC: ITB=5128:396), although CNNs  
367 often perform better in classification problems. This phenomenon is caused by the different  
368 morbidity associated with different types of inflammatory intestinal diseases and is a  
369 consequence of consecutive enrollment of patients. On the other hand, the CNN performed  
370 better than RF in the classification of CD and ITB. However, the weakness of the CNN was its  
371 interpretability, which meant that we could not figure out the classification features used by  
372 the CNN. Therefore, from the perspective of a physician, the result given by the CNN can be  
373 a reference but without supporting details; thus, it is difficult to accept the outcome of this  
374 approach as the final clinical diagnosis. However, further clinical trials might provide further  
375 evidence, and the CNN model could be used to assist physicians. Large-cohort studies may  
376 help improve the credibility of the CNN as well.

377

378 Regarding the three-class classification problem, RF could also extract some features.  
379 However, the attribution of each feature could be difficult to figure out because a single  
380 'either-or' feature such as damaged/undamaged mucosa was not sufficient for three-class  
381 classification. Linear combinations of extracted features might be used with variable

382 coefficients. The complexity of the three-class classifier might have caused its low  
383 performance. Although the CNN performed better than RF in the three-class classification  
384 problem according to Table 2 and Figure 2, the interpretability of the CNN will always be a  
385 barrier that is difficult to remove. However, the precision and recall rate of UC were extremely  
386 high, even when using the three-class classifier built by RF; thus, only patients who are  
387 suspected of having CD or ITB would need further classification. The two-class classifier of CD  
388 and ITB discussed above could then be applied.

389

390 **Due to the similarity of the endoscopic image of CD and ITB, some cases might**  
391 **could not be differentiated purely by endoscope and needed further examinations. The**  
392 **upper limit value of the proportion that could be differentiated of ITB and CD depended**  
393 **on both experience of doctors and the power of endoscopic diagnosis criteria. The**  
394 **classifiers built in this research could help unexperienced doctors to make more accuracy**  
395 **diagnosis while the RF classifiers could find some new features for diagnosis, thus might**  
396 **potentially improve the power of diagnosis criteria.**

397

398 There are several limitations of this study. First, Behcet's disease was not taken into  
399 consideration because it is often considered a systemic disease and with variable initial  
400 involvement, including genital aphthae, gastrointestinal involvement, skin lesions, vascular  
401 disease, neurologic disease, and arthritis [20], and it was not appropriate for inclusion in this  
402 study. Second, although we assumed the existence of a complete terminology, subjectivity

403 might still cause bias when endoscopists write the descriptions. In addition, endoscopy  
404 doctors were not completely blind to other clinical examination results, which could cause  
405 subjectivity bias. We have used the first endoscopic description, which had least additional  
406 information, to minimize these bias. However, to eliminate these bias, classification algorithms  
407 directly based on endoscopic images might be preferred, which requires much more  
408 computational resource. We are looking forward to further researches that cooperated by  
409 clinical doctors and technicians. Third, the numbers of patients with UC, CD and ITB involved  
410 in this research were not balanced (5128 UC, 875 CD and 396 ITB). This was mainly because  
411 of the different prevalences and morbidities of these three diseases. A study with  
412 nonconsecutive patients might have a more balanced sample but would have potential  
413 selection bias. Lastly, these research was conducted in Chinese. However, the NLP models  
414 could be easily applied on other languages.

415

## 416 **Conclusion**

417 This was the first study that applied RF and a CNN to realize the automatic classification of  
418 known inflammatory intestinal diseases. The performance of the classifiers was reasonable. In  
419 addition, RF found some new features that were potentially valuable for further research.  
420 Artificial intelligence through machine learning is very promising in helping unexperienced  
421 endoscopists differentiate inflammatory intestinal diseases.

422

423

424 **Ethics approval and consent to participate**

425 This research was approved by the Ethics Committee of Peking Union Medical College  
426 Hospital on September 31<sup>st</sup>, 2018 (IRB S-K894).

427 We hereby state that the study protocol conforms to the ethical guidelines of the 1975  
428 Declaration of Helsinki as reflected in a priori approval by the institution's human research  
429 committee.

430 The written/informed consent was exempted by the ethical committee of PUMCH because  
431 this was a retrospective study for about 10 years with little patients' personal or sensitive  
432 information.

433

434 **Consent for publication**

435 Not Applicable.

436

437 **Availability of data and materials**

438 The datasets of involved patients' endoscopic descriptions analyzed during the  
439 current study are not publicly available due to patient privacy and the requirement  
440 of the Ethics Committee of Peking Union Medical College Hospital. If there are any  
441 need, please contact the corresponding authors.

442

443 **Competing Interests**

444 The authors claim no conflict of interest.

445

446 **Funding**

447 This work was supported by CAMS Innovation Fund for Medical Sciences (No, 2017-I2M-3-  
448 017) , National Natural Science Foundation of China (No. 11801301) and Beijing Natural  
449 Science Foundation (No. Z190024).

450

451 **Authors' contributions**

452 Yuanren Tong and Yucong Lin designed the study. Yingyun Yang, Ji Li, Dong Wu, Aiming  
453 Yang, Jiaming Qian and Yue Li collected the data. Yuanren Tong cleaned the data. Keming Lu  
454 developed the main part of the classifiers. Sheng Yu modified the classifiers. Yuanren Tong  
455 and Yue Li conducted data interpretation. Yuanren Tong and Keming Lu wrote the manuscript.  
456 All authors contributed to editing and revising the manuscript critically. Yue Li and Sheng Yu  
457 directed the direction of the study.

458 **Acknowledgement**

459 We sincerely thank Dr. Wanying Zhang and Dr. Simon Ma for their suggestions towards the  
460 study.

461

462

- 464 1. Ji ML, Lee KM: **Endoscopic Diagnosis and Differentiation of Inflammatory Bowel**  
465 **Disease.** *Clinical Endoscopy* 2016, **49**(4):370-375.
- 466 2. Organization GWH: **Global tuberculosis report 2018.** 2018.
- 467 3. Choi CH, Jung SA, Lee BI, Lee KM, Kim JS, Han DS: **Diagnostic Guideline of Ulcerative**  
468 **Colitis.** *Korean J Gastroenterol* 2009, **53**(3):145-160.
- 469 4. Kim YS, Kim YH, Lee KM, Kim JS, Park YS: **Diagnostic Guideline of Intestinal Tuberculosis.**  
470 *Korean J Gastroenterol* 2009, **53**(3):177-186.
- 471 5. Ye BD, Jang BI, Jeon YT, Lee KM, Kim JS, Yang SK: **[Diagnostic guideline of Crohn's**  
472 **disease].** *Korean J Gastroenterol* 2009, **53**(3):161-176.
- 473 6. Li Y, Zhang LF, Liu XQ, Wang L, Wang X, Wang J, Qian JM: **The role of in vitro**  
474 **interferony-release assay in differentiating intestinal tuberculosis from Crohn's**  
475 **disease in China.** *Journal of Crohns & Colitis* 2012, **6**(3):317-323.
- 476 7. Ren M, Wang-Di L, Yao H, Chun-Hui O, Zhen-Hua Z, Chen Y, Shun-Hua L, Yu-Jun C, Zi-  
477 Ping L, Xiao-Ping W: **Computed tomographic enterography adds value to**  
478 **colonoscopy in differentiating Crohn's disease from intestinal tuberculosis: a**  
479 **potential diagnostic algorithm.** *Endoscopy* 2015, **47**(04):322-329.
- 480 8. He Y, Zhu Z, Chen Y, Chen F, Wang Y, Ouyang C, Yang H, Huang M, Zhuang X, Mao R *et*  
481 *al*: **Development and Validation of a Novel Diagnostic Nomogram to Differentiate**  
482 **Between Intestinal Tuberculosis and Crohn's Disease: A 6-year Prospective**  
483 **Multicenter Study.** *Am J Gastroenterol* 2019, **114**(3):490-499.
- 484 9. Lee YJ, S-K Y, J-S B, S-J M, H-S C, S-S H, K-J K, Lee GH, H-Y J, W-S H: **Analysis of**  
485 **colonoscopic findings in the differential diagnosis between intestinal tuberculosis**  
486 **and Crohn's disease.** *Endoscopy* 2006, **38**(06):592-597.
- 487 10. Delvaux M, Korman LY, Armengolmiro JR, Crespi M, Cass O, Hagenm 眉 ller F, Zwiebel  
488 FM: **The minimal standard terminology for digestive endoscopy: introduction to**  
489 **structured reporting.** *International Journal of Medical Informatics* 1998, **48**(1 欵?):217-  
490 225.
- 491 11. Groenen MJM, Hirs W, Becker H, Kuipers EJ, Henegouwen GPVB, Fockens P, Ouwendijk  
492 RJT: **Gastrointestinal Endoscopic Terminology Coding (GET-C): A WHO-Approved**  
493 **Extension of the ICD-10.** *Digestive Diseases & Sciences* 2007, **52**(4):1004-1008.
- 494 12. Haubrich WS: **Terminology Committee of the World Society of Digestive**  
495 **Endoscopy/OMEDZden 驛 kMaratkaTerminology, Definitions and Diagnostic Criteria**  
496 **in Digestive Endoscopy1984R 枚 hm Pharma GmbHWeiterstadt74.** *Gastrointestinal*  
497 *Endoscopy* 1985, **31**(3):231-232.
- 498 13. Lee DD, Seung HS: **Learning the parts of objects by non-negative matrix factorization.**  
499 *Nature* 1999, **401**(6755):788-791.
- 500 14. Pennington J, Socher R, Manning C: **Glove: Global Vectors for Word Representation.** In:  
501 *Proceedings of the 2014 Conference on Empirical Methods in Natural Language*  
502 *Processing (EMNLP): 2014*; 2014.
- 503 15. Kim Y: **Convolutional Neural Networks for Sentence Classification.** In: *Proceedings of*  
504 *the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP): oct*  
505 *2014; Doha, Qatar.* Association for Computational Linguistics; 2014: 1746-1751.
- 506 16. Kingma DP, Ba J: **Adam: A Method for Stochastic Optimization.** In: *arXiv e-prints*. 2014.

- 507 17. Lin T-Y, Goyal P, Girshick R, He K, Dollár P: **Focal Loss for Dense Object Detection**. In:  
508 *arXiv e-prints*. 2017.
- 509 18. Geoffrey Hinton PK: **Visualizing Data using t-SNE** Laurens van der Maaten MICC-**IKAT**.  
510 2004.
- 511 19. Fawcett T: **An introduction to ROC analysis**. *Pattern Recognition Letters* 2005, **27**(8):861-  
512 874.
- 513 20. Lee S, Kim B, Ti, Kim W: **Differential diagnosis of intestinal Behcet's disease and Crohn's**  
514 **disease by colonoscopic findings**. *Journal of Gastroenterology & Hepatology* 2009,  
515 **41**(01):9-16.

## 516 **Table and figure legends**

517 Table 1. Basic characteristics of the 6399 patients enrolled in the study

518 Table 2. Performances of Classifiers

519 Table 3. Differential features of UC and CD

520 Table 4. Differential features of UC and ITB

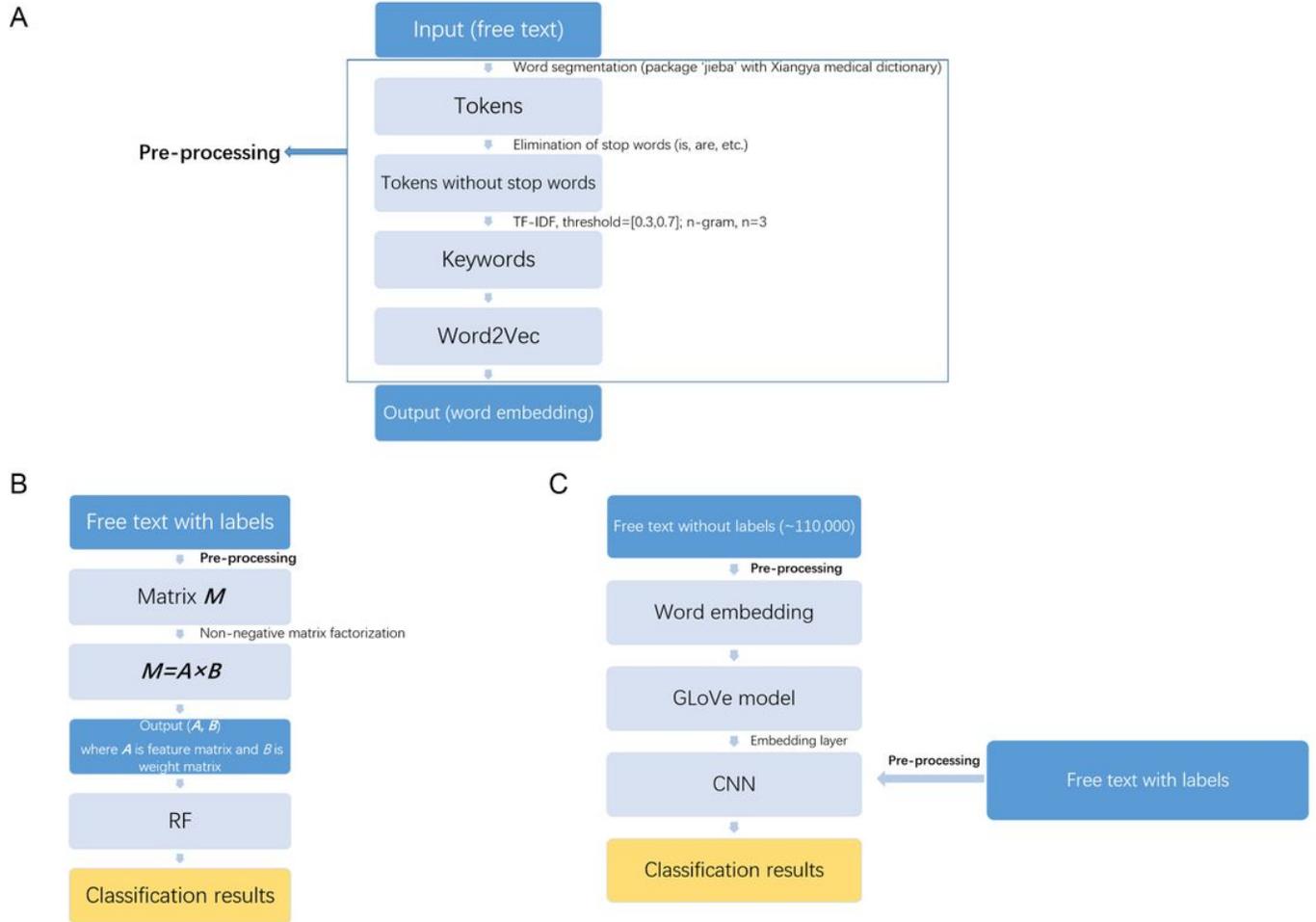
521 Table 5. Differential features of CD and ITB

522 Figure 1. The flow chart of data processing. A: the pre-processing which would be used  
523 later in RF and CNN; B: flow chart of RF; C: flow chart of CNN

524 Figure 2. The visualization of the result of two classifiers. A: the classifier built by RF; B: the  
525 classifier built by CNN. The distance between two dots was proportional to the similarities of  
526 the corresponding endoscopic descriptions. Different color of the dots represented different  
527 diseases. Purple: Ulcerative colitis (UC); yellow: Crohn's disease (CD); green, Intestinal  
528 tuberculosis (ITB).

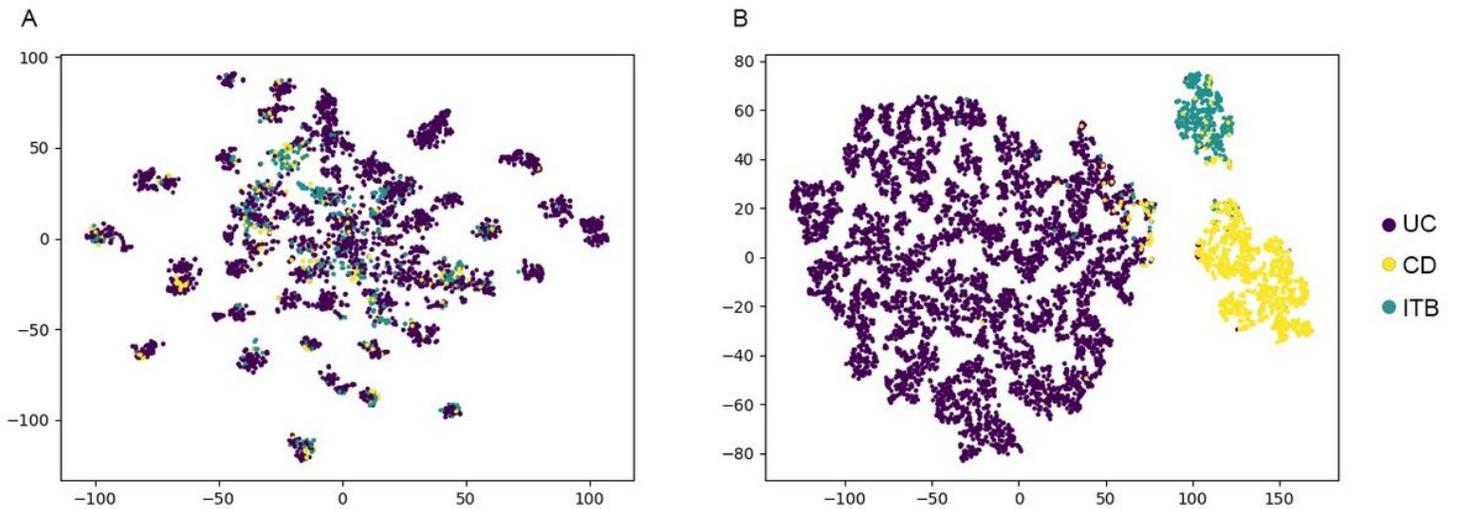
529

# Figures



**Figure 1**

The flow chart of data processing. A: the pre-processing which would be used later in RF and CNN; B: flow chart of RF; C: flow chart of CNN



## Figure 2

The visualization of the result of two classifiers. A: the classifier built by RF; B: the classifier built by CNN. The distance between two dots was proportional to the similarities of the corresponding endoscopic descriptions. Different color of the dots represented different diseases. Purple: Ulcerative colitis (UC); yellow: Crohn's disease (CD); green, Intestinal tuberculosis (ITB).

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementarymaterialsibdr2.pdf](#)