



# Radiomics combined with clinical features in distinguishing non-calcifying tuberculosis granuloma and lung adenocarcinoma in small pulmonary nodules

Qing Dong<sup>1</sup>, Qingqing Wen<sup>2</sup>, Nan Li<sup>3</sup>, Jinlong Tong<sup>4</sup>, Zhaofu Li<sup>5</sup>, Xin Bao<sup>6</sup>, Jinzhi Xu<sup>1</sup> and Dandan Li<sup>7</sup>

<sup>1</sup>Department of Thoracic Surgery at No. 4 Affiliated Hospital, Harbin Medical University, Harbin, China

<sup>2</sup>Icahn School of Medicine at Mount Sinai, New York, NY, United States of America

<sup>3</sup>Department of Pathology at No. 4 Affiliated Hospital, Harbin Medical University, Harbin, China

<sup>4</sup>Department of Medical Imaging at No. 4 Affiliated Hospital, Harbin Medical University, Harbin, China

<sup>5</sup>Heilongjiang Institute of Automation, Harbin, China

<sup>6</sup>Harbin Medtech Innovative Company, Harbin, China

<sup>7</sup>Department of Radiology at Cancer Hospital, Harbin Medical University, Harbin, China

## ABSTRACT

**Aim.** To evaluate the performance of radiomics models with the combination of clinical features in distinguishing non-calcified tuberculosis granuloma (TBG) and lung adenocarcinoma (LAC) in small pulmonary nodules.

**Methodology.** We conducted a retrospective analysis of 280 patients with pulmonary nodules confirmed by surgical biopsy from January 2017 to December 2020. Samples were divided into LAC group ( $n = 143$ ) and TBG group ( $n = 137$ ). We assigned them to a training dataset ( $n = 196$ ) and a testing dataset ( $n = 84$ ). Clinical features including gender, age, smoking, CT appearance (size, location, spiculated sign, lobulated shape, vessel convergence, and pleural indentation) were extracted and included in the radiomics models. 3D slicer and FAE software were used to delineate the Region of Interest (ROI) and extract clinical features. The performance of the model was evaluated by the Area Under the Receiver Operating Characteristic (ROC) Curve (AUC).

**Results.** Based on the model selection, clinical features gender, and age in the LAC group and TBG group showed a significant difference in both datasets ( $P < 0.05$ ). CT appearance lobulated shape was also significantly different in the LAC group and TBG group (Training dataset,  $P = 0.034$ ; Testing dataset,  $P = 0.030$ ). AUC were 0.8344 (95% CI [0.7712–0.8872]) and 0.751 (95% CI [0.6382–0.8531]) in training and testing dataset, respectively.

**Conclusion.** With the capacity to detect differences between TBG and LAC based on their clinical features, radiomics models with a combined of clinical features may function as the potential non-invasive tool for distinguishing TBG and LAC in small pulmonary nodules.

Submitted 24 August 2021

Accepted 6 September 2022

Published 19 October 2022

Corresponding author

Dandan Li,

hmu.cancer.hospital@gmail.com

Academic editor

Henkjan Huisman

Additional Information and  
Declarations can be found on  
page 11

DOI 10.7717/peerj.14127

© Copyright

2022 Dong et al.

Distributed under

Creative Commons CC-BY 4.0

OPEN ACCESS

**Subjects** Biotechnology, Oncology, Surgery and Surgical Specialties, Data Mining and Machine Learning, Tuberculosis

**Keywords** Radiomics, Non-calcified tuberculosis granuloma, Lung adenocarcinoma, Pulmonary nodules, Clinical features

## INTRODUCTION

Tuberculosis (TB) is an infectious disease that is caused by a single source (*MacNeil et al., 2020*). According to statistics, there are about 10 million new TB patients and 1.5 million deaths each year, more than any other infectious disease (*Thwaites & Nahid, 2020*). Among them, pulmonary TB is the most common, accounting for about 85% of all tuberculosis cases (*Reid et al., 2019*). Its pathological manifestation is chronic granulomatous inflammation (*Yuan & Sampson, 2018*). In 2020, 1,930 million new cancer cases and 10 million deaths were estimated worldwide, with approximately 2.2 million (11.4%) new lung cancer cases and 1.8 million (18%) deaths (*Sung et al., 2021*). LAC is the most common malignant tumor, its prognosis is much worse than tuberculosis, so early diagnosis and treatment are very important. However, it is difficult to distinguish TBG and LAC in chest images, and even nuclear medicine is nonspecific (*Fischer, Lassen & Højgaard, 2011; McWilliams et al., 2013*). Because both diseases can be shown as solid nodules or masses on imaging studies and have similar radiological features. The confirmative diagnosis of pulmonary nodules is usually biopsy or surgery (*Siegel, Miller & Jemal, 2019*). However, this invasive examination may lead to possible tissue damage (*Pisano et al., 2020*). Besides, unnecessary imaging studies may also delay treatment, or miss the best treatment time window (*Huo et al., 2019*). Therefore, it is expected in clinical practice that a method can be used to monitor pulmonary nodules noninvasively, and may also provide effective support for the diagnosis and treatment of pulmonary nodules. Radiomics is used to extract features from radiological images and make these features in a quantifiable manner. Its purpose is to better or more consistently discover radiological features, and provide objective features that cannot be provided by standard visual image interpretation for quantitative and qualitative density and morphological characteristics of pulmonary nodules (*Bi et al., 2019; Peikert, Bartholmai & Maldonado, 2020*). Radiomics can be used for auxiliary diagnosis of pulmonary nodules and prognosis prediction of lung cancer (*Mu et al., 2020; Hosny et al., 2018*). Importantly, radiomics has been applied to evaluate the molecular and clinical features of lung cancer because of its capacity of detecting atypical features in tumor lesions (*Grossmann et al., 2017*). In this study, we hypothesized that radiomics analysis could distinguish TBG and LAC in small pulmonary nodules based on imaging and clinical features. To test this idea, we extracted the features of small nodules from lung CT using radiomics technology, obtained the radiological model through statistical analysis, and combined it with clinical features. Our goal is to develop a non-invasive method of distinguishing benign and malignant pulmonary nodules using radiomics models in a combination of clinical features.

## MATERIALS & METHODS

### Patients selection

Our research had been approved by the Ethics Review Committee of No.4th Affiliated Hospital of Harbin Medical University (Institutional Review Board that approved number: KY2020-04). Since it was a retrospective study, additional informed consent was waived. Samples that meet all the following criteria were included: (1) Pulmonary tuberculosis or primary LAC confirmed by biopsy or surgical pathology. (2) Enhanced chest CT images that were collected within 1 month before surgery. (3) Isolated non-calcified pulmonary nodules. (4) The maximum diameter was less than 30 mm. Samples were excluded if they did not meet the above criteria. According to the above inclusion and exclusion criteria, we enrolled 280 patients (143 LAC, 137 TBG) who met the inclusion criteria from January 2017 to December 2020. Patients were randomly selected into training and testing data sets by FeAture Explorer (FAE) software based on the TBG or LAC group.

### Evaluation of pathology

All specimens were fixed with formalin and stained with hematoxylin and eosin (HE). In order to judge the biopsy results separately, two pathologists with more than 10 years of working experience were blind to the clinical information. All lesions were classified according to the international standard (*Rami-Porta et al., 2017*). Classification of Pulmonary Adenocarcinoma according to the latest IASLC/ATS/ERS criteria in previous study (*Eguchi et al., 2014*): (1). Preinvasive lesions (2). Minimally invasive adenocarcinoma ( $\leq 3$  cm lepidic predominant tumor with  $\leq 5$  mm invasion) (3) Invasive adenocarcinoma (4) Variants of invasive adenocarcinoma

### CT data collection

Scanning parameters: The second generation gemstone spectral CT (Discovery CT750 HD) of the US General Electric Company was used to perform dual-phase enhanced CT examination of 280 patients. Patients were in the supine position, scan range was from chest entrance to the diaphragm, to ensure full coverage of all lung tissue. A total of 75 mL non-ionic iodine contrast agent Ioversol (350 mgI/ml) was injected with a double-tube high-pressure syringe at a flow rate of 3.5 mL/s. After injection into the elbow vein, the thoracic aorta at the level of tracheal protuberance was automatically selected as the starting point for monitoring. The intelligent tracking technology of the contrast agent was used to determine the starting time of scanning. When the threshold reached 130 Hu, the scanning was automatically triggered. A venous phase scan started at 80 s. Other parameters were as follows: layer thickness was 0.625 mm, frame rotation time was 0.6 s, pitch was 1.375, and tube current was 600 mA.

### Image evaluation

The CT appearance including lesion size, location, burr, lobulation, vascular penetration, and pleural involvement was extracted by two radiologists with more than 10 years of imaging diagnosis experience. Other clinical features such as age, gender, and smoking history were obtained from the electronic health records. To keep a subjective clinical

judgment, the two radiologists were blind to both baseline information and biopsy results. If there were conflicting opinions, an agreement would be achieved after discussion. For example, an average value of lesion size was taken after discussion if there were conflicting opinions between radiologists.

### **Tumor segmentation**

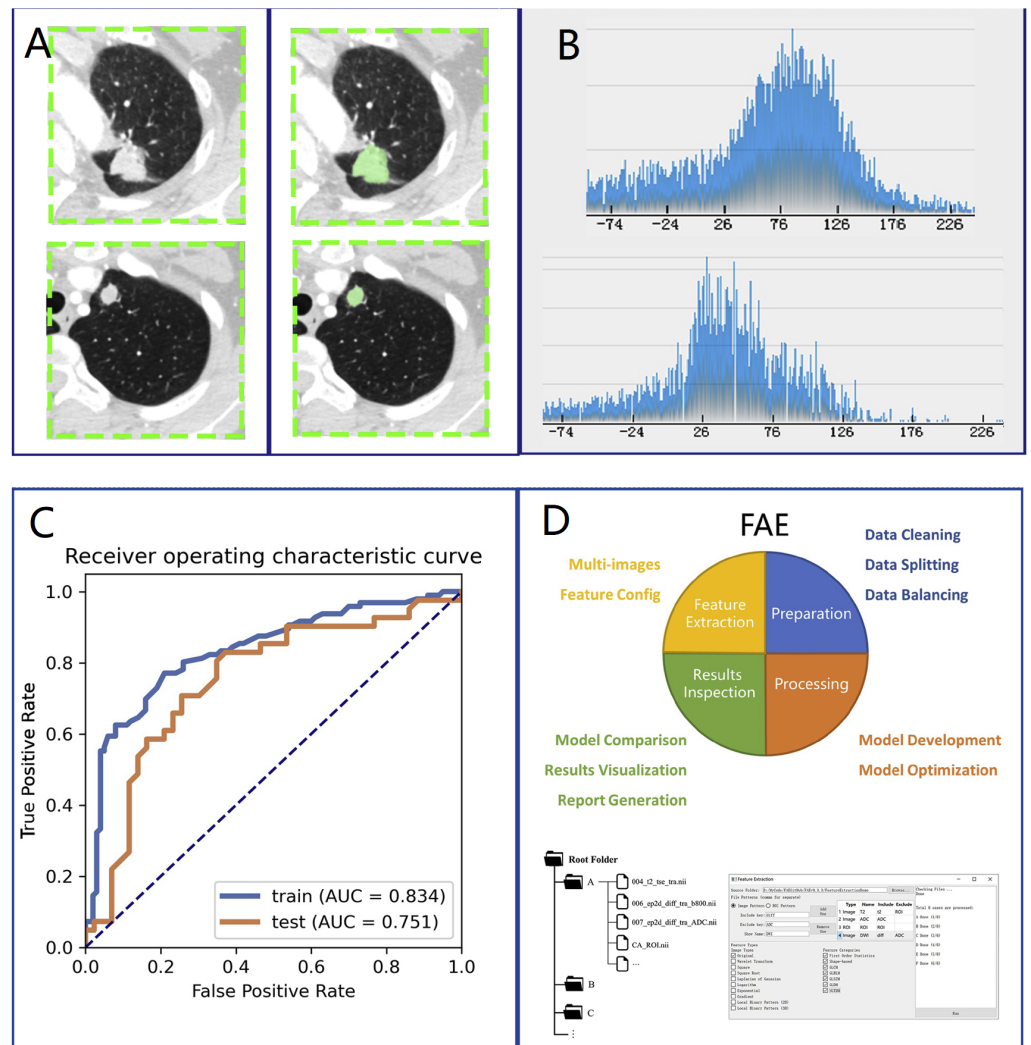
We loaded CT images into 3D slicer software (version 4.10.0) for manual segmentation (Fig. 1A, left side). The region of interest (ROI) on CT was delineated by a thoracic surgeon with 10 years of lung surgery experience (Fig. 1A, right side). The ROI was then confirmed by another senior radiologist with chest radiograph experience for more than 10 years.

### **Radiomics feature extraction and model building**

We selected 196 cases as the training dataset (96/100 = TBG/LAC) and 84 cases as the testing dataset (41/43 = TBG/LAC). 851 radiomics features were extracted from each ROI and divided into three main categories: (1) First-order features. (2) Shape characteristics. (3) Texture features, including gray level co-occurrence matrix (GLCM) features, grey-level run-length matrix (GLRLM) features, gray level size zone matrix (GLSZM) features, neighborhood grey tone difference matrix (NGTDM) features, and grey level dependence matrix (GLDM) features. Figure 2 showed how Grey Level Histogram worked. FAE applied uniformization automatically to the feature matrix when preprocessing CT data, where each feature vector subtracted its average value and then divided by its length. Since the dimensional feature space was very high, the similarity of each feature pair was compared. If the Pearson Correlation Coefficient (PCC) of one feature pair was greater than 0.99, one of them from the pair was removed. After this preprocessing procedure, the size of the feature space was reduced, and each feature was independent of another. Kruskal Wallis was utilized to explore the important features corresponding to labels. In the FAE software, Pearson and Kruskal Wallis methods were automatically selected in the FAE software and we applied them to the training dataset. To evaluate the relationship between features and labels, we calculated the  $F$  value. Afterward, we ranked the top 14 features according to the corresponding  $F$  value. These 14 features were chosen by the FAE software based on the highest  $F$  value. Eventually, Random Forest Model with the highest AUC value was chosen automatically by FAE software as a classifier from all existing models including Support Vector Machine (SVM), Latent Dirichlet Allocation (LDA), Autoencoder (AE), Random Forest, Logistic Regression-Lasso, Adaboost, Decision Tree, Gaussian Process, Naive Bayes. To determine the hyperparameters of the model (e.g., The number of features), we applied 10 times cross-validation on the training dataset. Therefore, hyperparameters were set according to the model performance on the validation dataset (Fig. 1).

### **Statistical analysis**

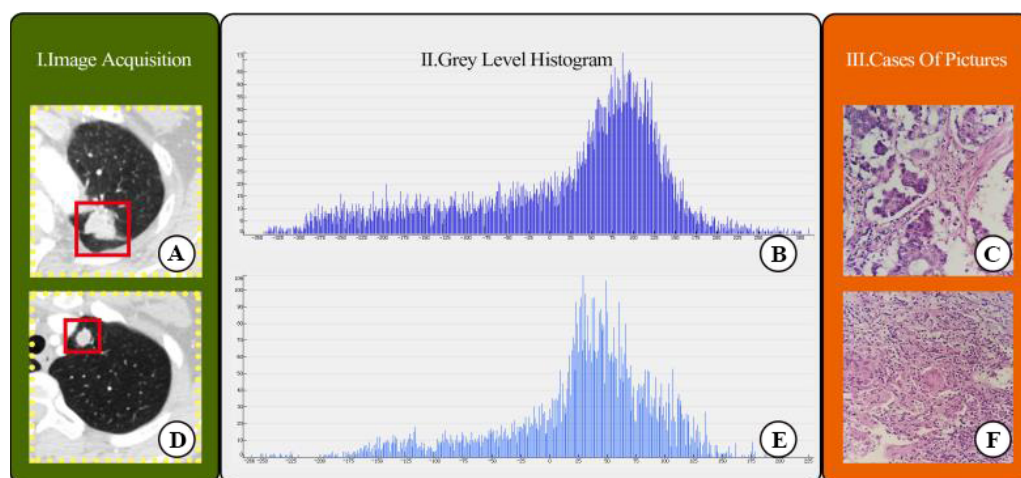
We used the Statistical Program for Social Science (SPSS, version 16.0) to test statistical differences in clinical features between LAC and TBG groups. The independence of categorical variables was examined by the Chi-square test and Fisher exact test. To test the continuous variables with normal distribution, a  $t$ -test was conducted ( $P < 0.05$  indicates statistical significance). We used the Chi-square test for categorical variables



**Figure 1** Research method. Overview of research methods: (A) Collection of chest CT data and ROI delineation. (B) Feature extraction. The image is the gray scale histogram of the lesion. (C) Data analysis. (D) Operation using FAE software.

Full-size DOI: [10.7717/peerj.14127/fig-1](https://doi.org/10.7717/peerj.14127/fig-1)

such as location, smoking, and other clinical features. The performance of the model and quantitative analysis were evaluated by the ROC curve and AUC (Fig. 1C), respectively. Sensitivity, specificity, positive predictive value (PPV), and negative predictive value (NPV) were calculated when the Youden index was maximized to its cut-point value. We estimated 95% confidence intervals for 1000 samples by bootstrapping. All the above processes were operated *via* FeAture Explorer Pro (FAEPro, V0.3.5, Fig. 1D) on Python (3.7.6) according to the software operation reference related literature (Song *et al.*, 2020).



**Figure 2** CT images showed lung adenocarcinoma (LAC) and non-calcified tuberculous granuloma (TB); (A) and (D) CT scan showed irregular solid nodules (red area) in the left upper lobe; (B) and (E) gray scale histogram of the nodule; (C) LAC with hematoxylin and eosin (H & E) stain,  $\times 400$ ; (F) TB with hematoxylin and eosin (H & E) stain,  $\times 400$ .

Full-size DOI: [10.7717/peerj.14127/fig-2](https://doi.org/10.7717/peerj.14127/fig-2)

## RESULTS

### Clinical features

Table 1 listed the statistical test results in the training dataset and testing dataset. There were 196 patients in the training dataset, including 96 males (age range: 40–79 years old, average age:  $64.53 \pm 9.21$  years old) and 100 females (age range: 33–72 years old, mean age:  $56.06 \pm 10.98$  years). The testing dataset included 84 patients with 42 males (age range: 41–79 years old, average age:  $63.71 \pm 10.22$  years old) and 42 females (age range: 33–73 years, mean age:  $58.65 \pm 10.71$  years). Patients' gender and age were significantly different in the LAC group and TBG group in both datasets (Training dataset, Gender:  $P = 0.001$ , Age:  $P = 0.006$ ; Testing dataset, Gender:  $P = 0.016$ , Age:  $P = 0.005$ ). However, TB and LAC were indistinguishable by some clinical features such as smoking status. For example, there was no statistical difference between smoking history and patients' LAC or TB status (Training dataset,  $P = 0.15$ ; Testing dataset,  $P = 0.536$ ). In CT appearance, the lobulated shape was found to show a significant difference in the LAC group and TBG group in the training dataset ( $P = 0.03$ ) and the testing dataset ( $P = 0.030$ ). The rest CT features did not show any statistical difference in two groups, including size (Training dataset,  $P = 0.60$ ; Testing dataset,  $P = 0.67$ ), location (Training dataset,  $P = 0.910$ ; Testing dataset,  $P = 0.43$ ), spiculated sign (Training dataset,  $P = 0.97$ ; Testing dataset,  $P = 0.79$ ), vessel convergence (Training dataset,  $P = 0.40$ ; Testing dataset,  $P = 0.43$ ), and pleural indentation (Training dataset,  $P = 0.34$ ; Testing dataset,  $P = 0.85$ ). These CT features were not distinguishable between LAC and TB in the model.

**Table 1** Clinical characteristics and CT findings in LAC and TB.

Characteristic	Training data set (n = 196)		P	Test data set (n = 84)		P
	LAC(100)	TB(96)		LAC(43)	TB(41)	
Gender			*0.001			*0.016
Male	37	59		16	26	
Female	63	37		27	15	
Age (mean ± SD, years)	64.53 ± 9.21	56.06 ± 10.98	*0.006	63.71 ± 10.22	58.65 ± 10.71	*0.005
Smoking history			0.148			0.536
Absence	69	75		29	25	
Presence	31	21		14	16	
Size (mean ± SD, mm)	19.81 ± 7.47	18.69 ± 5.44	0.595	20.71 ± 7.62	19.03 ± 9.01	0.667
Location			0.910			0.425
Upper and middle	68	66		28	30	
Lower	32	30		15	11	
Spiculated sign			0.967			0.791
Absence	57	55		25	25	
Presence	43	41		18	16	
Lobulated shape			*0.034			*0.030
Absence	36	49		14	23	
Presence	64	47		29	18	
Vessel convergence			0.400			0.425
Absence	44	48		22	21	
Presence	56	48		21	20	
Pleural indentation			0.337			0.884
Absence	30	35		13	13	
Presence	70	61		30	28	

**Notes.**

The differences were assessed with the Wilcoxon rank sum test or Pearson chi-squared test.

CT, computed tomography; LAC, lung adenocarcinoma; TB, pulmonary tuberculosis; SD, standard deviation.

\*P < 0.05.

## Feature selection and radiological model construction

Table 2 illustrated the prediction performance of the training dataset and testing dataset. The accuracy of the training data set was 0.781, AUC was 0.834 (95% Confidence Interval = 0.7712–0.887), NPV was 0.782, PPV was 0.779, sensitivity was 0.771, and specificity was 0.790. Accuracy of the testing dataset was 0.726, AUC was 0.751 (95% confidence interval = 0.6382–0.853), NPV was 0.794, PPV was 0.680, sensitivity was 0.829, and specificity was 0.628. Table 3 showed features with the 14 highest AUC values on the testing dataset (Table 3 and Fig. 3). In addition, the ROC curve was shown in Fig. 4 (Training dataset AUC = 0.834; Testing dataset AUC = 0.751).

## DISCUSSION

The article discussed a non-invasive diagnostic method for distinguishing non-calcifying tuberculosis granuloma from lung adenocarcinoma. The results of this study showed that age, gender, and lobulation were important predictors for distinguishing the LAC group

**Table 2** Clinical statistics in the diagnosis.

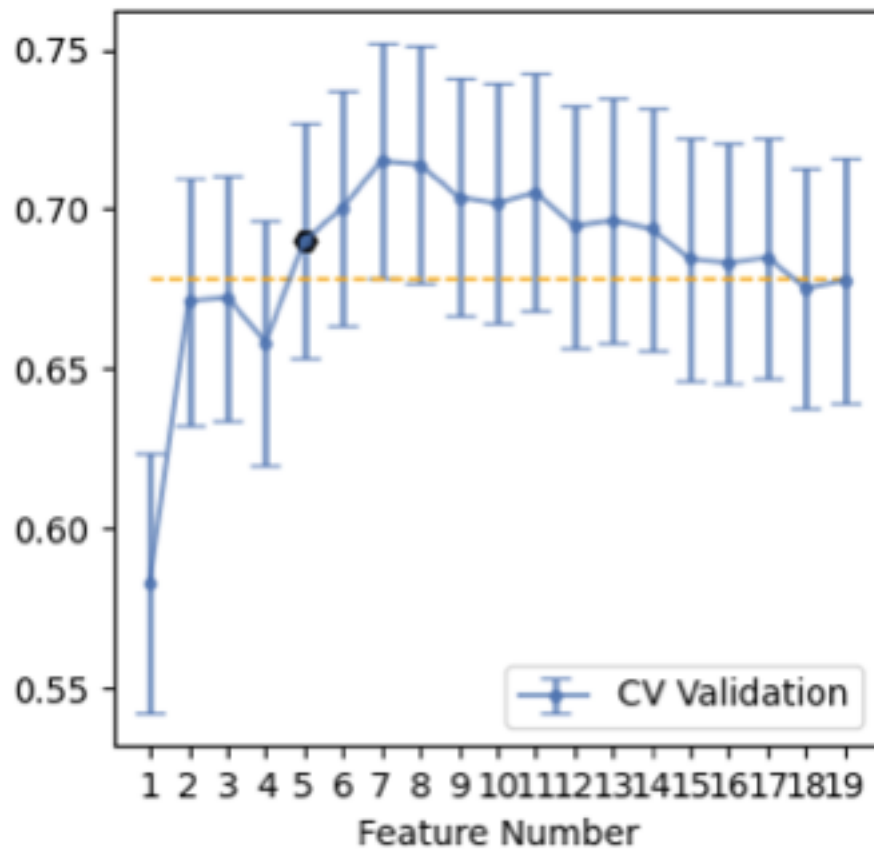
	Accuracy	AUC	AUC 95% CIs	NPV	PPV	Sensitivity	Specificity
Training data set	0.7806	0.8344	0.7712–0.8872	0.7822	0.7789	0.7708	0.7900
Test data set	0.7262	0.751	0.6382–0.8531	0.7941	0.68	0.8293	0.6279

**Table 3** The rank of selected features.

Features	Rank
original_firstorder_90Percentile	1
original_firstorder_Energy	2
original_firstorder_Mean	3
wavelet-HHL_firstorder_Median	4
wavelet-HHL_glcm_ClusterProminence	5
wavelet-HHL_glcm_Imc1	6
wavelet-HHL_glcm_Imc2	7
wavelet-HHL_gldm_DependenceEntropy	8
wavelet-HHL_glrIm_RunEntropy	9
wavelet-HHL_glszm_GrayLevelNonUniformityNormalized	10
wavelet-HHL_glszm_SizeZoneNonUniformityNormalized	11
wavelet-HHL_ngtdm_Busyness	12
wavelet-HHL_ngtdm_Strength	13
wavelet-LLH_glcm_MCC	14

and TB group (Cui et al., 2020). On the one hand, the average age of patients in the TB group was lower than that in the LAC group, which may be explained by the fact that LAC is a malignant tumor, which is common in elderly patients. On the other hand, the number of female patients in the LAC group was more than that in the TB group, whereas the number of male patients in the LAC group was more than that in the TB group. The gender imbalance in the two groups may lead to statistical differences. It could be explained by the fact that females are prone to LAC compared to males, and males are more susceptible to TB compared to females (Marçôa et al., 2018). Radiomics is a process that transforms the subjective evaluation of images into objective quantitative data. Many studies have shown that it can be used as a non-invasive method to predict the benign and malignant effects of pulmonary nodules (Feng et al., 2020b; Xu et al., 2019; Wilson & Devaraj, 2017). These objective data cannot be identified visually but can be determined in a computer-aided manner. The CT appearance ‘lobulated shape’ in this study was statistically different in both groups. This feature can reflect the heterogeneity within pulmonary nodules and help to identify benign and malignant nodules (Jiang et al., 2021). In this study, 196 cases were selected as the training dataset and 84 cases were chosen as the testing dataset. A total of 851 radiomics features were extracted from each ROI randomly and automatically by software, including 18 first-order features, 14 shape features, 24 gray level co-occurrence matrices, 16 gray area size matrices, 16 gray level travel matrices, five domain gray difference matrices, 14 gray level correlation matrices, and 744 wavelet features. The features were sorted

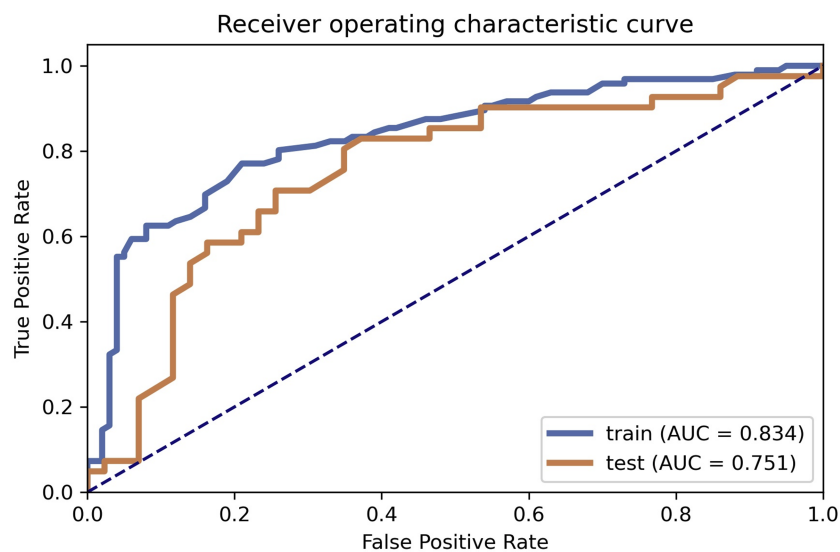




**Figure 3** Features selection. Fourteen features selection (above the yellow line).

Full-size DOI: [10.7717/peerj.14127/fig-3](https://doi.org/10.7717/peerj.14127/fig-3)

according to the corresponding  $F$  value, and the first 14 features are selected according to the verification performance. There are three firstorder features, including original \_ firstorder \_ 90Percentile, original \_ firstorder \_ Energy and original \_ firstorder \_ Mean. The first-order features stand for the difference in the distribution of individual prime parameter values, which reflects the difference in the density of lesions. This is the density difference in internal space between lung adenocarcinoma and non-calcified granuloma, which is difficult to identify from the eyes since it is a high-dimensional spatial feature. These features are related to gray matrix parameters. This indicates that the change of gray level in CT images of lung lesions may potentially contribute to the differential diagnosis of lung adenocarcinoma and non-calcified granuloma (Cui et al., 2020). Random forest was used as a classifier in the model because of its highest AUC value among all models. Lung cancer and granuloma were commonly found in the upper lobe in this study. This may be due to changes in lobulation caused by lung cancer infiltration. However, chronic inflammation may also have similar characteristics. This could explain the reason for the relatively low AUC in the results. The AUC of the training dataset and the testing dataset were 0.834 and 0.751, respectively. The AUC of the training dataset is 0.834 compared to 0.751 in the AUC of the testing dataset. The NPV, PPV, sensitivity, and specificity have high similarities when



**Figure 4** ROC curve selection. ROC curve of the model.

Full-size DOI: 10.7717/peerj.14127/fig-4

compared to previous studies of its kind (Feng et al., 2020a; Chen et al., 2020). It may not be appropriate to observe lung cancer for a long time without providing treatment, but suspected nodules that grow slowly are not easily identifiable with imaging studies without a sufficient waiting period. In addition, lung cancer and granuloma cannot be accurately distinguished in PET scans as well (Du et al., 2021). Although the gold standard for lung cancer diagnosis is the surgical biopsy, it is considered overtreatment if the nodule is a granuloma. On the contrary, conservative treatment may delay the timely treatment for lung cancer. Overall, it is difficult to distinguish benign and malignant pulmonary nodules merely using a lung CT scan. Physicians have been seeking a non-invasive examination to solve this problem. Radiomics, in combination with clinical features, shows its potential to be used as an effective tool to assist radiologists to distinguish benign and malignant pulmonary nodules. However, we have several limitations in this study. Firstly, it was a retrospective analysis. The sample size was relatively small and selection bias could be a potential issue. More high-quality samples are needed to prove the validity of the study in the future. Secondly, selected patients who had surgeries were more likely to be patients diagnosed with malignant tumors. Future research should maintain a relatively equal number of pathology results in both LAC and TB groups. Thirdly, different CT scans may affect the quality of image parameters. Therefore, thin-layer CT scanning (with a value of 0.625 mm) was adopted, and Radiomics normalization preprocessing was used to improve the quality of the data.

## CONCLUSIONS

In summary, radiomics combined with clinical features is a possible non-invasive tool to distinguish non-calcifying tuberculosis granuloma and lung adenocarcinoma in small

pulmonary nodules. The application of this combination has a great potential to decrease overdiagnosis and overtreatment in the future.

## ADDITIONAL INFORMATION AND DECLARATIONS

### Funding

This work was supported by the National Natural Science Foundation of China (No. 61263033), the International Science and Technology Cooperation Project of Hainan (No. KJHZ2015-4), the Higher School Scientific Research Project of Hainan Province (No. Hnky2015-80), and the Clinical Study on the Changes of Brain Net Efficiency after Preventive Brain Irradiation in Patients with Limited Stage Small Cell Lung Cancer. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

### Grant Disclosures

The following grant information was disclosed by the authors:

National Natural Science Foundation of China: 61263033.

International Science and Technology Cooperation Project of Hainan: KJHZ2015-4.

Higher School Scientific Research Project of Hainan Province: Hnky2015-80.

Clinical Study on the Changes of Brain Net Efficiency after Preventive Brain Irradiation in Patients with Limited Stage Small Cell Lung Cancer.

### Competing Interests

Xin Bao is employed by Harbin Medtech Innovative Company.

### Author Contributions

- Qing Dong conceived and designed the experiments, authored or reviewed drafts of the article, and approved the final draft.
- Qingqing Wen performed the experiments, authored or reviewed drafts of the article, and approved the final draft.
- Nan Li conceived and designed the experiments, prepared figures and/or tables, authored or reviewed drafts of the article, and approved the final draft.
- Jinlong Tong performed the experiments, prepared figures and/or tables, and approved the final draft.
- Zhaofu Li performed the experiments, authored or reviewed drafts of the article, and approved the final draft.
- Xin Bao analyzed the data, prepared figures and/or tables, and approved the final draft.
- Jinzhi Xu analyzed the data, authored or reviewed drafts of the article, and approved the final draft.
- Dandan Li conceived and designed the experiments, prepared figures and/or tables, and approved the final draft.

### Human Ethics

The following information was supplied relating to ethical approvals (*i.e.*, approving body and any reference numbers):

No. 4th Affiliated Hospital of Harbin Medical University granted Ethical approval to carry out the study within its facilities (KY2020-04).

### Data Availability

The following information was supplied regarding data availability:

The raw data is available as [Supplemental File](#).

### Supplemental Information

Supplemental information for this article can be found online at <http://dx.doi.org/10.7717/peerj.14127#supplemental-information>.

## REFERENCES

- Bi WL, Hosny A, Schabath MB, Giger ML, Birkbak NJ, Mehrtash A, Allison T, Arnaout O, Abbosh C, Dunn IF, Mak RH, Tamimi RM, Tempany CM, Swanton C, Hoffmann U, Schwartz LH, Gillies RJ, Huang RY, Aerts H. 2019. Artificial intelligence in cancer imaging: clinical challenges and applications. *CA: A Cancer Journal for Clinicians* **69**(2):127–157 DOI [10.3322/caac.21552](https://doi.org/10.3322/caac.21552).
- Chen X, Feng B, Chen Y, Liu K, Li K, Duan X, Hao Y, Cui E, Liu Z, Zhang C, Long W, Liu X. 2020. A CT-based radiomics nomogram for prediction of lung adenocarcinomas and granulomatous lesions in patient with solitary sub-centimeter solid nodules. *Cancer Imaging* **20**(1):45 DOI [10.1186/s40644-020-00320-3](https://doi.org/10.1186/s40644-020-00320-3).
- Cui EN, Yu T, Shang SJ, Wang XY, Jin YL, Dong Y, Zhao H, Luo YH, Jiang XR. 2020. Radiomics model for distinguishing tuberculosis and lung cancer on computed tomography scans. *World Journal of Clinical Cases* **8**(21):5203–5212 DOI [10.12998/wjcc.v8.i21.5203](https://doi.org/10.12998/wjcc.v8.i21.5203).
- Du D, Gu J, Chen X, Lv W, Feng Q, Rahmim A, Wu H, Lu L. 2021. Integration of PET/CT radiomics and semantic features for differentiation between active pulmonary tuberculosis and lung cancer. *Molecular Imaging and Biology* **23**(2):287–298 DOI [10.1007/s11307-020-01550-4](https://doi.org/10.1007/s11307-020-01550-4).
- Eguchi T, Kadota K, Park BJ, Travis WD, Jones DR, Adusumilli PS. 2014. The new IASLC-ATS-ERS lung adenocarcinoma classification: what the surgeon should know. *Seminars in Thoracic and Cardiovascular Surgery* **26**(3):210–222 DOI [10.1053/j.semtcvs.2014.09.002](https://doi.org/10.1053/j.semtcvs.2014.09.002).
- Feng B, Chen X, Chen Y, Liu K, Li K, Liu X, Yao N, Li Z, Li R, Zhang C, Ji J, Long W. 2020a. Radiomics nomogram for preoperative differentiation of lung tuberculoma from adenocarcinoma in solitary pulmonary solid nodule. *European Journal of Radiology* **128**:109022 DOI [10.1016/j.ejrad.2020.109022](https://doi.org/10.1016/j.ejrad.2020.109022).
- Feng B, Chen X, Chen Y, Lu S, Liu K, Li K, Liu Z, Hao Y, Li Z, Zhu Z, Yao N, Liang G, Zhang J, Long W, Liu X. 2020b. Solitary solid pulmonary nodules: a CT-based deep learning nomogram helps differentiate tuberculosis granulomas from lung adenocarcinomas. *European Radiology* **30**(12):6497–6507 DOI [10.1007/s00330-020-07024-z](https://doi.org/10.1007/s00330-020-07024-z).

- Fischer BM, Lassen U, Højgaard L. 2011.** PET-CT in preoperative staging of lung cancer. *The New England Journal of Medicine* **364**(10):980–981 DOI [10.1056/NEJMc1012974](https://doi.org/10.1056/NEJMc1012974).
- Grossmann P, Stringfield O, El-Hachem N, Bui MM, Velazquez ERios, Parmar C, Leijenaar RT, Haibe-Kains B, Lambin P, Gillies RJ, Aerts HJ. 2017.** Defining the biological basis of radiomic phenotypes in lung cancer. *eLife* **6**:e23421 DOI [10.7554/eLife.23421](https://doi.org/10.7554/eLife.23421).
- Hosny A, Parmar C, Coroller TP, Grossmann P, Zeleznik R, Kumar A, Bussink J, Gillies RJ, Mak RH, Aerts H. 2018.** Deep learning for lung cancer prognostication: a retrospective multi-cohort radiomics study. *PLOS Medicine* **15**(11):e1002711 DOI [10.1371/journal.pmed.1002711](https://doi.org/10.1371/journal.pmed.1002711).
- Huo J, Xu Y, Sheu T, Volk RJ, Shih YT. 2019.** Complication rates and downstream medical costs associated with invasive diagnostic procedures for lung abnormalities in the community setting. *JAMA Internal Medicine* **179**(3):324–332 DOI [10.1001/jamainternmed.2018.6277](https://doi.org/10.1001/jamainternmed.2018.6277).
- Jiang Y, Che S, Ma S, Liu X, Guo Y, Liu A, Li G, Li Z. 2021.** Radiomic signature based on CT imaging to distinguish invasive adenocarcinoma from minimally invasive adenocarcinoma in pure ground-glass nodules with pleural contact. *Cancer Imaging* **21**(1):1 DOI [10.1186/s40644-020-00376-1](https://doi.org/10.1186/s40644-020-00376-1).
- MacNeil A, Glaziou P, Sismanidis C, Date A, Maloney S, Floyd K. 2020.** Global epidemiology of tuberculosis and progress toward meeting global targets—worldwide, 2018. *Morbidity and Mortality Weekly Report* **69**(11):281–285 DOI [10.15585/mmwr.mm6911a2](https://doi.org/10.15585/mmwr.mm6911a2).
- Marçôa R, Ribeiro AI, Zão I, Duarte R. 2018.** Tuberculosis and gender—factors influencing the risk of tuberculosis among men and women by age group. *Pulmonology* **24**(3):199–202 DOI [10.1016/j.pulmoe.2018.03.004](https://doi.org/10.1016/j.pulmoe.2018.03.004).
- McWilliams A, Tammemagi MC, Mayo JR, Roberts H, Liu G, Soghrati K, Yasufuku K, Martel S, Laberge F, Gingras M, Atkar-Khattra S, Berg CD, Evans K, Finley R, Yee J, English J, Nasute P, Goffin J, Puksa S, Stewart L, Scott T, Johnston MR, Manos D, Nicholas G, Goss GD, Seely JM, Amjadi K, Tremblay A, Burrowes P, MacEachern P, Bhatia R, Tsao M-S, Lam S. 2013.** Probability of cancer in pulmonary nodules detected on first screening CT. *The New England Journal of Medicine* **369**(10):910–919 DOI [10.1056/NEJMoa1214726](https://doi.org/10.1056/NEJMoa1214726).
- Mu W, Jiang L, Zhang J, Shi Y, Gray JE, Tunali I, Gao C, Sun Y, Tian J, Zhao X, Sun X, Gillies RJ, Schabath MB. 2020.** Non-invasive decision support for NSCLC treatment using PET/CT radiomics. *Nature Communications* **11**(1):5228 DOI [10.1038/s41467-020-19116-x](https://doi.org/10.1038/s41467-020-19116-x).
- Peikert T, Bartholmai BJ, Maldonado F. 2020.** Radiomics-based management of indeterminate lung nodules? Are we there yet? *American Journal of Respiratory and Critical Care Medicine* **202**(2):165–167 DOI [10.1164/rccm.202004-1279ED](https://doi.org/10.1164/rccm.202004-1279ED).
- Pisano C, O'Connor J, Krick S, Russell DW. 2020.** A fatal case of pneumocephalus during computed tomography-guided lung biopsy. *American Journal of Respiratory and Critical Care Medicine* **201**(12):e83–e84 DOI [10.1164/rccm.201902-0280IM](https://doi.org/10.1164/rccm.201902-0280IM).

- Rami-Porta R, Asamura H, Travis WD, Rusch VW. 2017.** Lung cancer—major changes in the American Joint Committee on Cancer eighth edition cancer staging manual. *CA: A Cancer Journal for Clinicians* **67(2)**:138–155 DOI [10.3322/caac.21390](https://doi.org/10.3322/caac.21390).
- Reid M, Arinaminpathy N, Bloom A, Bloom BR, Boehme C, Chaisson R, Chin DP, Churchyard G, Cox H, Ditiu L, Dybul M, Farrar J, Fauci AS, Fekadu E, Fujiwara PI, Hallett TB, Hanson CL, Harrington M, Herbert N, Hopewell PC, Ikeda C, Jamison DT, Khan AJ, Koek I, Krishnan N, Motsoaledi A, Pai M, Raviglione MC, Sharman A, Small PM, Swaminathan S, Temesgen Z, Vassall A, Venkatesan N, Van Weezenbeek K, Yamey G, Agins BD, Alexandru S, Andrews JR, Beyeler N, Bivol S, Brigden G, Cattamanchi A, Cazabon D, Crudu V, Daftary A, Dewan P, Doepel LK, Eisinger RW, Fan V, Fewer S, Furin J, Goldhaber-Fiebert JD, Gomez GB, Graham SM, Gupta D, Kamene M, Khaparde S, Mailu EW, Masini EO, McHugh L, Mitchell E, Moon S, Osberg M, Pande T, Prince L, Rade K, Rao R, Remme M, Seddon JA, Selwyn C, Shete P, Sachdeva KS, Stallworthy G, Vesga JF, Vilc V, Goosby EP. 2019.** Building a tuberculosis-free world: the lancet commission on tuberculosis. *Lancet* **393(10178)**:1331–1384 DOI [10.1016/S0140-6736\(19\)30024-8](https://doi.org/10.1016/S0140-6736(19)30024-8).
- Siegel RL, Miller KD, Jemal A. 2019.** Cancer statistics, 2019. *CA: A Cancer Journal for Clinicians* **69(1)**:7–34 DOI [10.3322/caac.21551](https://doi.org/10.3322/caac.21551).
- Song Y, Zhang J, Zhang YD, Hou Y, Yan X, Wang Y, Zhou M, Yao YF, Yang G. 2020.** Feature Explorer (FAE): a tool for developing and comparing radiomics models. *PLOS ONE* **15(8)**:e0237587 DOI [10.1371/journal.pone.0237587](https://doi.org/10.1371/journal.pone.0237587).
- Sung H, Ferlay J, Siegel RL, Laversanne M, Soerjomataram I, Jemal A, Bray F. 2021.** Global cancer statistics 2020: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA: A Cancer Journal for Clinicians* **71(3)**:209–249 DOI [10.3322/caac.21660](https://doi.org/10.3322/caac.21660).
- Thwaites G, Nahid P. 2020.** Triumph and tragedy of 21st century tuberculosis drug development. *The New England Journal of Medicine* **382(10)**:959–960 DOI [10.1056/NEJMe2000860](https://doi.org/10.1056/NEJMe2000860).
- Wilson R, Devaraj A. 2017.** Radiomics of pulmonary nodules and lung cancer. *Translational Lung Cancer Research* **6(1)**:86–91 DOI [10.21037/tlcr.2017.01.04](https://doi.org/10.21037/tlcr.2017.01.04).
- Xu Y, Lu L, LN E, Lian W, Yang H, Schwartz LH, Yang ZH, Zhao B. 2019.** Application of radiomics in predicting the malignancy of pulmonary nodules in different sizes. *American Journal of Roentgenology* **213(6)**:1213–1220 DOI [10.2214/AJR.19.21490](https://doi.org/10.2214/AJR.19.21490).
- Yuan T, Sampson NS. 2018.** Hit generation in TB drug discovery: from genome to granuloma. *Chemical Reviews* **118(4)**:1887–1916 DOI [10.1021/acs.chemrev.7b00602](https://doi.org/10.1021/acs.chemrev.7b00602).