

InCAR: a comprehensive resource for lncRNAs from Cancer Arrays

Yueyuan Zheng^{1,#}, Qingxian Xu^{1,#}, Mengni Liu^{1,#}, Huanjing Hu^{1,#}, Yubin Xie¹, Zhixiang Zuo^{1,*}, Jian Ren^{1,*}

1. State Key Laboratory of Oncology in South China, Cancer Center, Collaborative Innovation Center for Cancer Medicine, School of Life Sciences, Sun Yat-sen University, Guangzhou 510060, China

[#]These authors contributed equally to this work

*Correspondence to: Jian Ren (renjian@sysucc.org.cn); Zhixiang Zuo (zuozhx@sysucc.org.cn).

SUPPLEMENTARY METHODS

Integration of lncRNA annotations

All of the gene annotations from the Ensembl project and RefSeq database were combined using the Cuffcompare program in the Cufflinks suite (1-3). We obtained lncRNAs and protein-coding genes according to the following protocol steps (**Supplementary Fig. S3**). First, we merged the fully matched (class code “=”) or contained transcripts (class code “c”). Second, the gene-transcript mapping files in different databases were downloaded from the Ensembl and Refseq database to correct and supplement the results, respectively. Third, the transcripts were divided into protein-coding and non-coding transcripts according to the “Transcript biotype”. Fourth, small RNA and pseudogene transcripts were filtered out. Only transcripts longer than 200nt were retained as candidate lncRNA transcripts. We then predicted the coding ability of the lncRNA transcripts with the CPAT (4) and CNCI (5) programs and removed those transcripts which both predicted as protein-coding. Finally, the protein-coding and lncRNA transcripts were separately combined according to their gene names and genomic coordinates.

Probe re-annotation

The following steps were taken to re-annotate the probe sequences (**Supplementary Fig. S3**). First, all of the probe sequences from the different platforms were aligned to the human genome (GRCh38) using a BLAST-like alignment tool (BLAT) (6). Only the alignments with no more than one mismatch, no gaps and a similarity score larger than 90 were preserved. In general, the probes (50 to 60 nucleotides) from Agilent, Illumina and other companies were designed to locate target genes or transcripts. The microarrays from Affymetrix utilized a probe set containing a group of 25mer probes to represent a gene or transcript. Thus, for the Affymetrix data we combined probes that specifically corresponded to the same probe set and ensured there were at least 3 perfectly matching and adjacent probes in each probe set. Second, we mapped the probes to coding genes or lncRNAs according to their genomic coordinates. If any probes were targeted to both coding genes and lncRNAs, we only preserved the annotation of the coding genes.

Differential expression analysis

We categorized differential expression with one of the following conditions: cancer versus normal controls, high stage versus low stage, high grade versus low grade, metastasis versus primary, drug-treatment versus control, and “other features”. Other features referred to certain cancer-specific features, such as estrogen receptor, progesterone receptor and *HER2* status in breast cancer as well as the smoking and drinking status in lung cancer and esophageal cancer. The differential expression analysis under each condition was performed with the Limma

package (7). We utilized a robust rank aggregation algorithm to integrate the lncRNA profiles in an unbiased manner (8). The aggregation rank score (AR score) represents the integrated rank from the meta-analysis of the fold-change in the different microarray studies. All results were scaled by cancer subtype and presented in the heatmaps, allowing users to explore the expression of lncRNAs of interest in the different conditions.

Survival analysis

All of the data we collected contained four types of survival data: overall survival, relapse-free survival, metastasis-free survival and progression-free survival. For each study with survival information, we accessed the associations between lncRNA expression and survival via a univariate Cox regression analysis (R 'survival' package v. 3.4) (9). Specifically, genes with hazard ratios >1 and a p value <0.05 were considered as poorly prognostic lncRNAs, whereas genes with hazard ratios of 0-1 and a p value < 0.05 were lncRNAs with good prognosis value. To check the prognostic landscape of lncRNAs in each cancer type, and to reduce the differences between the different studies, we utilized the Z scores in the Cox regression to perform a meta-analysis. Z scores were directly associated with P values, and were independent of different follow-up time scales and the range of input variables, allowing a direct comparison across multiple studies and platforms. We calculated the prognostic meta Z scores of lncRNAs in each cancer type by using Lipták's weighted meta Z test, in which weights refer to the square roots of the sample sizes (10,11). lncRNAs with positive meta Z scores indicated poor prognosis value in most studies; otherwise they were a bad prognosis. The results were shown in heatmaps in order to facilitate comparative analysis.

Advanced annotations and other analysis

Advanced annotations contained the transcript information, structure, coding potential and conservative score. Briefly, all transcript information was extracted from NCBI and Ensembl, and the transcript structure was predicted by RNAfold (12). CPAT and CNCI were used to assess the coding potential of the transcripts. The conservative scores of the transcripts were calculated from Phastcon100 (13). In addition, cancer-related lncRNAs validated by low-throughput experiments were manually collected from Lnc2Cancer, LncRNADisease, EVLncRNAs and MNDR v2.0. These data were integrated into InCAR as validation datasets.

Moreover, our database also included co-expression analysis, pathway enrichment analysis and ceRNA analysis modules to help researchers further explore the function of lncRNAs of interest. To determine the co-expression patterns of each lncRNA and protein-coding genes in a study, we calculated pairwise expression correlations across all of the protein-coding genes. The first 200 lncRNA-coding gene pairs with strong correlations (correlation coefficient $|r|>0.3$) were selected to represent the co-expression network. Users are able to further determine

the potential function of such lncRNA through the pathway enrichment analysis of co-expressed protein-coding genes. In addition, we obtained the lncRNA-miRNA and miRNA-coding gene interaction pairs from starBase2 and constructed a ceRNA interaction network based on the lncRNAs and protein-coding genes having positive expression correlations (14).

Determination of cancer subtypes by lncRNA expression

Initially, to normalize the cross-platform datasets, some filters were employed on each dataset: i) lncRNAs shared by 2-3 platforms with the largest samples were selected as the standard and other expression platforms were merged subsequently. ii) For each study, batch effects were minimized using the Combat function (15). iii) MergeMaid was applied to remove lncRNAs of low reliability across platforms and the preserved data was scaled by mean centering on genes and samples, respectively (16). iv) 500 lncRNAs with the highest expression median difference were screened for the clustering process. Secondly, we exploited the Consensus Clustering method to classify the cancer subtypes based on lncRNA expression (17). The number of clusters in each cancer was determined by the relative change in area under the CDF curve (the delta area < 0.05). Then the differentially expressed lncRNAs in each cluster were identified using a moderated F-test, and these lncRNAs were merged for the next iteration. The iteration stopped when the lncRNAs were unchanged from the previous run. The resulting lncRNAs were regarded as the gene signatures for each cancer subtype.

Web interface implementation

In order to visualize the analysis results, multiple statistical diagrams were embedded in the web server. For example, the heatmaps were constructed by DataTables, the secondary structures of transcripts were shown using FornaContainer (18), the expression network were presented by Highcharts, the boxplot together with the bubble chart were demonstrated by Echarts, and the circos plots were displayed though BioCircos.js (19). Furthermore, all of the analyses in the lnCAR website were performed in R.

In addition, we allowed the genomic coordinates of the regions of interest as an input to search the expression and prognosis of any user-defined lncRNA. For each cancer, we filtered probes that targeted protein-coding transcripts in each platform and merged all the remaining probes as a candidate lncRNA library. Then the regions submitted by users were transformed into “bed” format and intersected with the candidate lncRNA probes by applying the BEDTools’ intersectBed (v. 2.26.0) (20). The differential expression analysis and survival analysis of each targeted probes were presented in the result pages.

Reference

1. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, *et al.* Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nature biotechnology* **2010**;28:511-5
2. Zerbino DR, Achuthan P, Akanni W, Amode MR, Barrell D, Bhai J, *et al.* Ensembl 2018. *Nucleic acids research* **2018**;46:D754-D61
3. Haft DH, DiCuccio M, Badretdin A, Brover V, Chetvernin V, O'Neill K, *et al.* RefSeq: an update on prokaryotic genome annotation and curation. *Nucleic acids research* **2018**;46:D851-D60
4. Wang L, Park HJ, Dasari S, Wang S, Kocher JP, Li W. CPAT: Coding-Potential Assessment Tool using an alignment-free logistic regression model. *Nucleic acids research* **2013**;41:e74
5. Sun L, Luo H, Bu D, Zhao G, Yu K, Zhang C, *et al.* Utilizing sequence intrinsic composition to classify protein-coding and long non-coding transcripts. *Nucleic acids research* **2013**;41:e166
6. Kent WJ. BLAT--the BLAST-like alignment tool. *Genome research* **2002**;12:656-64
7. Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, *et al.* limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic acids research* **2015**;43:e47
8. Kolde R, Laur S, Adler P, Vilo J. Robust rank aggregation for gene list integration and meta-analysis. *Bioinformatics* **2012**;28:573-80
9. Therneau TM, Grambsch PM. Modeling survival data: extending the Cox model. Springer Science & Business Media; 2013.
10. Zaykin DV. Optimally weighted Z-test is a powerful method for combining probabilities in meta-analysis. *Journal of evolutionary biology* **2011**;24:1836-41
11. Gentles AJ, Newman AM, Liu CL, Bratman SV, Feng W, Kim D, *et al.* The prognostic landscape of genes and infiltrating immune cells across human cancers. *Nature medicine* **2015**;21:938-45
12. Lorenz R, Hofacker IL, Stadler PF. RNA folding with hard and soft constraints. *Algorithms for molecular biology : AMB* **2016**;11:8
13. Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome research* **2005**;15:1034-50
14. Li JH, Liu S, Zhou H, Qu LH, Yang JH. starBase v2.0: decoding miRNA-ceRNA, miRNA-ncRNA and protein-RNA interaction networks from large-scale CLIP-Seq data. *Nucleic acids research* **2014**;42:D92-7
15. Leek JT, Johnson WE, Parker HS, Jaffe AE, Storey JD. The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics* **2012**;28:882-3
16. Cope L, Zhong X, Garrett E, Parmigiani G. MergeMaid: R tools for merging and cross-study validation of gene expression data. *Statistical applications in genetics and molecular biology* **2004**;3:Article29
17. Wilkerson MD, Hayes DN. ConsensusClusterPlus: a class discovery tool with confidence assessments and item tracking. *Bioinformatics* **2010**;26:1572-3
18. Kerpedjiev P, Hammer S, Hofacker IL. Forna (force-directed RNA): Simple and effective

- online RNA secondary structure diagrams. *Bioinformatics* **2015**;31:3377-9
19. Cui Y, Chen X, Luo H, Fan Z, Luo J, He S, *et al.* BioCircos.js: an interactive Circos JavaScript library for biological data visualization on web applications. *Bioinformatics* **2016**;32:1740-2
20. Quinlan AR. BEDTools: The Swiss-Army Tool for Genome Feature Analysis. *Current protocols in bioinformatics* **2014**;47:11 2 1-34

SUPPLEMENTARY FIGURES

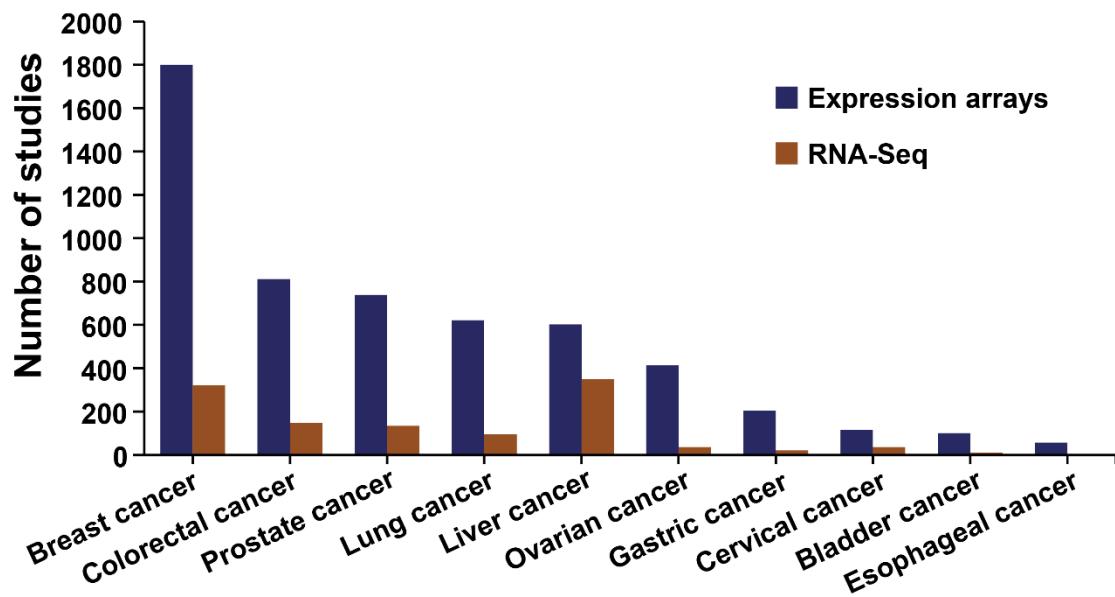


Fig. S1 – Comparison of the number of published gene expression datasets and RNA-Seq datasets from GEO database in cancer . All the datasets in GEO database were limited to human studies published before September 2017.

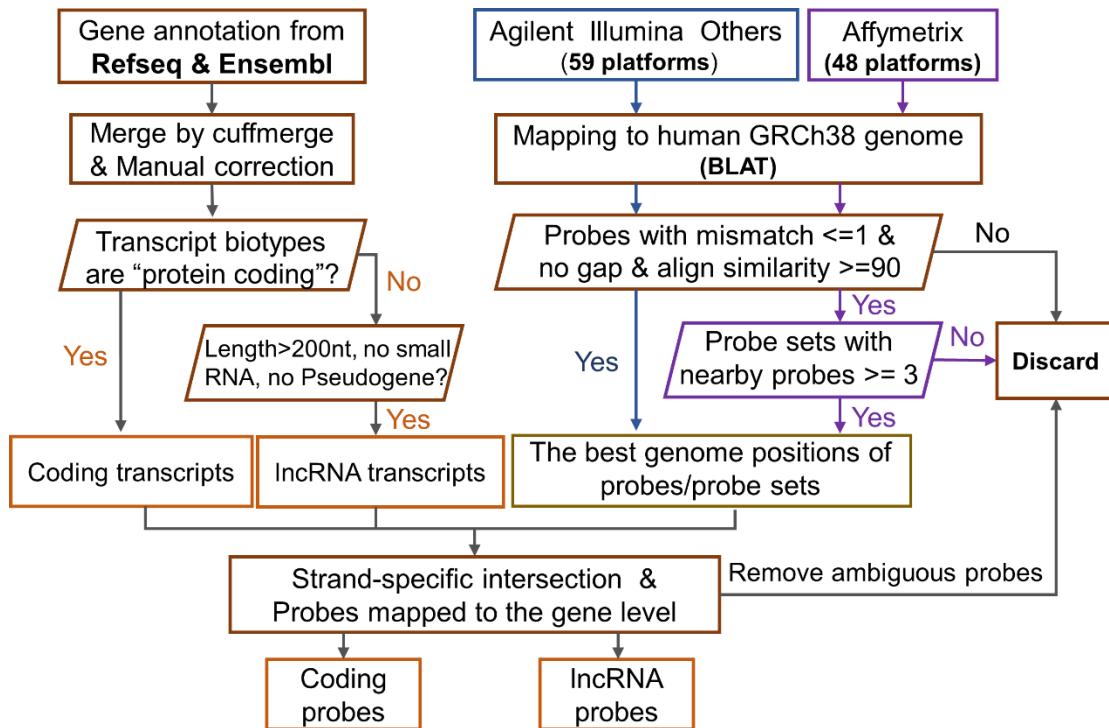


Fig. S2 - Computational pipeline for re-annotation of microarray probes

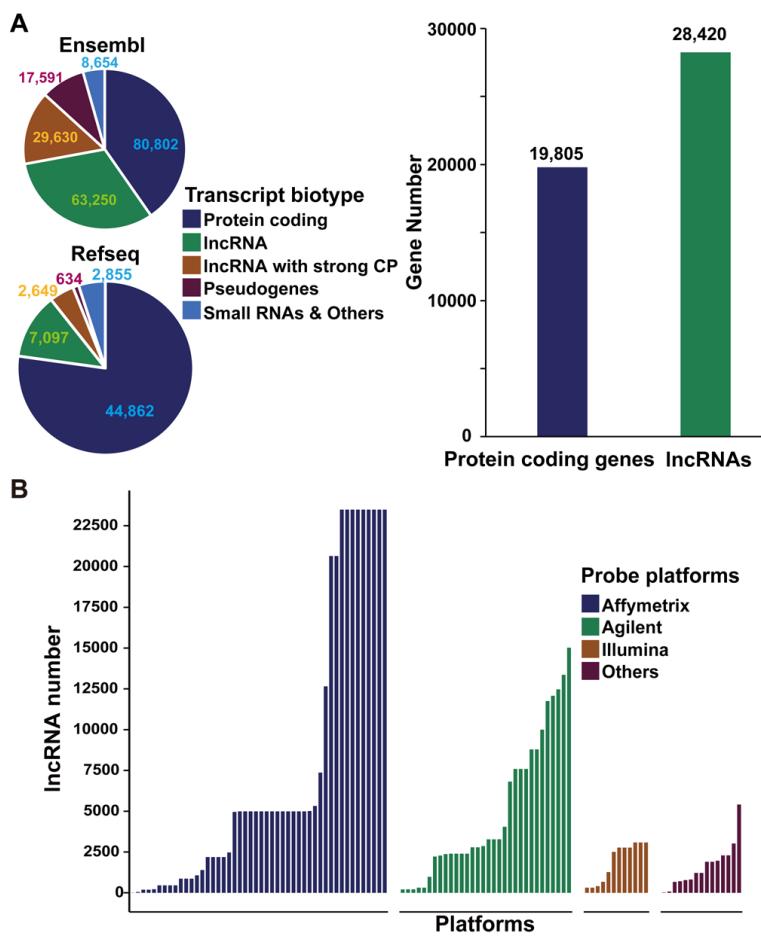


Fig. S3 - Re-annotation of the array probes. (A) The percentage of the transcript type distribution in protein-coding gene, lncRNA, psedogenes, small RNAs and others for the transcripts from the Ensembl and RefSeq annotations (left panel), and the number of protein coding genes and lncRNAs in Ensembl and RefSeq annotations (right panel). (B) The number of lncRNAs probes reannotated from different array platforms including Affymetrix, Agilent and Illumina.

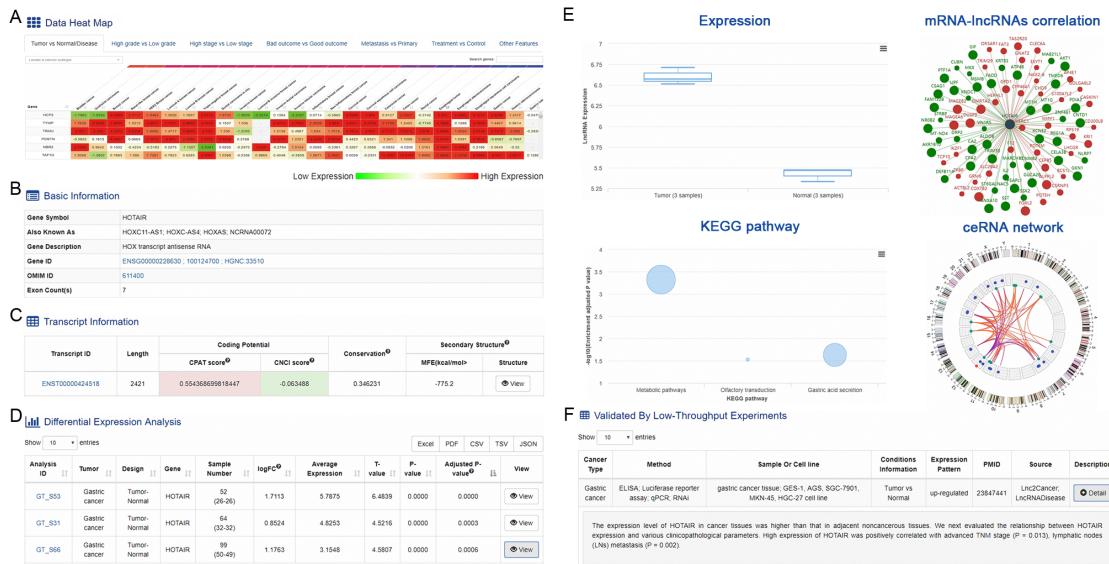


Fig. S4 - A snapshot of the example for differential expression analysis in InCAR. (A) The lncRNA expression characteristics in different cancer subtypes under the same conditions. Red cells indicate increased expression while green cells indicate decreased expression. (B) The basic information on lncRNA gene. (C) The information on the lncRNA transcripts, including the coding potential, degree of conservation and secondary structure. (D-E) The expression of lncRNA in different studies of gastric cancer compared to normal tissues. Click the “view” button to get the expression pattern of lncRNAs in different tissues, the related coding genes, the potential biological pathways and related ceRNA network. (F) The validated results in the low-throughput experiments.

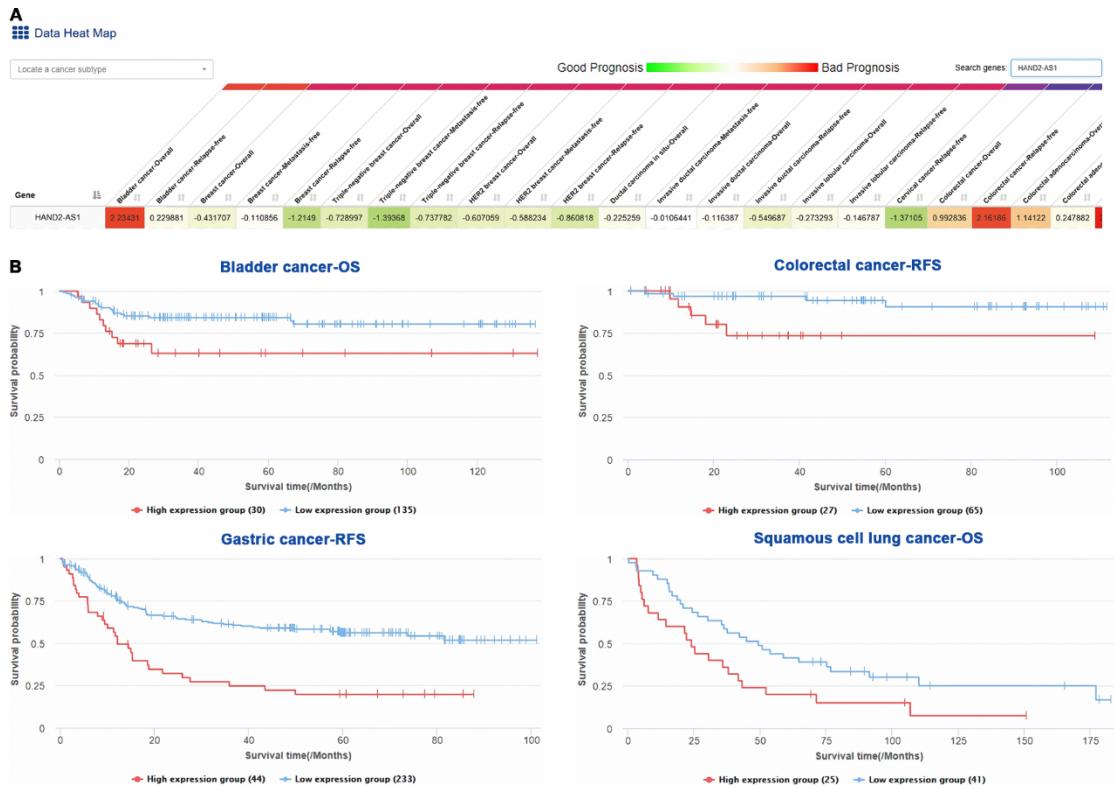


Fig. S5 - A snapshot of the example for survival analysis in InCAR. (A) The prognostic impact of *HAN2-AS1* in different cancer subtypes was shown in heatmap. (B) Survival curves of *HAN2-AS1* in bladder cancer, colorectal cancer, gastric cancer and squamous cell lung cancer indicate a poor prognosis. OS represents over-all survival, while the RFS means relapse-free survival.

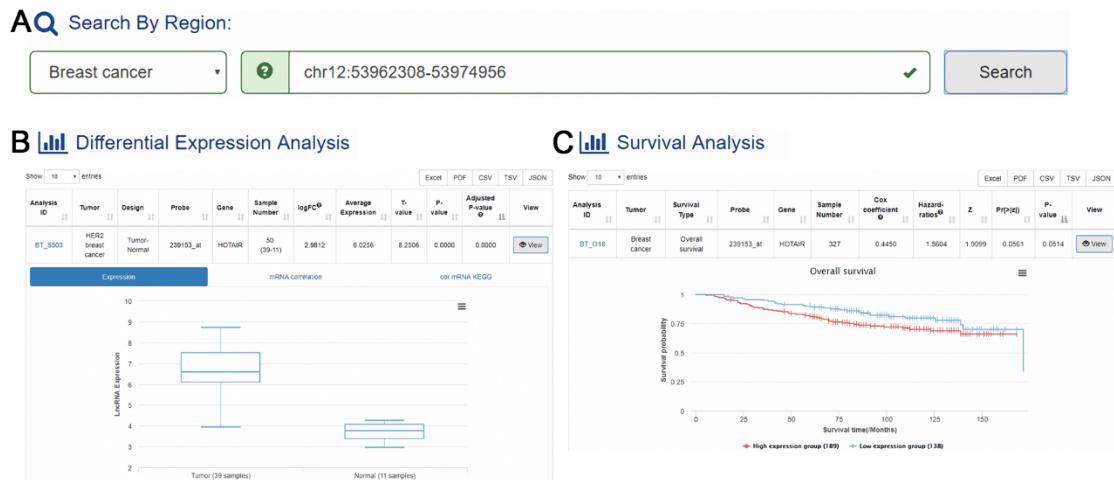


Fig. S6 - A schematic workflow of My-lncRNA interface. (A) Search by chromosome regions of interest. Take “chr12:53962308-53974956” as an example. (B) Results of target probes by

differential expression analysis in different studies. (C) Prognostic impact of target probes by survival analysis.

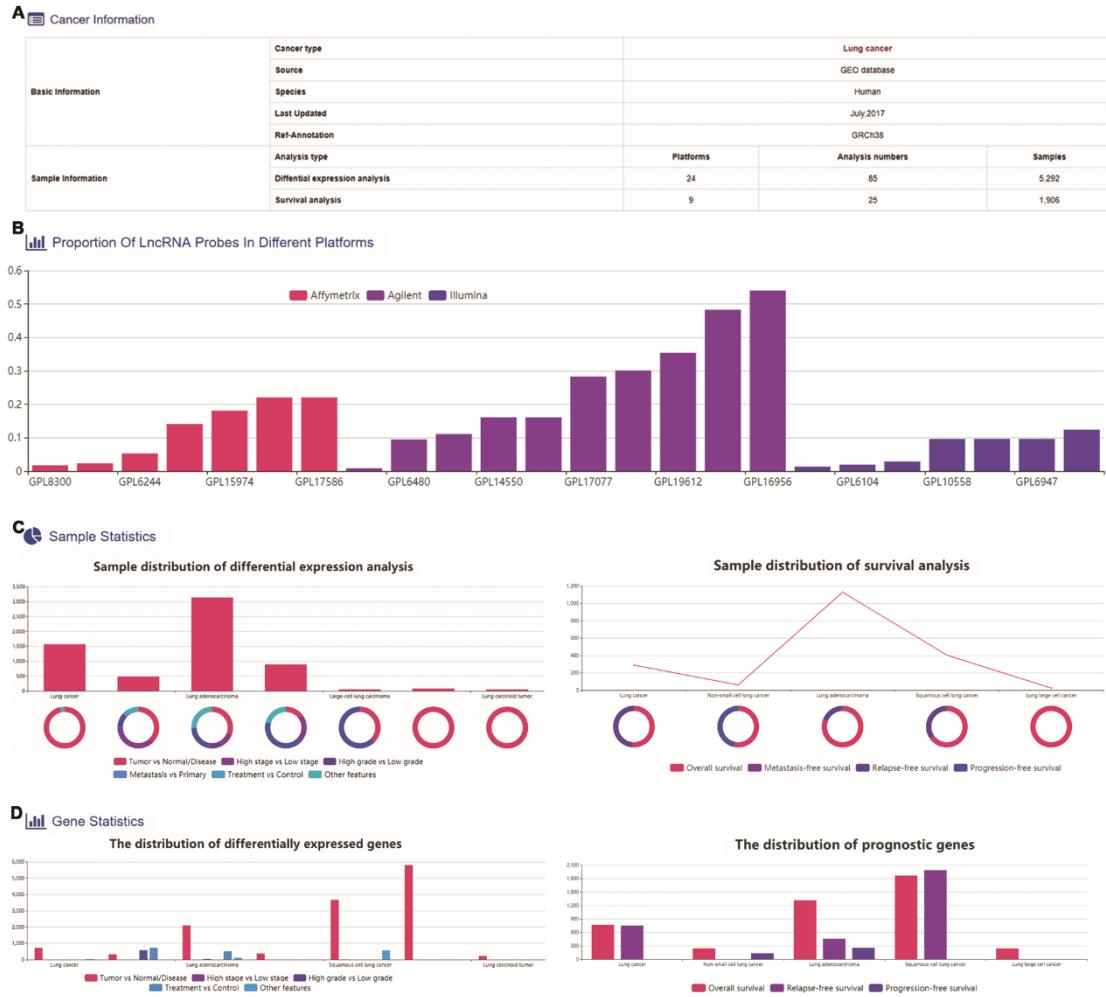


Fig. S7 - A snapshot of the browse interface. Take Lung cancer as an example. (A) Summary information of the cancer. (B) The proportion of lncRNA probes in different platforms. (C) The sample distribution of differential expression analysis and survival analysis. (D) The distribution of differentially expressed genes and prognostic genes.

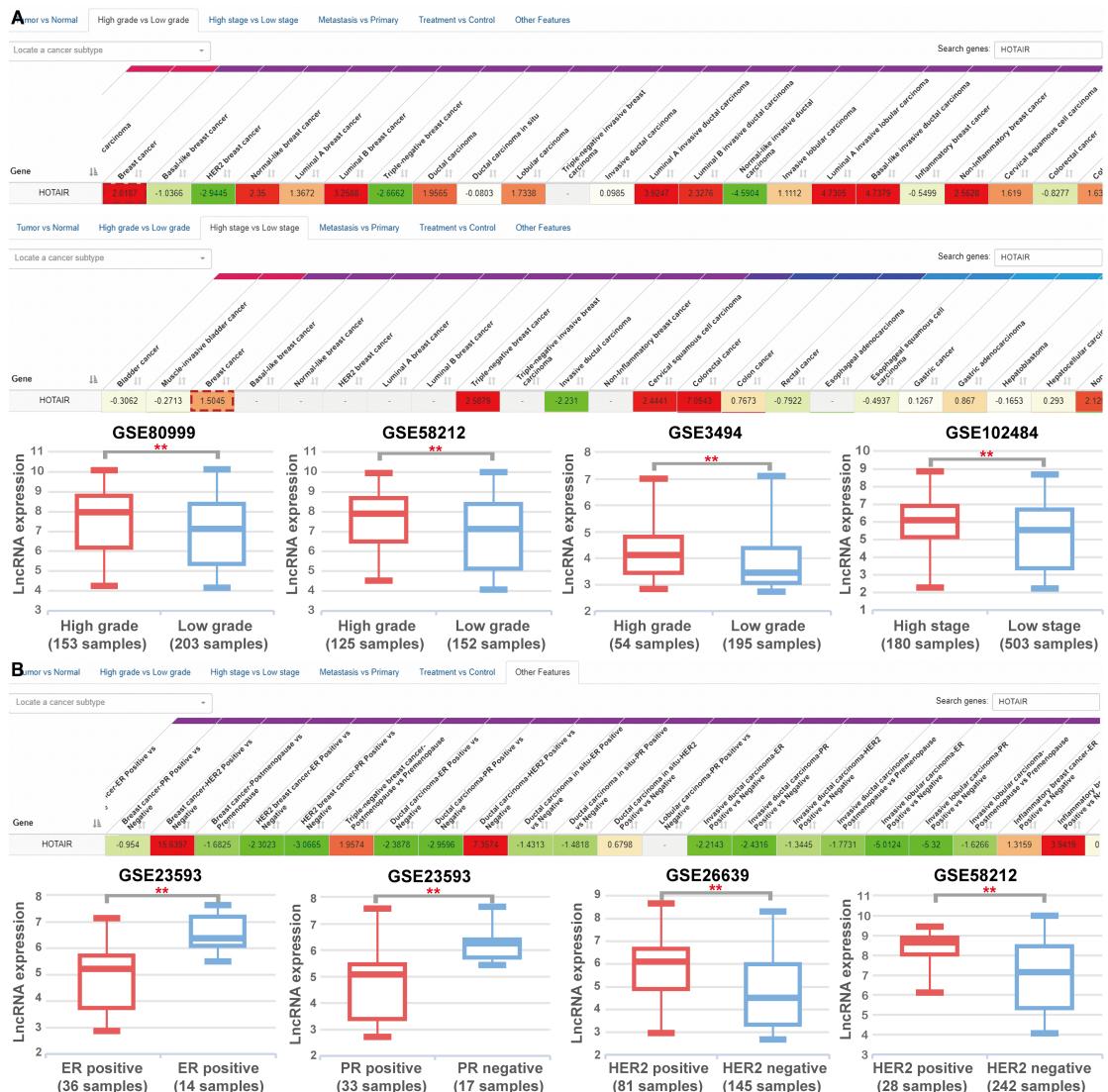


Fig S8 – Differential expression of lncRNA HOTAIR between different breast cancer clinical factors. (A) Increased *HOTAIR* expression is related to advanced tumor stage and higher tumor grade (B) *HOTAIR* appears to be overexpressed in ER-negative, PR-negative and HER2-positive patients.

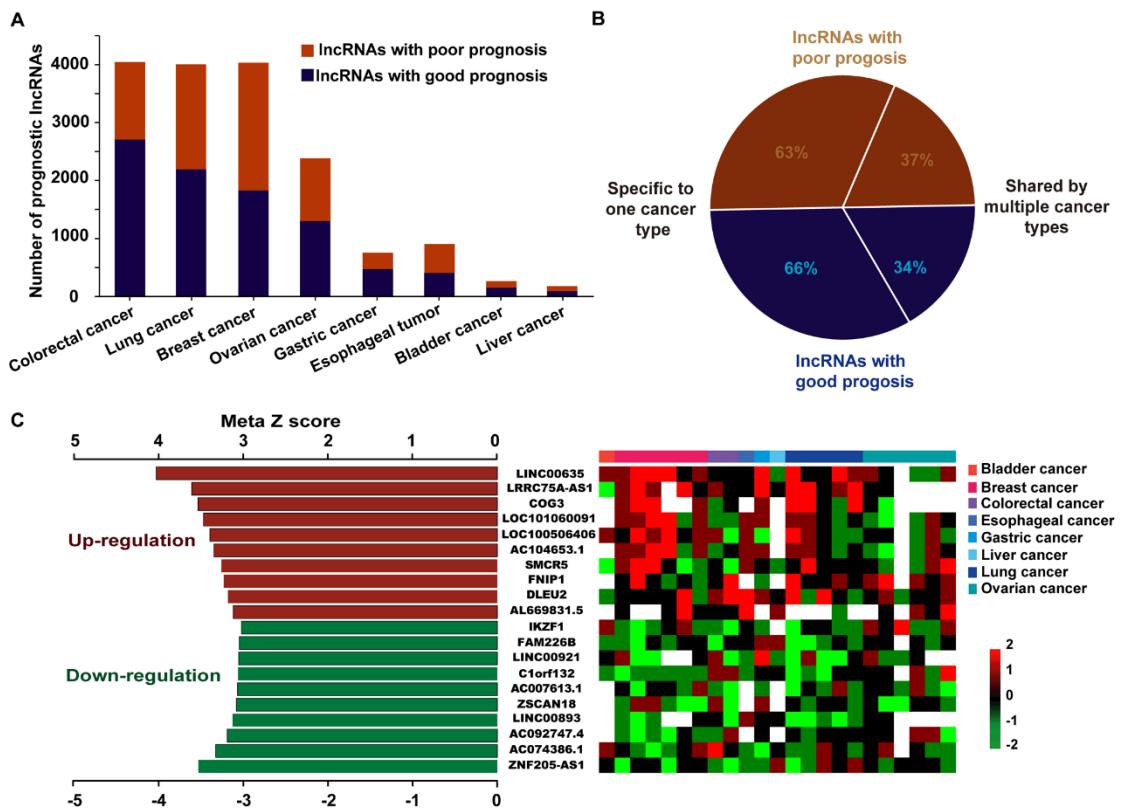


Fig. S9 - A large number of lncRNAs with potential clinical significance in various cancer types. (A) Bar plot showing the number of prognostic lncRNAs in different cancer types. (B) The proportion of prognostic lncRNAs specific to one cancer type and shared by multiple cancer types for poor prognosis and good prognosis. (C) The top 10 lncRNAs that were significantly associated with poor prognosis or good prognosis. Left panel is showing the Meta Z score of the P values from the survival analysis of different datasets of different cancer types. Right panel is showing the P values from the survival analysis of different datasets.

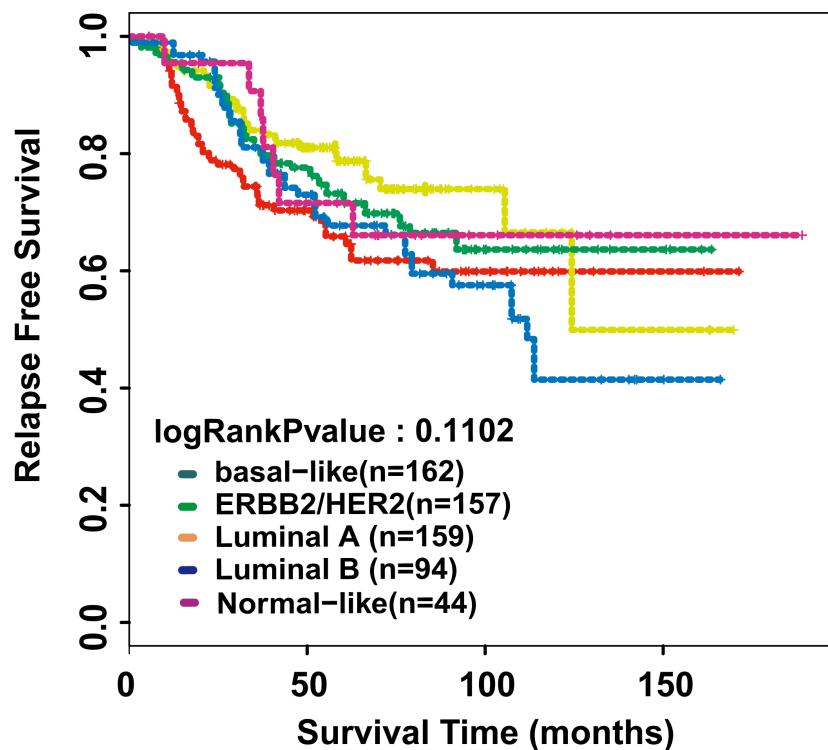


Fig. S10 Relapse-free survival curves for four groups defined by known molecular subtypes of breast cancer

Table S1 Search keywords for different cancer types

Cancer	Search keywords
Bladder cancer	(bladder cancer) OR (urothelial carcinoma of the bladder) OR (transitional cell carcinoma of urinary bladder) OR (Squamous cell carcinoma of urinary bladder) OR (adenocarcinoma of urinary bladder) OR (Small-cell carcinoma of urinary bladder) OR (bladder sarcoma)
Breast cancer	(breast cancer) OR (breast carcinoma) OR (breast ductal carcinoma)
Cervical cancer	(Uterine cancer) OR (cervical cancer) OR (endometrial carcinoma) OR (endometrial stromal sarcoma) OR (Squamous Cell Carcinoma of the Cervix) OR (Cervical Adenocarcinoma) OR (adenosquamous carcinoma of the cervix) OR (Small cell cervical cancer) OR (Small cell neuroendocrine cervical carcinoma) OR (Glass cell carcinoma of the cervix) OR (Villoglandular adenocarcinoma of the cervix)
Colorectal cancer	(rectum cancer) OR (rectal carcinoma) OR (carcinoma of rectum) OR (rectal cancer) OR (colon cancer) OR (colorectal cancer) OR (colorectal carcinoma) OR (carcinoma of colon)
Esophagus cancer	(esophageal cancer) OR (esophageal adenocarcinoma) OR (esophageal squamous cell carcinoma)
Gastric cancer	(gastric cancer) OR (Stomach cancer) OR (gastric adenocarcinoma) OR (gastrointestinal stromal tumor)
Liver cancer	(Hepatocellular carcinoma) OR (Liver cancer) OR (hepatic cancer) OR (intrahepatic cholangiocarcinom) OR (Hepatoblastoma) OR (hepatome)
Lung cancer	(lung carcinoma) OR (small-cell lung carcinoma) OR (non-small-cell lung carcinoma) OR (Lung cancer) OR (lung Adenocarcinoma) OR (lung tumor) OR (lung squamous cell carcinoma)
ovarian cancer	(ovarian cancer) OR (ovarian carcinoma) OR (ovarian epithelial carcinoma) OR (Malignant mixed müllerian tumor) OR (Mucinous tumors) OR (ovarian epithelial carcinoma) OR (Serous ovarian carcinoma) OR (Clear-cell ovarian carcinomas) OR (Clear-cell ovarian adenocarcinomas) OR (ovarian squamous cell carcinomas) OR (ovarian borderline tumor)
prostate cancer	(prostate cancer) OR (prostate carcinoma)

