

## Original Article

# Screening of diagnostic biomarkers for lung cancer by bioinformatics analysis

Bao Liu<sup>1</sup>, An Yan<sup>1</sup>, Leiguang Ye<sup>1</sup>, Yina Gao<sup>1</sup>, Fang Liu<sup>1</sup>, Weibin Jing<sup>2</sup>, Limin Zhang<sup>3</sup>, Yan Yu<sup>1</sup>, Li Zhong<sup>1</sup>

<sup>1</sup>Department of Oncology, Harbin Medical University Cancer Hospital, Haerbin 150081, China; <sup>2</sup>Department of Burn, Hei Long Jiang Provincial Hospital, Haerbin 150036, China; <sup>3</sup>Head and Neck Surgery, Harbin Medical University Cancer Hospital, Haerbin 150081, China

Received May 11, 2015; Accepted March 5, 2016; Epub February 15, 2017; Published February 28, 2017

**Abstract:** Lung cancer was a leading cause of cancer-related death, this study aimed to explore target genes and specific biomarkers associated with the lung cancer for early diagnosis and treatment. Firstly, the gene expression profile GSE32863 was downloaded from Gene Expression Omnibus database (GEO), which included 15 lung cancer tissue samples and 10 normal lung tissue samples. In addition, the miRNA microarray profile GSE17681 with 10 lung cancer and 10 normal tissue samples was also downloaded. Then the probe-level data were pro-processed by Significance Analysis of Microarray (SAM), and the differentially expressed genes (DEGs) and miRNAs were screened with multi-test package in R language. The selected DEGs were further subjected to functional enrichment analysis using DAVID online tool. Following that, the verified target genes based on screened miRNAs were selected from miRTarBase and miRecords databases. Then miRNA-target gene regulation network was constructed to identify symbolic miRNAs of lung cancer. miR-29a with target genes such as FGG and COL4A1, miR-7 with SLC7A5 and miR-222 with ICAM-1 were found differentially expressed in lung cancer compared with the normal samples. The discovery of meaningful miRNAs and their differentially expressed target genes has the potential to be used in clinic for diagnosis and treatment of lung cancer.

**Keywords:** Lung cancer, differentially expressed gene, miRNA

### Introduction

Lung cancer was the most common cause of cancer-related death worldwide especially in industrialized countries, despite improvements in diagnosis and therapy, the overall 5-year survival was still 15% [1], possibly because lung cancer was often diagnosed at advanced stage and treatment options were limited. Smoking was the major risk factor for lung cancer [2], which can alter the activity of chemo-preventive drugs [3, 4], stimulate the clearance of selected targeted anticancer therapies [5], reduce the efficacy of cancer treatment, and increase the risk of second primary tumors. Although other factors, such as environmental exposure (e.g., chemicals, physical agents, and radiation), clinical history of lung diseases (e.g., chronic bronchitis, emphysema, pneumonia, and tuberculosis; ref. [6]), familial tumor history [7], or diet [8, 9], may also be associated with the development of lung cancer [10].

As detection of early lung cancer based on symptoms was not very effective, so studying the pathogenesis of lung cancer and exploring biomarker and new therapeutic targets had important practical significance [11].

Several molecular markers associated with lung cancer progression have been identified, including TGF, MET, TP53, HIF1A, APC, KRAS, EGFR and so on [12].

In recent years, MicroRNAs (MiRNAs) have become a hot research topic. MiRNAs were post-transcriptional regulators of gene expression [13]. The primary miRNA (primiRNA) was transcribed from the genome by RNA polymerase II and processed by Drosha, yielding the precursor miRNA transcript [14]. It was reported expression of at least 20-30% of human protein-coding genes were modulated by miRNAs [15].

Studies found miR-486, miR-30d, miR-1 and miR-499 from the serum may serve as non-invasive predictors for the overall survival of NSCLC [16]. With the penetrating study of miRNA in complex diseases, it was found that miRNA plays a great role during disease development [17, 18]. A recent report indicated disease may be cured by importing exogenous synthetic miRNA, therefore it was of great significance to identify the disease-related miRNAs and take them as targets for future therapies [19].

With the development of DNA microarray technology, it was now possible to screen many genes with expression alterations simultaneously. As a global approach, DNA microarray analysis was applied to investigate physiological mechanisms in health and disease [20]. Gene expression profiling based on microarrays was a robust and straightforward way to study the molecular features of cancer at a system level.

In this study, we aimed to identify target genes and specific biomarkers for identification and treatment of lung cancer with DNA microarray, and construct a miRNA co-expression network using miRNA micro-array data obtained by high-throughput means, so as to better understand the molecular mechanisms and explore new therapy strategies.

### Materials and methods

#### *Data resources and preprocessing*

The gene expression profile GSE32863 [21] was downloaded from GEO (Gene Expression Omnibus (<http://www.ncbi.nlm.nih.gov/geo/>)) database, which included 15 lung cancer tissue samples from previously untreated patients and 10 normal lung tissue samples, and miRNA microarray profile GSE17681 [22] with 10 lung cancer and 10 normal tissue samples was also downloaded. The test platform was GPL570 (Affymetrix Human Genome U133 Plus 2.0 Array). The probe-level data in CEL files were firstly converted into expression measures by ReadAffy of Affy package of R software. Then, we imputed the missing data [23] and normalized the quartile data [24]. The probes without gene annotation or having more than one gene annotation were filtered out; the average value of multiple probes corresponding to the same

gene was calculated as a unique value of the gene.

#### *Screening of differentially expressed genes (DEGs)*

The multistep package [25] in R language was used to identify DEGs between normal and lung cancer tissues. The Benjamin and Hochberg (BH) method [26] was used to adjust the raw *P*-values into false discovery rate (FDR) so as to avoid the multi-test problem which may cause too many false positive results. The  $FDR < 0.05$  and  $|\log_{2}FC| > 1$  were used as the cut-off thresholds.

#### *Disease-related miRNA screening*

Pearson correlation coefficients (that was co-expression intensity values) for each two miRNAs could be calculated by formula, and then miRNA co-expression networks were filtered out through a given threshold 0.7. A specific miRNA co-expression network was constructed with expression profiles of miRNA.

#### *Enrichment analysis for DEGs and disease-related miRNA*

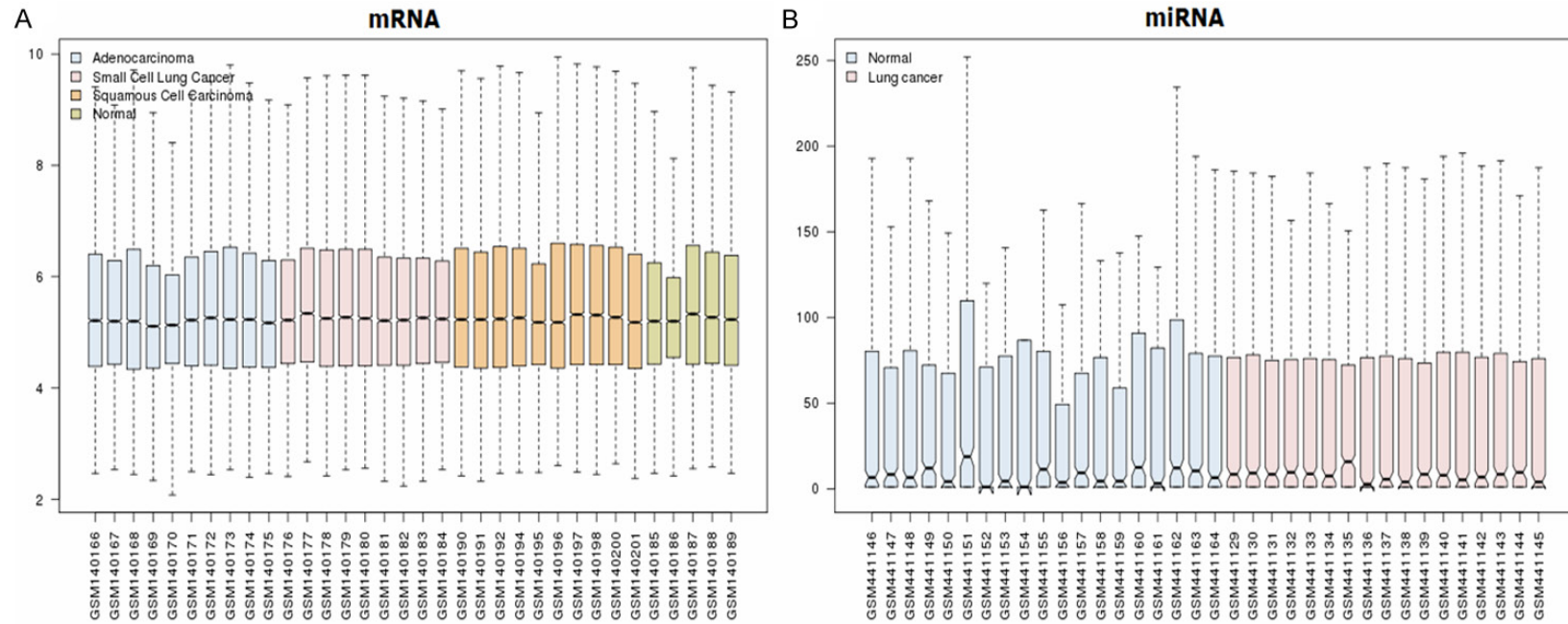
DAVID (the Database for Annotation, Visualization and Integrated Discovery) was a software with built-in rich graphical display, clustering the significant gene collections according to their functions, and it had abundant public database links [27]. The functional enrichment analysis of specific DEGs was performed by DAVID based on hypergeometric distribution algorithm ( $P < 0.05$ ).

If the set of disease-related miRNA screened was *M*, a known disease-associated miRNA set was *N*, fisher test was used to test whether the screened disease-related miRNA could enrich the already known miRNA based on *P*-value, then DAVID was used to do GO [28] functions and KEGG pathway enrichment analysis, so as to validate the function of the screened miRNA.

#### *Target gene selection of disease-related miRNAs*

According to different algorithms, each miRNA had different corresponding target genes. In order to find target genes with high confidence level, the mutual genes from two miRNA data-

## Biomarkers for lung cancer



**Figure 1.** Box graph of the standardized expression profile data. A. The blue, pink, orange and green columns represent gene expression profiles of adenocarcinomas, small cell lung cancer, squamous cell carcinomas and normal tissue samples, respectively. B. The blue and pink columns represent miRNA expression profiles of normal and unclassified lung cancer tissue samples, respectively.

## Biomarkers for lung cancer

**Table 1.** The top ten miRNAs with significant differences

| miRNA           | MDA_observe | MDA_predict | P_value |
|-----------------|-------------|-------------|---------|
| hsa-miR-32      | 0.835903    | -0.07087    | 0       |
| mmu-miR-486     | 0.847641    | -0.02171    | 0       |
| hsa-miR-520d-5p | 0.850673    | -0.01774    | 0       |
| hsa-miR-1265    | 0.822457    | -0.04277    | 0       |
| hsa-miR-16      | 0.857515    | 0.00105     | 0       |
| hsa-miR-369-5p  | 0.868527    | 0.018522    | 0       |
| hsa-miR-31      | 0.855548    | 0.018323    | 0       |
| hsa-miR-301a    | 0.846672    | 0.014547    | 0       |
| hsa-miR-19a     | 0.860351    | 0.037397    | 0       |
| hsa-miR-1298    | 0.850156    | 0.031328    | 0       |

MAD (mean absolute distance): observe represents the MAD value calculated under normal state, MAD (mean absolute distance): predict represents the average MAD values of miRNA calculated after 1000 times of disturbance.

**Table 2.** Gene Ontology enrichment terms of DEGs

| Term  | P-value  |
|---|----------|
| GO:0006928~cell motion  | 7.03E-04 |
| GO:0022604~regulation of cell morphogenesis                       | 0.00209  |
| GO:0045664~regulation of neuron differentiation                   | 0.002209 |
| GO:0051129~negative regulation of cellular component organization | 0.0028   |
| GO:0042330~taxis  | 0.004293 |
| GO:0006935~chemotaxis   | 0.004293 |
| GO:0050767~regulation of neurogenesis                             | 0.004891 |
| GO:0007626~locomotory behavior                                    | 0.005407 |
| GO:0051960~regulation of nervous system development               | 0.008129 |

bases (miRecords and miRTarBase) were screened as the target genes of miRNAs.

MiRecords was a database collecting target interactions of animal miRNAs, including the target genes verified by artificial collection experiments [29]. MiRTarBase was also a database which comprehensively collects miRNA targets based on experimental verification [30].

### Correlation analysis between DEGs and disease-related miRNAs

Target genes of disease-related miRNAs were firstly corresponded to the screened DEGs. According to the corresponding DEGs, unclassified miRNAs associated with lung cancer were categorized. Then miRNA-target gene regulation network was constructed to determine symbolic miRNAs of lung cancer. The reliability score of interaction relations more than 0.4 was regarded as an index to screen the reliable interaction relationship.

## Results

### DEGs and disease-related miRNAs analysis

After data preprocessing, we identified genes differentially expressed from the standardized data (**Figure 1**). A total of 1562 genes were screened as DEGs under condition of  $FDR < 0.05$  and  $|\log FC| > 1$ , which including 897 down-regulated genes (for example SUN1, SPTA1) and 665 up-regulated genes (for example SH3GL3, HTT).

25 miRNAs with significant differences were screened and the top10 miRNAs were listed in **Table 1**. By enrichment analysis we found that the screened miRNA related to lung cancer were enrichment of known miRNA (Fisher exact test,  $P = 1.769e-07$ ).

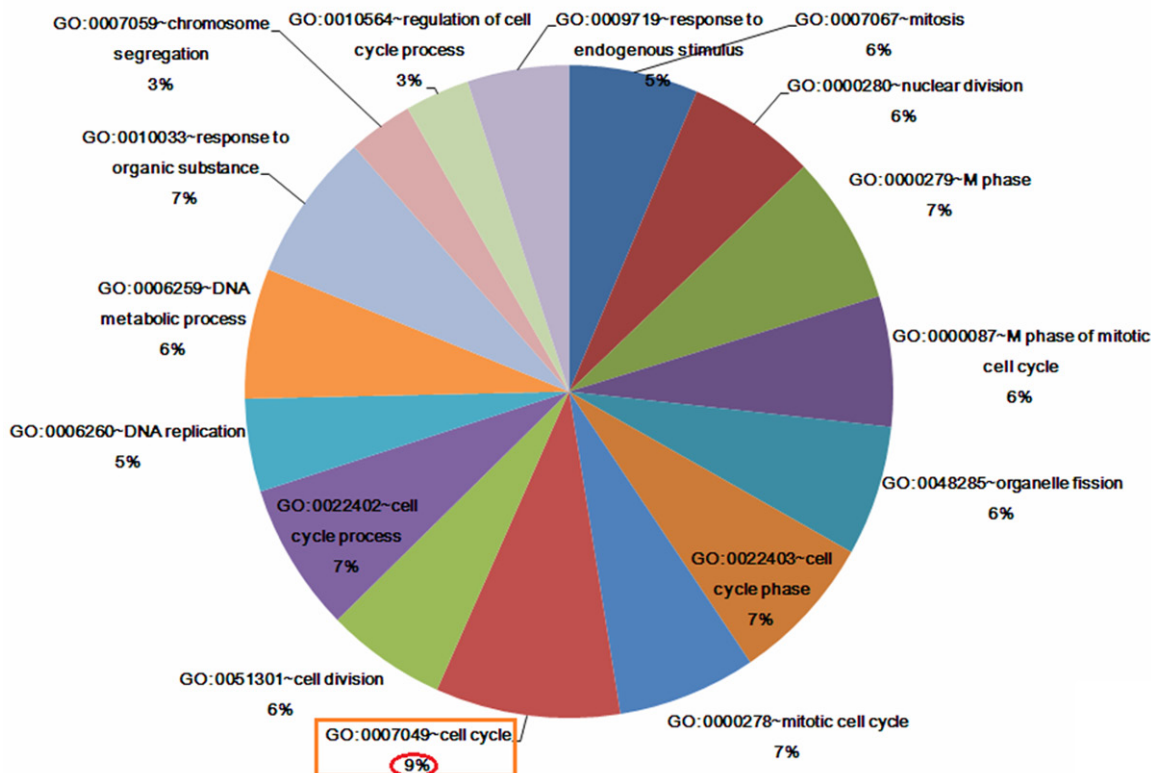
### Enrichment analysis of specific DEGs

DAVID was utilized to do enrichment analysis of specific DEGs, finally 9 GO functional nodes were obtained, such as regulation of cell morphogenesis, regulation of neuron differentiation and negative regulation of cellular component organization and so on, the results were listed in **Table 2**. The Gene Ontology enrichment map is shown in **Figure 2**.

### Target gene screening of miRNAs and correlation analysis with DEGs

The miRNAs information and the corresponding target genes were downloaded from miRecords and miRTarBase database. As shown in **Table 3**, the 38 screened target genes of differentially expressed miRNAs and the common part between target genes and DEGs were identified successfully. Then the regulation network for miRNA-target gene was constructed, the results were presented in **Figure 3**.

## Biomarkers for lung cancer



**Figure 2.** Gene ontology enrichment map. The function node in the frame was enriched with the most differentially co-expressed genes.

We could see that hsa-miR-29a was differentially expressed with target genes such as FGG and COL4A1, has-miR-7 with its target gene SLC7A5 and hsa-miR-222 with target gene ICAM-1 were identified. In addition, enrichment analysis of target gene function found that the disease-related miRNAs were mainly associated with cell adhesion and metabolic, regulation of cell morphogenesis, cell cycle process, regulation of cell proliferation and so on, which indicated that there was potential link between lung cancer and these function changes of miRNAs.

### Discussion

Lung cancer was the leading cause of cancer-related mortality worldwide and has become the largest threat to human health [31, 32]. Therefore, there was an urgent need to explore the mechanism of lung cancer and explore novel potential diagnostic and therapeutic targets. In this study, we linked DEGs with the disease-related miRNAs and concluded that has-miR-7 with SLC7A5, hsa-miR-222 with ICAM-1 and hsa-miR-29a with target genes such as

FGG and COL4A1, etc. were differentially expressed in lung cancer, which indicated that these target genes may be considered as potential biomarkers for the early diagnosis and future therapy of lung cancer caused by smoking.

MiRNA was a class of non-coding small RNA molecules, with length of about 22 nucleotides. It combined with 3 untranslated region of its target gene through the silence complex RNA-induced, leading to degradation of target mRNA or preventing translation of target miRNA. Since the first discovery of miRNA in *C. elegans* (*Caenorhabditis elegans*) in 1993, a large number of studies have shown that miRNA can part in cell growth, differentiation, proliferation, apoptosis, stress response and other biological processes by fine regulation of gene expression. With the deep study of complex diseases, it was found miRNA plays an important role in disease [19, 33].

It was reported Hsa-miR-7 played an important role in apoptosis (ACIN1, BAD, CASP8, CRYAA, GLO1, HSPA5, INHA, PCSK6, RELA, UBE4B),

## Biomarkers for lung cancer

**Table 3.** Differentially expressed miRNAs and corresponding target genes

| miRNA           | Target   |
|-----------------|----------|
| hsa-miR-199b-3p | KRT7     |
| hsa-miR-222     | ICAM1    |
| hsa-miR-29a     | FGG      |
| hsa-miR-29b     | FGG      |
| hsa-miR-29c     | FGG      |
| hsa-miR-7       | SNCA     |
| hsa-miR-29a     | BCL2     |
| hsa-miR-29b     | BCL2     |
| hsa-miR-29c     | BCL2     |
| hsa-miR-210     | CASP8AP2 |
| hsa-miR-210     | CBX1     |
| hsa-miR-126     | CCNE2    |
| hsa-miR-29b     | COL1A1   |
| hsa-miR-29c     | COL1A1   |
| hsa-miR-29a     | COL4A1   |
| hsa-miR-29b     | COL4A1   |
| hsa-miR-29c     | COL4A1   |
| hsa-miR-126     | DNMT1    |
| hsa-miR-29b     | DNMT1    |
| hsa-miR-210     | NCAM1    |
| hsa-miR-29a     | PPM1D    |
| hsa-miR-29a     | RAN      |
| hsa-miR-7       | SLC7A5   |
| hsa-miR-129-5p  | SOX4     |
| hsa-miR-29a     | SPARC    |
| hsa-miR-29c     | SPARC    |
| hsa-miR-29c     | SRSF10   |
| hsa-miR-29c     | TDG      |
| hsa-miR-19a     | TGFBR2   |
| hsa-miR-30e     | UBE2I    |
| hsa-miR-222     | KIT      |
| hsa-miR-29a     | GLUL     |

cell cycle (NUSAP1, STK11) and cell proliferation and differentiation (ALOX15B, TBX2) with its target genes [34]. In addition, its target gene SLC7A5 was part of a two-protein complex with SLC3A2, the heavy chain of a neutral amino-acid transporter implicated in nutrient transport at the blood-brain barrier and has been noted to be differentially expressed between adeno- and squamous cell lung carcinomas [35, 36]. SLC7A5 protein expression was positive in SCLC (small cell lung cancer) tumors [37], and suppression of SLC7A5 resulted in an increasing percentage of transfected cells in the G1 phase over time and a concomitant

decrease in the percentage of cells in the G2/M phase. Therefore we concluded that hsa-miR-7 cooperated with SLC7A5 to regulate in lung cancer.

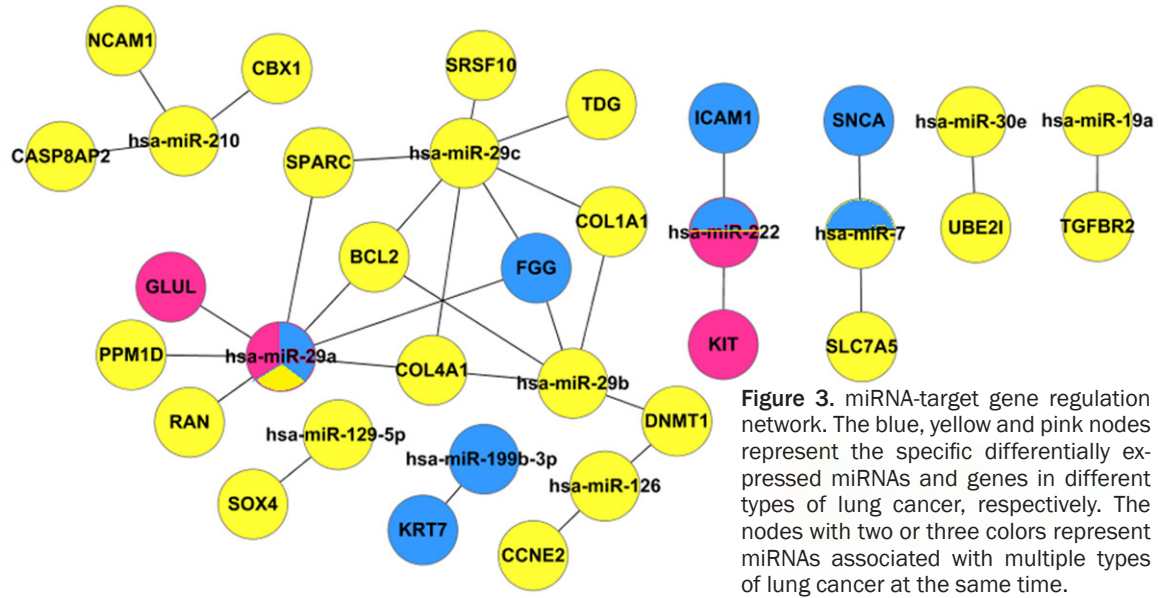
Hsa-miR-222 also found to play a critical role in cell cycle and multiple biological processes [38]. It has been reported to be involved in lung cancer metastases in vitro [39] and significantly down-regulated in the diseased serum samples [40]. ICAM-1 (Intercellular adhesion molecule-1), one target gene of hsa-miR-222 identified in this study, was a cell adhesion molecule with a key role in inflammation and immunosurveillance. It has been implicated in carcinogenesis by facilitating instability of the tumor environment [41] and increases significantly in NSCLC (non-small cell lung cancer) patients compared to healthy individuals [42]. The specific miRNAs and differentially expressed target genes found in lung cancer may be unique biomarkers.

Moreover, from GO enrichment analysis of DEGs we found interleukin 6 (IL-6) and interleukin 8 (IL-8), which were expressed in pre-malignant epithelial cells, and their expression were associated with a poor prognosis in lung cancer patients [43, 44]. Evidence showed that inflammatory mediators contributed to the pathogenesis of many human cancers, including lung cancer, as under inflammatory stress, IL-6 and IL-8 participated in tumorigenesis by acting directly on lung epithelial cells via signaling through the nuclear factor of kappa light polypeptide gene enhancer in B-cells 1 (NFkB1) pathway [45]. In addition, IL-6 and IL-8 were expressed by lung cancer cells and acted in an autocrine and/or paracrine fashion to stimulate cancer cell proliferation [46, 47], migration, and invasion [48].

To conclude, in the present study, we explored the DEGs and differentially expressed miRNAs of lung cancer caused by smoking, and compared the DEGs with the target genes of the miRNAs, the results showed that IL6, IL8, COL4A1, ICAM-1 and so on may have potential to be used as biomarkers for the early diagnosis and future target treatment for lung cancer.

Although great efforts have been made in the research of lung cancer, yet it is not enough to fully understand the pathogenesis of lung cancer and effective diagnostic and therapeutic strategies are needed to be developed.

## Biomarkers for lung cancer



**Figure 3.** miRNA-target gene regulation network. The blue, yellow and pink nodes represent the specific differentially expressed miRNAs and genes in different types of lung cancer, respectively. The nodes with two or three colors represent miRNAs associated with multiple types of lung cancer at the same time.

### Disclosure of conflict of interest

None.

**Address correspondence to:** Yan Yu, Department of Oncology, Harbin Medical University Cancer Hospital, 150 Haping Road, Nangang District, Harbin 150081, Heilongjiang Province, China. Tel: +86-45186298285; Fax: +8645186298285; E-mail: 153020426@163.com

### References

- [1] Boyle P, Chapman CJ, Holdenrieder S, Murray A, Robertson C, Wood WC, Maddison P, Healey GH, Fairley G, Barnes AC and Robertson JF. Clinical validation of an autoantibody test for lung cancer. *Ann Oncol* 2011; 22: 383-389.
- [2] Wood ME, Kelly K, Mullineaux LG and Bunn PA Jr. The inherited nature of lung cancer: a pilot study. *Lung Cancer* 2000; 30: 135-144.
- [3] Mayne ST and Lippman SM. Cigarettes: a smoking gun in cancer chemoprevention. *J Natl Cancer Inst* 2005; 97: 1319-1321.
- [4] The effect of vitamin E and beta carotene on the incidence of lung cancer and other cancers in male smokers. The Alpha-Tocopherol, Beta Carotene Cancer Prevention Study Group. *N Engl J Med* 1994; 330: 1029-35.
- [5] Hamilton M, Wolf JL, Rusk J, Beard SE, Clark GM, Witt K and Cagnoni PJ. Effects of smoking on the pharmacokinetics of erlotinib. *Clin Cancer Res* 2006; 12: 2166-2171.
- [6] Wu AH, Fontham ET, Reynolds P, Greenberg RS, Buffler P, Liff J, Boyd P, Henderson BE and Correa P. Previous lung disease and risk of lung cancer among lifetime nonsmoking wom-

en in the United States. *Am J Epidemiol* 1995; 141: 1023-1032.

- [7] Wu AH, Mimi CY, Thomas DC, Pike MC and Henderson BE. Personal and family history of lung disease as risk factors for adenocarcinoma of the lung. *Cancer Res* 1988; 48: 7279-7284.
- [8] Alavanja MC, Brown CC, Swanson C and Brownson RC. Saturated fat intake and lung cancer risk among nonsmoking women in Missouri. *J Natl Cancer Inst* 1993; 85: 1906-1916.
- [9] De Stefani E, Deneo-Pellegrini H, Mendilaharsu M, Carzoglio JC and Ronco A. Dietary fat and lung cancer: a case-control study in Uruguay. *Cancer Causes Control* 1997; 8: 913-921.
- [10] Samet JM, Avila-Tang E, Boffetta P, Hannan LM, Olivo-Marston S, Thun MJ and Rudin CM. Lung cancer in never smokers: clinical epidemiology and environmental risk factors. *Clin Cancer Res* 2009; 15: 5626-5645.
- [11] Pastorino U. Lung cancer screening. *Br J Cancer* 2010; 102: 1681-1686.
- [12] Wang L, Xiong Y, Sun Y, Fang Z, Li L, Ji H and Shi T. HLungDB: an integrated database of human lung cancer research. *Nucleic Acids Res* 2010; 38: D665-D669.
- [13] Chen CZ, Li L, Lodish HF and Bartel DP. MicroRNAs modulate hematopoietic lineage differentiation. *Science* 2004; 303: 83-86.
- [14] Carthew RW and Sontheimer EJ. Origins and mechanisms of miRNAs and siRNAs. *Cell* 2009; 136: 642-655.
- [15] Krichevsky AM. MicroRNA profiling: from dark matter to white matter, or identifying new players in neurobiology. *ScientificWorldJournal* 2007; 7: 155-166.

## Biomarkers for lung cancer

- [16] Hu Z, Chen X, Zhao Y, Tian T, Jin G, Shu Y, Chen Y, Xu L, Zen K and Zhang C. Serum MicroRNA signatures identified in a genome-wide serum MicroRNA expression profiling predict survival of non-small-cell lung cancer. *J Clin Oncol* 2010; 28: 1721-1726.
- [17] Bhayani MK, Calin GA and Lai SY. Functional relevance of miRNA\* sequences in human disease. *Mutat Res* 2012; 731: 14-19.
- [18] Kinet V, Dirx E and De Windt LJ. Quaero muneris: exploring microRNA function in cardiovascular disease. *J Mol Cell Cardiol* 2012; 52: 1-2.
- [19] Iorio MV and Croce CM. MicroRNA dysregulation in cancer: diagnostics, monitoring and therapeutics. A comprehensive review. *EMBO Mol Med* 2012; 4: 143-159.
- [20] Spies M, Dasu MR, Svrakic N, Nestic O, Barrow RE, Perez-Polo JR and Herndon DN. Gene expression analysis in burn wounds of rats. *Am J Physiol Regul Integr Comp Physiol* 2002; 283: R918-R930.
- [21] Rohrbeck A, Neukirchen J, Roskopf M, Pardillos GG, Geddert H, Schwalen A, Gabbert HE, von Haeseler A, Pitschke G and Schott M. Gene expression profiling for molecular distinction and characterization of laser captured primary lung cancers. *J Transl Med* 2008; 6: 69.
- [22] Keller A, Leidinger P, Borries A, Wendschlag A, Wucherpfennig F, Scheffler M, Huwer H, Lenhof HP and Meese E. miRNAs in lung cancer-studying complex fingerprints in patient's blood cells by microarray experiments. *BMC Cancer* 2009; 9: 353.
- [23] Troyanskaya O, Cantor M, Sherlock G, Brown P, Hastie T, Tibshirani R, Botstein D and Altman RB. Missing value estimation methods for DNA microarrays. *Bioinformatics* 2001; 17: 520-525.
- [24] Fujita A, Sato J, Rodrigues L, Ferreira C and Sogayar M. Evaluating different methods of microarray data normalization. *BMC bioinformatics* 2006; 7: 469.
- [25] Dudoit S, Shaffer JP and Boldrick JC. Multiple hypothesis testing in microarray experiments. *Sta Sci* 2003; 71-103.
- [26] Benjamini Y and Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)* 1995; 289-300.
- [27] Huang da W, Sherman BT, Lempicki RA. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat Protoc* 2008; 4: 44-57.
- [28] Harris M, Clark J, Ireland A, Lomax J, Ashburner M, Foulger R, Eilbeck K, Lewis S, Marshall B and Mungall C. The Gene Ontology (GO) database and informatics resource. *Nucleic Acids Res* 2004; 32: D258-261.
- [29] Xiao F, Zuo Z, Cai G, Kang S, Gao X and Li T. miRecords: an integrated resource for microRNA-target interactions. *Nucleic Acids Res* 2009; 37: D105-D110.
- [30] Hsu SD, Lin FM, Wu WY, Liang C, Huang WC, Chan WL, Tsai WT, Chen GZ, Lee CJ and Chiu CM. miRTarBase: a database curates experimentally validated microRNA-target interactions. *Nucleic Acids Res* 2011; 39: D163-D169.
- [31] Parkin DM, Bray F and Devesa S. Cancer burden in the year 2000. The global picture. *Eur J Cancer* 2001; 37: 4-66.
- [32] Pirozynski M. RETRACTED: 100 years of lung cancer. *Respir Med* 2006; 100: 2073-2084.
- [33] Conrad R, Barrier M and Ford LP. Role of miRNA and miRNA processing factors in development and disease. *Birth Defects Res C Embryo Today* 2006; 78: 107-117.
- [34] Gilad S, Lithwick-Yanai G, Barshack I, Benjamin S, Krivitsky I, Edmonston TB, Bibbo M, Thurm C, Horowitz L and Huang Y. Classification of the four main types of lung cancer using a microRNA-based diagnostic assay. *J Mol Diagn* 2012; 14: 510-517.
- [35] Kaira K, Oriuchi N, Imai H, Shimizu K, Yanagitani N, Sunaga N, Hisada T, Tanaka S, Ishizuka T and Kanai Y. Prognostic significance of L-type amino acid transporter 1 expression in resectable stage I-III non-small cell lung cancer. *Br J Cancer* 2008; 98: 742-748.
- [36] Kido Y, Tamai I, Uchino H, Suzuki F, Sai Y and Tsuji A. Molecular and functional identification of large neutral amino acid transporters LAT1 and LAT2 and their pharmacological relevance at the blood-brain barrier. *J Pharm Pharmacol* 2001; 53: 497-503.
- [37] Miko E, Margitai Z, Czimmerer Z, Várkonyi I, Dezső B, Lányi Á, Bacsó Z and Scholtz B. miR-126 inhibits proliferation of small cell lung cancer cells by targeting SLC7A5. *FEBS Lett* 2011; 585: 1191-1196.
- [38] Arnold CP, Tan R, Zhou B, Yue SB, Schaffert S, Biggs JR, Doyonnas R, Lo MC, Perry JM and Renault VM. MicroRNA programs in normal and aberrant stem and progenitor cells. *Genome Res* 2011; 21: 798-810.
- [39] Lin Q, Mao W, Shu Y, Lin F, Liu S, Shen H, Gao W, Li S and Shen D. A cluster of specified microRNAs in peripheral blood as biomarkers for metastatic non-small-cell lung cancer by stem-loop RT-PCR. *J Cancer Res Clin Oncol* 2012; 138: 85-93.
- [40] Guo L, Ji X, Yang S, Hou Z, Luo C, Fan J, Ni C and Chen F. Genome-wide analysis of aberrantly expressed circulating miRNAs in patients with coal workers' pneumoconiosis. *Molecular Biol Rep* 2013; 40: 3739-47.



## Biomarkers for lung cancer

- [41] Thanopoulou E, Kotzamanis G, Pateras IS, Ziras N, Papalambros A, Mariolis-Sapsakos T, Sigala F, Johnson E, Kotsinas A and Scorilas A. The single nucleotide polymorphism g. 1548A > G (K469E) of the ICAM-1 gene is associated with worse prognosis in non-small cell lung cancer. *Tumor Biol* 2012; 33: 1429-1436.
- [42] Guney N, Soydinc HO, Derin D, Tas F, Camlica H, Duranyildiz D, Yasasever V and Topuz E. Serum levels of intercellular adhesion molecule ICAM-1 and E-selectin in advanced stage non-small cell lung cancer. *Med Oncol* 2008; 25: 194-200.
- [43] Liu Y, Sun W, Zhang K, Zheng H, Ma Y, Lin D, Zhang X, Feng L, Lei W and Zhang Z. Identification of genes differentially expressed in human primary lung squamous cell carcinoma. *Lung Cancer* 2007; 56: 307-317.
- [44] Wolff KD, Follmann M and Nast A. The Diagnosis and Treatment of Oral Cavity Cancer. *Dtsch Arztebl Int* 2012; 109: 829.
- [45] Engels EA. Inflammation in the development of lung cancer: epidemiological evidence. *Expert Rev Anticancer Ther* 2008; 8: 605-615.
- [46] Perwez Hussain S and Harris CC. Inflammation and cancer: an ancient link with novel potentials. *Int J Cancer* 2007; 121: 2373-2380.
- [47] Lin W-W and Karin M. A cytokine-mediated link between innate immunity, inflammation, and cancer. *J Clin Invest* 2007; 117: 1175-1183.
- [48] Kamohara H, Ogawa M, Ishiko T, Sakamoto K and Baba H. Leukemia inhibitory factor functions as a growth factor in pancreas carcinoma cells: involvement of regulation of LIF and its receptor expression. *Int J Oncol* 2007; 30: 977.