

BMJ Open VOC biomarkers identification and predictive model construction for lung cancer based on exhaled breath analysis: research protocol for an exploratory study

Wenwen Li,^{1,2} Wei Dai,³ Mingxin Liu,³ Yijing Long,² Chunyan Wang,² Shaohua Xie,^{3,4} Yuanling Liu,² Yinchenxi Zhang,² Qiuling Shi,⁵ Xiaoqin Peng,^{3,4} Yifeng Liu,^{3,4} Qiang Li,³ Yixiang Duan²

To cite: Li W, Dai W, Liu M, *et al.* VOC biomarkers identification and predictive model construction for lung cancer based on exhaled breath analysis: research protocol for an exploratory study. *BMJ Open* 2019;**9**:e028448. doi:10.1136/bmjopen-2018-028448

► Prepublication history and additional material for this paper are available online. To view these files, please visit the journal online (<http://dx.doi.org/10.1136/bmjopen-2018-028448>).

WL and WD contributed equally.

Received 8 December 2018

Revised 26 March 2019

Accepted 19 June 2019



© Author(s) (or their employer(s)) 2019. Re-use permitted under CC BY-NC. No commercial re-use. See rights and permissions. Published by BMJ.

For numbered affiliations see end of article.

Correspondence to

Professor Qiang Li; liqiang@sichuancancer.org and Dr Yixiang Duan; yduan@scu.edu.cn

ABSTRACT

Introduction Lung cancer is the most common cancer and the leading cause of cancer death in China, as well as in the world. Late diagnosis is the main obstacle to improving survival. Currently, early detection methods for lung cancer have many limitations, for example, low specificity, risk of radiation exposure and overdiagnosis. Exhaled breath analysis is one of the most promising non-invasive techniques for early detection of lung cancer. The aim of this study is to identify volatile organic compound (VOC) biomarkers in lung cancer and to construct a predictive model for lung cancer based on exhaled breath analysis.

Methods and analysis The study will recruit 389 lung cancer patients in one cancer centre and 389 healthy subjects in two lung cancer screening centres. Bio-VOC breath sampler and Tedlar bag will be used to collect breath samples. Gas chromatography-mass spectrometry coupled with solid phase microextraction technique will be used to analyse VOCs in exhaled breath. VOC biomarkers with statistical significance and showing abilities to discriminate lung cancer patients from healthy subjects will be selected for the construction of predictive model for lung cancer.

Ethics and dissemination The study was approved by the Ethics Committee of Sichuan Cancer Hospital on 6 April 2017 (No. SCCHEC-02-2017-011). The results of this study will be disseminated in presentations at academic conferences, publications in peer-reviewed journals and the news media.

Trial registration number ChiCTR-DOD-17011134; Pre-results.

INTRODUCTION

In China, lung cancer incidence and related mortality have been increasing annually for the last 30 years.^{1,2} Lung cancer is now the most common cancer and the leading cause of cancer death, accounting for 24.9% of all cancer deaths in China in 2010.^{1,2} Lung cancer is also the most common cancer and

Strengths and limitations of this study

- To construct a predictive model for early prediction of lung cancer based on exhaled biomarkers.
- This study will select exhaled biomarkers of lung cancer based on the largest sample size from multiple centres ever analysed by gas chromatography-mass spectrometry.
- Well-selected healthy controls with high risk but negative of lung cancer on chest CT will be recruited in this study.
- Gas chromatography-mass spectrometry can only obtain a limited amount of volatile organic compounds (VOCs) in exhaled breath and some VOCs cannot be detected.

the leading cause of cancer death worldwide.³ Only 16.8% of all patients with lung cancer survive for 5 years or more after diagnosis,⁴ mainly because lung cancer is often staged as locally advanced or metastatic disease at the initial diagnosis.^{4,5} However, the 5-year survival rate of patients with stage I lung cancer is greater than 60%,^{5,6} and the estimated 10-year survival rate of patients with stage I lung cancer detected on CT screening can reach 88%.⁷ Therefore, early detection is the key to improving survival rates of patients with lung cancer.

Currently, early detection techniques for lung cancer have many limitations, and a simple, reliable and non-invasive early lung cancer screening technique is urgently needed. Sputum cytology has a very low sensitivity.⁸ Chest radiography is radioactive and has a very high rate of false-negative results, especially in the detection of early stage lung cancers.⁹ Screening with low-dose CT can reduce mortality from lung cancer by 20%,

and lung cancer screening using low-dose CT for high-risk individuals is now recommended.^{10 11} However, there are also many disadvantages of low-dose CT, such as a high false-positive rate, overdiagnosis, the risk of radiation exposure and high cost, which limit its application in population-based screening.^{12 13}

Exhaled breath analysis is completely non-invasive and has great potential to become a screening and diagnostic method for early detection of cancer.^{14–17} The majority of previous studies focusing on lung cancer were conducted on small samples.^{18 19} Potential volatile organic compounds (VOCs) biomarkers for lung cancer have been discussed and summarised.^{15 19} However, to date, there are no unified VOC biomarkers for lung cancer, and the sets of VOCs employed vary between studies. Therefore, breath analysis is still in an early stage of clinical application. Several factors account for this situation, such as large variation in sample size, diverse sample collection approaches, different analytical techniques and different data processing methods.^{19–21} Our group has investigated exhaled VOC biomarkers in diabetes and breast cancer.^{22–25} We have rich experience in breath sample collection, exhaled VOCs analysis and data processing. In addition, we have a stable and reliable source of lung cancer patients and healthy subjects. In this study, we aim to identify exhaled VOC biomarkers of lung cancer and establish a predictive model for lung cancer. Our research hypothesis is that the predictive model will reach 80% sensitivity and 80% specificity through cross validation. In subsequent studies, this model will be validated in a clinical setting and population-based screening.

METHOD AND ANALYSIS

Main centres

Sichuan Cancer Hospital, Sichuan University, Nanchong Central Hospital and Chengdu Longquanyi District Center for Disease Control and Prevention. Lung cancer patients will be recruited from Sichuan Cancer Hospital. Healthy subjects will be recruited from two lung cancer screening centres, including Nanchong Central Hospital and Chengdu Longquanyi District Center for Disease Control and Prevention. Breath sample analysis will be conducted in Sichuan University.

Dates of the study

From 1 March 2017 to 28 February 2020.

Design

Inclusion criteria

Lung cancer patients and healthy subjects are both aged from 50 to 74 years. Lung cancer patients should have a pathological diagnosis of primary lung cancer based on the 2015 WHO Classification of lung tumours.²⁶ The pathological stages were based on the 8th edition of the TNM (Primary Tumor, Regional Lymph Nodes, Distant Metastasis) classification for lung cancer.²⁷ All the recruited patients should not receive any cancer treatment

before breath sampling. Healthy subjects recruited from two lung cancer screening centres should be negative of lung cancer on chest CT based on a previous project 'Early Diagnosis and Early Treatment of Rural Cancer' in China from 2014.

Exclusion criteria

Patients and controls with diabetes, other malignancies, active asthma, severe liver dysfunction, end-stage renal disease and acute inflammation will be excluded.

Breath sampling procedures

Each participant recruited into this study will be given an information leaflet explaining the research and will sign an informed consent form. Participants are required to fast for at least 8 hour and rest for at least 10 min in a separate room with good ventilation before breath sampling. Exhaled gas will be collected in the morning from 07:00 to 09:00. So, fasting from 23:00 the day before can meet the requirement. For lung cancer subjects who will not receive surgery, breath sample collection will be performed after pathological diagnosis and prior to cancer treatment. For patients undergoing surgery, breath samples will be collected the day before the surgery. If the postoperative pathology is not primary lung cancer, the patient will be excluded.

All the subjects will be asked to take a normal inhalation followed by a normal exhalation via their mouth. Deep inhalation before sampling and nasal ventilation during sampling will not be allowed. Each subject will exhale three times to complete a breath sample collection. Exhaled breath gas will be collected with Bio-VOC breath sampler (Markes Int. UK). Subjects will exhale gently into the Bio-VOC syringe until they feel mild resistance. The Bio-VOC sampler collects end-tidal breath gas by expelling the dead space air from the nasopharynx and upper airway using a three-way valve. The end-tidal breath gas will then be transferred to Tedlar bag (500 mL) through the three-way valve for storage and transportation. Breath samples will be kept at -40°C until analysis (within 7 days). The storage stability of VOCs in Tedlar bag at -40°C has been assessed (online supplementary figure 1). And the results indicated that VOCs could remain stable within 7 days in Tedlar bags at -40°C .

VOCs analysis

Solid phase microextraction (SPME) technique will be used in this study to preconcentrate VOCs in breath samples prior to gas chromatography-mass spectrometry (GC-MS) analysis. Divinylbenzene/carboxen/polydimethylsiloxane (DVB/CAR/PDMS) fibre will be used to extract exhaled VOCs. Samples will be analysed using a Thermo Scientific TRACE 1300 gas chromatograph coupled to a TSQ8000 triple quadrupole mass spectrometer. VF-624ms capillary column (60 m \times 0.25 mm \times 0.25 μm , Agilent Technology) will be used to separate VOCs. The analytical process will be as follows: first, the breath sample kept at -40°C will be incubated at 37°C for 5 min;

then a DVB/CAR/PDMS fibre will be used to preconcentrate VOCs for 30 min at 37°C, and finally the fibre will be desorbed thermally at the front inlet of the gas chromatograph. Split mode and a specific liner for SPME fibre will be used for the gas chromatograph. Electron ionisation source (70 eV) will be used for the mass spectrometer.

Endpoint

Exhaled VOC biomarkers of lung cancer and the accuracy of the predictive model for lung cancer.

Quality assurance

Standardisation of breath sampling

In order to minimise interobserver and intraobserver error caused by researchers during breath sampling, a fixed number of well-trained staff will be appointed to collect breath samples from patients and controls. Operation requirements that may bring about errors in breath sampling and lead to biased results are listed in online supplementary table 1. All staff participating in breath sampling will be trained and tested for the process and requirements of sampling. In addition, all researchers involved in this study will go through the certification process, including informed consent signing, breath sampling, samples storage and transportation.

Ambient room air

Breath samples will be collected from patients and controls in a separate room with good ventilation. Ambient room air will be collected simultaneously with each batch of breath samples. The ambient air is used to monitor possible VOC contamination, because it may be inhaled by subjects and become a prominent component of the exhaled breath, leading to anomalous results. Samples with VOCs that appear only in one batch of samples and simultaneously in the corresponding ambient room air at significant concentrations will be regarded as contamination and will be excluded.

Data management and monitoring

All data collected in this study, including demographic information, clinicopathological information and raw GC-MS data, will be uploaded to the Clinical Trial Management Public Platform (<http://www.medresman.org>). Data quality will be checked regularly by the principal investigator of this study. Data monitoring will be performed regularly by the Ethics Committee of Sichuan Cancer Hospital.

Statistical analysis and plans

Sample size calculation

Exhaled samples of 40 lung cancer patients and 40 healthy subjects were collected and analysed as preliminary data for sample size calculation. Thirty-one VOCs were measured in at least 70% of samples. Among them, ethanol, isoprene, acetone, isopropanol, benzene, ethylbenzene, octanal, nonanal and decane were reported in literature as exhaled biomarkers of lung cancer. The significant level is set as 0.029 based on False Discovery Rate (FDR) procedure for multiple

comparisons correction.²⁸ The FDR procedure is as follows: let $P_{(1)}, \dots, P_{(n)}$ be the ordered p values for testing hypotheses $H_0 = \{H_{(1)}, \dots, H_{(n)}\}$, and then H_0 is rejected if $P_{(i)} \leq ia/n$ for any $i=1, \dots, n$. In this study, $P_{(18)}=0.00068$ is the largest one that rejected H_0 , so $P=0.05 \times 18/31=0.029$. Finally, 350 subjects for each group are expected to identify the reported nine biomarkers (with additional 15 markers) that may be significantly different between lung cancer patients and controls, with a power of 90% or greater. From our previous work, 10% breath samples may be invalid due to sampling or storage faults. Therefore, the sample size of subjects for this study is estimated to be 778 ($700/0.9$), with 389 for lung cancer group and 389 for the control group. Mann-Whitney U test (SPSS, IBM, V.20.0) was used to evaluate the statistical differences of 31 VOCs in preliminary data. The sample size was calculated by G-power software based on normal distribution of the parent.

Statistical analysis for this study will include the following

Comparisons of VOCs between lung cancer patients and controls will be performed with the Mann-Whitney U test. Principal component analysis, linear discriminant analysis or independent component analysis, and so on will be used to reduce the dimension of datasets. Afterwards, multivariate analysis (supervised machine learning method such as PLSDA, OPLSDA, sPLSDA) and cross validation (10-fold or leave-one-out) will be used to identify exhaled biomarkers. Subjects will be randomly assigned to training set or testing set. A logistic regression model will be constructed based on exhaled biomarkers through the training set and trained by the testing set. Then the model will be verified through leave-one-out cross validation. Receiver operating characteristic (ROC) curve will be constructed and the area under the curve of the ROC curve will be used to evaluate the diagnostic accuracy of the model. Correct classification rates of the model will be calculated afterwards. The influence of potential confounders, for example, age, sex, smoking status, alcohol drinking, will also be investigated through univariate analysis. Logistic regression will be applied to evaluate the impact of potential factors on the identified biomarkers. A two-sided p value of <0.05 will be considered to indicate statistical significance. All analyses will be performed using SPSS (IBM, V.20.0) and WEKA software V.3.8.3. Data analyst will be blinded to group allocation during the data processing.

Patient and public involvement statement

Patients, healthy subjects and the public were not involved in the study design, recruitment and conduct. We do not have a plan to inform the result to the study participants unless they apply for it.

DISSEMINATION

The results of this study will be disseminated through various channels, including presentations at academic

conferences, publications in journals and in the news media.

Author affiliations

¹West China School of Public Health and West China Fourth Hospital, Sichuan University, Chengdu, China

²Research Center of Analytical Instrumentation, The College of Life Sciences, Sichuan University, Chengdu, China

³Department of Thoracic Surgery, Sichuan Cancer Hospital and Institute, Sichuan Cancer Center, School of Medicine, University of Electronic Science and Technology of China, Chengdu, China

⁴Graduate School, Chengdu Medical College, Chengdu, China

⁵Department of Symptom Research, The University of Texas MD Anderson Cancer Center, Houston, Texas, USA

Acknowledgements The authors appreciate all the lung cancer patients, patient advisers and healthy subjects who participate in this study. We are also grateful to related staff in Nanchong Central Hospital (Jun Bie) and Chengdu Longquanyi District Center for Disease Control and Prevention (Yang Shi, Honghai Ruan), who provide a lot of help in recruiting and managing healthy subjects.

Contributors YD and QL were involved in the study conception. QL and YD were involved in acquisition of funding and review of full protocol. WL and WD were involved in the study design and protocol writing. QS and CW were involved in statistical analysis plan. YJL was involved in sorting out operation requirements for breath sampling. ML, SX, YLL, YZ, XP and YFL were involved in drafting of the protocol.

Funding Sichuan Science and Technology Program (grant number 2017SZ0013).

Competing interests None declared.

Patient consent for publication Obtained.

Ethics approval This study was approved by the Ethics Committee of Sichuan Cancer Hospital on 6 April 2017 (No. SCCHEC-02-2017-011).

Provenance and peer review Not commissioned; externally peer reviewed.

Open access This is an open access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited, appropriate credit is given, any changes made indicated, and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>.

REFERENCES

- Chen W, Zheng R, Baade PD, *et al.* Cancer statistics in China, 2015. *CA Cancer J Clin* 2016;66:115–32.
- Chen W, Zheng R, Zeng H, *et al.* Epidemiology of lung cancer in China. *Thorac Cancer* 2015;6:209–15.
- Torre LA, Siegel RL, Jemal A. *Lung cancer statistics*. In: Ahmad A, Gadgeel S, *Lung cancer and personalized medicine*: Springer, Cham, 2016:1–19.
- Howlader N, Noone A-M, Krapcho M, *et al.* *SEER cancer statistics review (CSR) 1975–2011*. Bethesda, MD: National Cancer Institute, 2014.
- Hoffman PC, Mauer AM, Vokes EE. Lung cancer. *Lancet* 2000;355:479–85.
- Ou SH, Zell JA, Ziogas A, *et al.* Prognostic factors for survival of stage I nonsmall cell lung cancer patients: a population-based analysis of 19,702 stage I patients in the California Cancer Registry from 1989 to 2003. *Cancer* 2007;110:1532–41.
- Henschke CI, Yankelevitz DF, Libby DM, *et al.* Survival of patients with stage I lung cancer detected on CT screening. *N Engl J Med* 2006;355:1763–71.
- Gledhill A, Bates C, Henderson D, *et al.* Sputum cytology: a limited role. *J Clin Pathol* 1997;50:566–8.
- Sone S, Takashima S, Li F, *et al.* Mass screening for lung cancer with mobile spiral computed tomography scanner. *Lancet* 1998;351:1242–5.
- Aberle DR, Adams AM, Berg CD, *et al.* Reduced lung-cancer mortality with low-dose computed tomographic screening. *N Engl J Med* 2011;365:395–409.
- Moyer VA. U.S. Preventive Services Task Force. Screening for lung cancer: U.S. Preventive Services Task Force recommendation statement. *Ann Intern Med* 2014;160:330–8.
- Usman Ali M, Miller J, Peirson L, *et al.* Screening for lung cancer: a systematic review and meta-analysis. *Prev Med* 2016;89:301–14.
- Patz EF, Pinsky P, Gatsonis C, *et al.* Overdiagnosis in low-dose computed tomography screening for lung cancer. *JAMA Intern Med* 2014;174:269–74.
- Davis MD, Fowler SJ, Montpetit AJ. Exhaled breath testing - a tool for the clinician and researcher. *Paediatr Respir Rev* 2019;29.
- Hakim M, Broza YY, Barash O, *et al.* Volatile organic compounds of lung cancer and possible biochemical pathways. *Chem Rev* 2012;112:5949–66.
- Haick H, Broza YY, Mochalski P, *et al.* Assessment, origin, and implementation of breath volatile cancer markers. *Chem Soc Rev* 2014;43:1423–49.
- Pereira J, Porto-Figueira P, Cavaco C, *et al.* Breath analysis as a potential and non-invasive frontier in disease diagnosis: an overview. *Metabolites* 2015;5:3–55.
- Nardi-Agmon I, Peled N. Exhaled breath analysis for the early detection of lung cancer: recent developments and future prospects. *Lung Cancer* 2017;8:31–8.
- Saalberg Y, Wolff M. VOC breath biomarkers in lung cancer. *Clin Chim Acta* 2016;459:5–9.
- Mathew TL, Pownraj P, Abdulla S, *et al.* Technologies for clinical diagnosis using expired human breath analysis. *Diagnostics* 2015;5:27–60.
- Smolinska A, Hauschild AC, Fijten RR, *et al.* Current breathomics—a review on data pre-processing techniques and machine learning in metabolomics breath analysis. *J Breath Res* 2014;8:027105.
- Yan Y, Wang Q, Li W, *et al.* Discovery of potential biomarkers in exhaled breath for diagnosis of type 2 diabetes mellitus based on GC-MS with metabolomics. *RSC Adv* 2014;4:25430–9.
- Li J, Peng Y, Liu Y, *et al.* Investigation of potential breath biomarkers for the early diagnosis of breast cancer using gas chromatography-mass spectrometry. *Clin Chim Acta* 2014;436:59–67.
- Li W, Liu Y, Lu X, *et al.* A cross-sectional study of breath acetone based on diabetic metabolic disorders. *J Breath Res* 2015;9:016005.
- Li W, Liu Y, Liu Y, *et al.* Exhaled isopropanol: new potential biomarker in diabetic breathomics and its metabolic correlations with acetone. *RSC Adv* 2017;7:17480–8.
- Travis WD, Brambilla E, Nicholson AG, *et al.* The 2015 World Health Organization classification of lung tumors: impact of genetic, clinical and radiologic advances since the 2004 classification. *J Thorac Oncol* 2015;10:1243–60.
- Goldstraw P, Chansky K, Crowley J, *et al.* The IASLC Lung cancer staging project: proposals for revision of the TNM stage groupings in the forthcoming (eighth) edition of the TNM Classification for lung cancer. *J Thorac Oncol* 2016;11:39–51.
- Rom DM, Simes RJ. An improved Hochberg procedure for multiple tests of significance. *Br J Math Stat Psychol* 2013;66:751–4.