

Prioritizing Cancer lncRNA Modulators via Integrated lncRNA-mRNA Network and Somatic Mutation Data

Dianshuang Zhou

Harbin Medical University <https://orcid.org/0000-0002-9205-4641>

Xin Li

Harbin Medical University

Shipeng Shang

Harbin Medical University

Hui Zhi

Harbin Medical University

Peng Wang

Harbin Medical University

Yue Gao

Harbin Medical University

Shangwei Ning (✉ ningsw@ems.hrbmu.edu.cn)

Harbin Medical University

Research Article

Keywords: computational framework, lncRNA-mRNA network, mutations, cancer lncRNA modulators

Posted Date: June 22nd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-636414/v1>

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

[Read Full License](#)

Version of Record: A version of this preprint was published at Current Bioinformatics on April 21st, 2022.
See the published version at <https://doi.org/10.2174/1574893617666220421095601>.

1 **Prioritizing cancer lncRNA modulators via integrated**
2 **lncRNA-mRNA network and somatic mutation data**

3 Dianshuang Zhou¹⁺, Xin Li¹, Shipeng Shang¹, Hui Zhi¹, Peng Wang¹, Yue Gao¹, Shangwei Ning^{1*}

4 ¹ College of Bioinformatics Science and Technology, Harbin Medical University, 194 Xuefu Road,
5 Harbin 150081, China

6

7

8 **ABSTRACT:**

9 **Background:** Long noncoding RNAs (LncRNAs) represent a large category of
10 functional RNA molecules that play a significant role in human cancers. lncRNAs can
11 be genes modulators to affect the biological process of multiple cancers.

12 **Methods:** Here, we developed a computational framework that uses lncRNA-
13 mRNA network and mutations in individual genes of 9 cancers from TCGA to prioritize
14 cancer lncRNA modulators. Our method screened risky cancer lncRNA regulators
15 based on integrated multiple lncRNA functional networks and 3 calculation methods in
16 network.

17 **Results:** Validation analyses revealed that our method was more effective than
18 prioritization based on a single lncRNA network. This method showed high predictive
19 performance and the highest ROC score was 0.836 in breast cancer. It's worth noting
20 that we found that 5 lncRNAs scores were abnormally high and these lncRNAs
21 appeared in 9 cancers. By consulting the literatures, these 5 lncRNAs were
22 experimentally supported lncRNAs. Analyses of prioritizing lncRNAs reveal that these
23 lncRNAs are enriched in various cancer-related biological processes and pathways.

24 **Conclusions:** Together, these results demonstrated the ability of this method
25 identifying candidate lncRNA molecules and improved insights into the pathogenesis
26 of cancer.

27 **Keywords:** computational framework, lncRNA-mRNA network, mutations, cancer
28 lncRNA modulators,

30

31 **Introduction:**

32 Cancer is a major public health problem across the world and is a leading cause of death
33 in many countries(1). Cancer is a complex disease involving DNA abnormalities,
34 transcriptomic alterations and epigenetic aberrations(2) and whole genome sequencing
35 efforts have uncovered the genomic landscapes of common forms of human cancers(3).
36 The cancer Genome atlas (TCGA) has provided a mass of data of human samples and
37 discover molecular alterations at the DNA and RNA levels(4).

38 Mutations are important markers of cancer genes and the somatic mutation
39 landscapes and signatures of major cancer types have been reported and stockpiled by
40 international cancer genome projects, such as TCGA and ICGC(5). In recent years,
41 increasing experimentally supported evidence has suggested that lncRNAs as genes
42 modulators affect the process of cancers. For instance, in lung cancer, MALAT1
43 actively regulates a set of metastasis-associated genes expression, including MIA2,
44 HNF4G and CA2. Moreover, MALAT1 can be used as a valuable prognostic marker
45 and a promising therapeutic target(6). TUG1 can regulate the expression of LIMK2b
46 and then promoted cell growth and chemoresistance of small cell lung cancer(7).
47 NEAT1 transforms the epigenetic landscape of target gene promoters to facilitate
48 transcription and promotes carcinogenic growth(8). These data indicate that researching
49 lncRNAs is an important step in the understanding of cancer mechanisms. However,
50 there is currently no systematic way to explore the lncRNAs which can be genes
51 modulators in multiple cancers.

52 Network analysis is often used to explore the function of lncRNAs and the
53 relationship between lncRNAs and diseases. For instance, Guo et al. developed a bi-
54 colored network to annotate lncRNA function(9). Currently, Data sets obtained from
55 many lncRNA-related databases can be used for in-depth exploration of lncRNA.
56 TANRIC platform systematically collected data resource that records the expression of
57 lncRNA in 20 human cancers(10). starBase v2.0, RAID V2.0 and NPInter V3.0 stored
58 a huge amount of lncRNA data resources, including lncRNA-mRNA and lncRNA-
59 miRNA interactions(11-13). LncRNADisease and Lnc2Cancer databases manually
60 curated credible lncRNA-disease or lncRNA-cancer data, meanwhile, the dysfunction
61 pattern of lncRNAs were annotated (14, 15). However, our knowledge of using these
62 public data to build a network of lncRNAs to further explore cancer lncRNA modulators
63 remains limited.

64 Here we exploit a cancer lncRNA prioritization computational method based on a
65 lncRNA-mRNA function network and somatic mutation data for 9 cancers (BLCA,
66 BRCA, COAD, GBM, LIHC, LUAD, OV, PRAD, STAD) from TCGA. This study
67 systematically analyzes the result of method, including topological property of our
68 functional network, the performance of method and functional enrichment analysis and
69 biological characteristics of risky cancer lncRNA modulators. We found that most of
70 top ranking lncRNA modulators have been identified cancer lncRNAs by Lnc2Cancer
71 database, then these lncRNA modulators also enriched in cancer-related Gene Ontology
72 (GO) terms and pathways. Our method can identify the lncRNA modulators which
73 cannot be identified by differential expression and can be a valuable complement to

74 experimental studies used in future studies.

75 **MATERIALS AND METHODS**

76 **Data sources**

77 **Cancer somatic mutation data and golden standard**

78 We used 9 types of cancer somatic mutation data from TCGA: Bladder Urothelial
79 Carcinoma (BLCA), Breast invasive carcinoma (BRCA), Colon adenocarcinoma
80 (COAD), Glioblastoma multiforme (GBM), Liver hepatocellular carcinoma (LIHC),
81 Lung adenocarcinoma (LUAD), Ovarian serous cystadenocarcinoma (OV), Stomach
82 adenocarcinoma (STAD).

83 The golden standard of gene set was obtained from Cancer Genome Census(CGC)
84 database, which includes 616 cancer genes(16). In order to verify the accuracy of the
85 model, we choose the 9 cancers with the most experimentally confirmed data in the
86 Lnc2Cancer database, which includes 148 lncRNAs for 9 cancers.

87 **The lncRNA-mRNA functional network**

88 In our study, the lncRNA-mRNA functional network was a fusion of lncRNA-mRNA
89 co-expression network, lncRNA-mRNA ceRNA network and lncRNA-protein
90 interaction network. We have used gene expression data from TCGA and lncRNA
91 expression data from TANRIC. Pearson correlation was used to construct our lncRNA-
92 mRNA co-expression network and the lncRNA-mRNA pair was selected if it meets
93 following criteria: $\text{corr}(\text{lncRNA}, \text{mRNA}) > 0.8$, $\text{fdr} < 0.05$. We used lncRNA-miRNA
94 pairs and mRNA-miRNA pairs to construct our lncRNA-mRNA ceRNA network.
95 lncRNA-miRNA interaction data and mRNA-miRNA interaction data were
96 downloaded from starBase v2.0, NPInter v3.0 and RAID v2.0. According to number of

97 miRNA shared with lncRNA and mRNA, we used hypergeometric distribution to
98 construct our lncRNA-mRNA ceRNA network. For lncRNA-protein interaction
99 network, we also downloaded lncRNA-protein interaction data from starBase v2.0,
100 NPInter v3.0 and RAID v2.0. The weight of lncRNA-mRNA ceRNA network and
101 lncRNA-protein interaction network was defined as 1. Then we coalesced 3 network
102 and defined the weight of lncRNA-mRNA pairs as $(\text{corr}_{lm} + W_{\text{celm}} + W_{\text{lplm}})/3$, where
103 corr_{lm} , W_{celm} , W_{lplm} indicated Pearson correlations for lncRNA-mRNA, the weight
104 in ceRNA network and the weight in interaction network, respectively.

105 **Scoring scheme of genes**

106 We defined N_i for the number of non-synonymous mutations of a gene from the somatic
107 mutation data. Meanwhile, we screened the differential expression genes according to
108 the gene expression data and got the *P-value* of Student's t-test of each gene. The *P-*
109 *value* was transformed by a standardization method which named *scale function* in R:

$$110 \quad np_i = 1 - \frac{1}{1 + e^{\text{scale}(-\log_{10}(p_i))}}$$

111 Where np_i represents the normalization score of differential expression gene(i), p_i
112 represents the P-value of differential expression gene(i).

113 We formulated the score of gene to use mutation occurrences and np_i :

$$114 \quad G_i = M_i * (1 + np_i)$$

115 Where G_i represents the final score of gene(i), M_i represents the mutational
116 occurrences of gene(i).

117 **Scoring scheme of lncRNAs**

118 We designed three different ways to use gene's score of direct neighbors and edge
119 weights in our network for prioritizing lncRNAs: the first computational method named

120 “Smax” was defined that lncRNA’s score is the biggest score of its direct gene’s score
121 multiply by edge weights:

$$122 \quad L_j = \max(G_i * w(i,j))$$

123 the second computational method named “Ssum” was defined that lncRNA’s score is
124 the sum score of its direct gene’s score multiply by edge weights:

$$125 \quad L_j = \sum_{i=1}^N G_i * w(i,j)$$

126 Where L_i represents the score of lncRNA(j), $w(i,j)$ represents the edge weights of
127 lncRNA(j) and gene(i).

128 the third computational method named “NWsum” was defined that lncRNA’s score is
129 the sum score of its direct gene’s score divide by the number of gene’s direct neighbors:

$$130 \quad L_j = \sum_{i=1}^N \frac{G_i}{n_i}$$

131 Where n_i represents the number of gene(i)’s direct neighbors.

132 For example, if a lncRNA had 5 direct neighbors and the number of 5 genes’ direct were
133 2,4,1,1,3, we can obtain the score for this lncRNA:

$$134 \quad L = \frac{G(1)}{2} + \frac{G(2)}{4} + \frac{G(3)}{1} + \frac{G(4)}{1} + \frac{G(5)}{3}$$

135 **A summary of validated cancer-related lncRNAs**

136 Literature mining is an effective way to collect “gold-standard” for a large number of
137 disease-related molecules because of experimental methods, such as Western blot,
138 Luciferase reporter assay. In this study, we used a set of validated lncRNAs for 9
139 cancers from Lnc2Cancer database(<http://www.bio-bigdata.net/lnc2cancer/>), which
140 contains cancer-related lncRNAs based on experiment by thousands of articles. Due to
141 the emergence of a lot of new data after the database published, we added the latest data

142 to test our method through manually collecting lncRNA-cancer associations.

143 **Calculation of network topology property and survival analysis**

144 Degree centrality and betweenness centrality are two important indicators in the nature
145 of network topology. Generally, the larger the node degree of a node is, the higher the
146 degree of centrality of the node is, and the more important the node is in the network;
147 betweenness centrality is equal to the number of shortest paths from each node to all
148 others that pass through this node, as an important global geometric quantity,
149 betweenness centrality reflects the role and influence of the corresponding node in the
150 entire network. Degree centrality and betweenness centrality was calculated using the
151 R package “igraph” (17).

152 A Kaplan-Meier survival analysis was performed using the clinical data from
153 TCGA, and statistical significance was assessed using the log-rank test. The survival
154 curve was drawn using the R package "survival". All analyses were performed on the
155 R 3.6.0 framework.

156

157 **Result**

158 **Overview**

159 A general workflow of method is given in Fig. 1. To prioritize lncRNA molecules, the
160 first step was to score genes based on somatic mutation data and gene expression data.
161 The second step was to integrate three lncRNA-mRNA networks: lncRNA-mRNA co-
162 expression, lncRNA-mRNA ceRNA network, lncRNA-protein interaction network.
163 The third step was to score lncRNAs based on genes' score and integrated lncRNA-
164 mRNA network. We used three methods to score lncRNAs and then sorted them
165 according to lncRNAs' score. The higher the ranking, the more likely it is to become a
166 risk cancer lncRNA.

167

168 **Global properties of the lncRNA-mRNA function network**

169 There was an average of 378768 edges, including 19883 coding genes and 12150
170 lncRNAs in our network (Fig. 2A). In this network, we first calculated the degree
171 centrality and found a small number of nodes had very high degree. The degree
172 distribution of the network obeys a power law distribution (Fig. 2B). Then we compared
173 the degree and betweenness (Fig. 2C) of cancer lncRNA with candidate lncRNA. The
174 degree and betweenness of cancer lncRNA were significantly higher than the candidate
175 lncRNA and the P-value by Wilcoxon test was less than 0.001. These results indicated
176 "cancer lncRNA nodes" is a key factor in the network and plays a regulatory role for a
177 large number of genes, which is consistent with the previous research results(18).

178 **Performance of method**

179 First, we use genes' score which were calculated based on somatic mutation data and
180 gene expression data for 9 cancers from TCGA and a protein-protein interaction
181 network which named STRING v10 to appraise the method when used to prioritizing
182 cancer genes. The golden standard was a high-confidence gene set form the Cancer
183 Genome Census database (CGC), including 616 cancer driver genes.

184 To assess the performance of our method, ROC analysis was executed for each
185 type of cancer and the AUC value used to determine the quality of the method. For
186 instance, the AUC values which were calculated based on 3 methods and cancer genes
187 annotated by CGC were 0.909, 0.896, 0.843, respectively (Fig. 3A). This indicated that
188 our method showed high predictive performance when using PPI network and mutated
189 genes to predict driving genes.

190 For prioritizing cancer lncRNA, we used genes' score and lncRNA-mRNA
191 functional network. To evaluate the performance of our method, we select the
192 experimentally supported cancer lncRNAs as the golden standard from Lnc2Cancer
193 database. We first compared the values of AUC using a single subnet with the integrated
194 network, finding that the values of AUC which used lncRNA-mRNA co-expression
195 network significantly below the values of AUC which used lncRNA-mRNA ceRNA
196 network and lncRNA-protein network. Meanwhile, the values of AUC which used
197 integrated network showed the best performance (Fig. 3B). This result showed that
198 using a single network for scoring lncRNA may not be the best method and consolidated
199 multiple networks can improve prediction performance. We also compared the
200 performance of the three methods, finding that all three methods show high predictive

201 performance. The median values of AUC for 3 methods were more than 0.77. NWsum
202 algorithm got the highest AUC values and the median values of AUC was 0.822 (Fig.
203 3C).

204 The robustness analysis was performed by deleting some network nodes. We
205 randomly deleted 10% and 20% of the network nodes and their connected edges. Using
206 incomplete networks and NWsum algorithm, we recalculate the AUC values, finding
207 that the AUC values of all cancer had no prominent change. The most variable cancer
208 was Glioblastoma multiforme (GBM), however, the level of change did not exceed 10%.
209 These results showed that although the absence of the network would have a certain
210 impact on the performance of the method, the effect was small (Fig. 3D). This was
211 because that the key nodes in the network had very high degrees and betweenness, even
212 if a part of the network was missing, it would not affect the result. All the results
213 indicated our network was robust.

214 **Analysis of high-risk cancer lncRNA modulators**

215 Duo to the performance of NWsum algorithm was slightly stronger than the other two
216 algorithms, we use the lncRNA predicted by the NWsum algorithm to analyze.

217 **Top-rank lncRNAs consist with known cancer lncRNAs**

218 Currently, there have been some studies of cancer related lncRNAs by experimental
219 methods, although the number is not large, it is still the most powerful standard to verify
220 the performance of our method. According to the scores calculated by the NWsum
221 algorithm, we found that 5 lncRNAs scores were abnormally high and the top 5
222 lncRNAs which were MALAT1, NEAT1, FENDRR, CRNDE, TUG1 all appeared in 9
223 cancers. By consulting the literatures, these five lncRNAs were experimentally

224 supported lncRNAs and associated with the development of multiple cancers
225 (Supplementary Table S1). LncRNA MALAT1 had been confirmed with all cancers we
226 researched and many pathogenic mechanisms had been described, including
227 influencing the metastasis phenotype of lung cancer cells(6), inducing epithelial-to-
228 mesenchymal transition due to the dysfunction of MALAT1 in bladder cancer(19),
229 promoting tumor growth in colorectal cancer(20) and so on. So far, the relationships
230 between LncRNA FENDRR and cancers were not so much. It had been confirmed that
231 decreased expression of FENDRR regulates gastric cancer cell metastasis(21) and
232 FENDRR as a tumor suppressor gene in NSCLC inhibits cell proliferation and induces
233 apoptosis. Meanwhile, researcher found FENDRR was able to regulate heart and body
234 wall development in the mouse(22). So lncRNAs FENDRR may prove to be risk factors
235 for other cancers in future studies.

236 **Common character and special character cancer lncRNA modulators**

237 Previous studies have shown that some lncRNAs function in many types of cancer, and
238 some lncRNAs have certain tumor specificities. We chose the top 20 lncRNAs and their
239 corresponding scores (excluding the 5 lncRNAs mentioned above) and standardized
240 scores for each cancer. At last, we got 75 lncRNAs and their standardized scores for 9
241 cancers. Based on the score of these lncRNAs, a global risk-evaluation score profile
242 was constructed and clustered by the k-means method (Fig.4A). 18 lncRNAs scored
243 higher in 9 cancers and 13 lncRNAs only scored higher in colon adenocarcinoma
244 (COAD). To characterize the function of these lncRNAs, a gene set enrichment analysis
245 web server called “Enrichr” was used for pathway annotations, including Gene

246 Ontology biological process, KEGG pathway, Reactome pathway, WikiPathways(23).
247 Common character lncRNAs were annotated the pathways which associated with
248 multiple cancers. The biological process of cell (mitochondrion localization and
249 apoptotic cell clearance), known cancer pathway (WNT pathway) and cancer-related
250 protein (SMO and GTP) were highly represented. To the special character cancer
251 lncRNAs, BMP signaling pathway which was a developmental pathway and a potential
252 therapeutic target was been enriched (Fig.4B). Specifically, these lncRNAs enriched to
253 the corresponding diseases (Neoplasm of the colon, intestinal polyposis and Neoplasm
254 of the rectum) by using a tool named Human Phenotype Ontology (Fig.4C).

255 In Common character lncRNAs, 5 lncRNAs (ZNF518A, HCG18, FGD5-AS1,
256 TSIX, RP11-553L6.5) were not verified by literatures to this day. We extracted the
257 interaction genes and miRNAs of these 5 lncRNAs and compared with the confirmed
258 evidence. 399 miRNAs and 2855 genes were extracted, respectively. The
259 experimentally supported miRNAs for 9 cancers were from HMDD 2.0 database(24)
260 and genes for 9 cancers were from Cancer Genome Census database (CGC). The
261 number of experimentally supported miRNAs and genes were 134 and 616 (Fig.4D).
262 We found 96 miRNAs and 166 genes exhibited both interaction of association and
263 confirmed evidence, including miR-106a, let-7b, FUS, EWS, etc. (Fig.4E) High
264 overlap indicated they had a significant association with 9 cancers ($P < 0.001$, $P < 0.001$).
265 We thus speculated that these 5 lncRNAs were likely to be a high-risk clinical factor
266 and further experiments need to be carried out.

267 **Biological characteristics of prioritizing lncRNA modulators**

268 lncRNAs are diverse and involved in a variety of biological processes. In the process
269 of transcription, some special features of lncRNA are formed and these features are
270 closely linked with the function of lncRNA. We divided the top 2000 lncRNAs into 3
271 levels (Lv1: 1~50, Lv2: 51~500, Lv3: 501~2000) to observe the differences in
272 biological characteristics. Previous studies have shown that exons are one of the
273 important genomic characteristics of lncRNAs(25). The average number of exons for
274 Lv1 was 1.74-fold that of Lv2(Wilcoxon test, $P = 0.039$) and the average number of
275 exons for Lv2 was 1.78-fold that of Lv3($P < 0.001$) (Fig.5A). One of the important
276 functions of lncRNA is its ability which bind to miRNA competitively to regulate gene
277 expression. The number of miRNA binding sites on lncRNA reflects the ability of
278 lncRNAs to bind miRNAs. The density of miRNA binding sites was computed by two
279 prediction methods: miRanda and TargetScan. The average number of miRNA binding
280 sites for Lv1 was 1.42-fold that of Lv2 ($P < 0.001$) and the average number of miRNA
281 binding sites for Lv2 was 1.32-fold that of Lv3 ($P < 0.001$) (Fig.5B). Previous research
282 shows some specific human lncRNAs which different evolutionary conservation
283 beyond primates but have proven to be both functional and therapeutically relevant(26).
284 The UCSC phyloP score was used to calculate the conservation score of lncRNAs. Duo
285 to the conserved state belongs to the corresponding coding gene instead of lncRNA, we
286 excluded antisense lncRNAs when we calculated the conservation scores. The average
287 scores for Lv1 were 3.2-fold that of Lv2 ($P = 0.003$) and the average scores for Lv2 was
288 1.42-fold that of Lv3 ($P = 0.013$) (Fig.5C). Some famous lncRNA, such as HOTAIR
289 and PCAT7, were significantly upregulated in various cancers. We got the lncRNA

290 expression profile from TANRIC database and calculated the expression values of 3
291 levels. The average expression values for Lv1 were 13.2-fold that of Lv2 ($P < 0.001$)
292 and the average expression values for Lv2 was 3.29-fold that of Lv3 ($P < 0.001$)
293 (Fig.5D). In general, top lncRNAs have more significant biological characteristics than
294 others candidate lncRNAs.

295 **Case study: Stomach Cancer**

296 Stomach cancer is one of the most common cancers in the digestive system and the
297 second most common cause of cancer death in the world. Researchers found some
298 lncRNAs affected the occurrence and development of stomach cancer. Hao Li et.al.
299 presented overexpression of H19 promoted the features of GC including proliferation,
300 migration, invasion and metastasis(27). Yongchao Liu et.al. presented lncRNA
301 GAS5/YBX1/p21 pathway may become a useful therapeutic method since the YBX1
302 protein level was reduced by down-regulation of GAS5 and decreased p21 expression,
303 thus abolished G1 phase cell cycle arrest in stomach cancer(28). So, we propose a case
304 study of stomach cancer to investigate whether our methods can find lncRNA molecules
305 of stomach cancer. We selected the top 20 lncRNAs for sample analysis.

306 First, we compared the differential expression lncRNAs in samples of stomach
307 cancer by using Student's t-test with the top 20 lncRNAs which we selected
308 (Supplementary Table S2). Thirteen of the top 20 lncRNAs were differential expression.
309 LncRNA TINCR which cannot be identified by differential expression was identified
310 by our method. This indicated our method could identify cancer lncRNAs that were not
311 identified by differential expression method.

312 For the top 20 lncRNAs that were not be verified by experiment, we use the way
313 of literature retrieval to mine their functions one by one. ZNF518A was recorded in top
314 20 Co-Expressed gene's Protein-protein interaction Network Plot of Cancer Cell
315 Metabolism Gene DB(29). LINC00657 could suppresses hepatocellular carcinoma cell
316 growth and breast cancer cell growth(30, 31). HCG18 may serve as a potential
317 biomarker in breast cancer and lung cancer and as an immune-related lncRNA to affect
318 the survival of glioma patients(32, 33). XIST and TSIX, as cancer-immune biomarkers,
319 manipulated PD-L1 expression in BC cell lines, leading to lymph node metastasis. For
320 other lncRNAs that were not be verified by experiment, there was currently no content
321 related to tumor or tumor pathogenesis. They may become new directions for future
322 research.

323 To verify whether the top 20 lncRNAs are closely related to gastric cancer, a
324 functional annotation of these lncRNAs was carried out by “Enrichr” web server. The
325 results manifested these lncRNAs enriched many critical biological functions and
326 pathways, including S33 mutants of beta-catenin, Beta-catenin phosphorylation
327 cascade, Toll Like Receptor signaling, TGF beta signaling pathway, etc. The result of
328 functional annotation proved these lncRNAs affected the occurrence and development
329 of stomach cancer.

330 Furthermore, we investigated whether top 20 lncRNAs can act as an independent
331 prognostic factor. Clinical information for stomach Cancer was obtained from TCGA.
332 We used lncRNA expression data and clinical information to performed a Kaplan–
333 Meier survival analysis. Then we obtained 2 lncRNAs which significantly related to the

334 patient's survival (NEAT1, $P = 0.0309$; LINC00657, $P = 0.0296$). Our results displayed
335 that high or low expression of lncRNA could significantly divide patients into two
336 categories based on survival time. Drawing a comparison between the miRNAs and
337 genes which were interaction with the 2 lncRNAs and stomach cancer related miRNAs
338 and genes, many experimentally confirmed gastric cancer miRNAs and genes appear
339 in miRNAs and genes that interact with the 2 lncRNAs ($P < 0.001$). The stomach cancer
340 miRNAs and genes linked by the 2 lncRNAs are displayed using Cytoscape (34).

341 In summary, our method can accurately identify stomach cancer-associated
342 lncRNA molecules, furthermore, our method can identify lncRNAs which can't be
343 identified by the method of differential expression.

344

345 **Discussion**

346 lncRNA is a regulator of gene expression to act the translation of genetic codons into
347 protein sequences (35). In previous studies, researchers were more inclined to develop
348 methods which explore driver coding genes for cancers, such as PolyPhen-2,
349 MutSigCV and HotNet2. However, lncRNA as a key factor that can influence tumor
350 initiation, growth and metastasis of cancer cells (36), a few methods which prioritize
351 candidate cancer lncRNA molecules were exploited. At the same time, there were
352 currently not many data for detecting lncRNA mutation sites, so we attempted to use
353 the data from somatic mutations of TCGA to screen for cancer risk lncRNAs. We first
354 utilized expression data and public lncRNA, miRNA, gene interaction data to construct
355 our lncRNA-mRNA network. Then the gene's score was obtained by combining the
356 occurrences of somatic mutations and the differential expression level. Here, we present
357 an approach for identifying cancer lncRNA molecules using our lncRNA-mRNA
358 functional network and three computing method in direct neighbors for nine cancers,
359 providing a view for predicting cancer lncRNAs.

360 Using a co-expression network alone could not obtain good predictive results. We
361 found that the correlation coefficient of co-expression between most experiments
362 supported lncRNA and gene were not very high. In fact, lncRNA-mRNA co-expression
363 does not mean that the two have a regulatory relationship. When we choose a higher
364 Pearson correlation coefficient, the experimentally supported lncRNAs will be filtered
365 out in lncRNA-mRNA co-expression network, so we join the lncRNA-mRNA ceRNA
366 network and the lncRNA-protein interaction network. The ROC analysis results

367 indicated using lncRNA-protein interaction network could get the best result in three
368 networks. This shows that the experimentally supported interaction network has higher
369 accuracy than the network obtained by the calculation method. In our integrated
370 network, the degree and median of experimentally supported cancer lncRNAs were
371 significantly higher than other lncRNAs.

372 In our predicted results, in addition to using ROC analysis to evaluate the overall
373 prediction results, we also chose the top ranked lncRNA for analysis. By comparing
374 biological characteristics, lncRNA ranks and characteristic scores had a consistent trend.
375 For more accurate analysis, we identified 5 lncRNAs (MALAT1, NEAT1, FENDRR,
376 CRNDE, TUG1) which ranked in the top 5 of 9 cancers. These 5 lncRNAs had been
377 documented in many literatures and were closely related to various cancers. In STAD
378 study, we found this method make accurate and complement lncRNA found by
379 differential expression as the high false positives and missing parts cancer lncRNA of
380 the differential expression method. A mass of experimentally supported cancer genes
381 and cancer miRNAs which connect with our predicted lncRNA also show the accuracy
382 of our method.

383 Currently, experimentally supported cancer lncRNAs were small number. Some
384 lncRNAs have been confirmed in common cancers only. Using golden standard to
385 evaluate predictions is an effective method, so we select 9 cancers which have more
386 confirmed lncRNAs. With the continuous improvement of experimental technology, an
387 increasing number of confirmed cancer lncRNAs will be found and our method will
388 also apply to other cancers. Patients with the same tumor may have different reactions

389 using the same drug. The relationship of lncRNA and drugs is also the research direction
390 of many researchers. Ehsan Malek et al. present lncRNA expression is associated with
391 drug resistance and base on the function and structure of lncRNA to avoid drug
392 resistance by using appropriate drugs(37). In future research, we make an attempt to
393 use a data from a patient's somatic mutation to obtain more effective therapies.
394

395 **Ethics approval and consent to participate**

396 Not applicable

397 **Authors' contributions**

398 Conception and design of the work, SW Ning and DS Zhou; acquisition and analysis
399 of data, DS Zhou; writing, reviewing and editing the paper, DS Zhou, X Li, SP Shang,
400 H Zhi, P Wang, Y Gao. All authors have read and agreed to the published version of the
401 manuscript and to have agreed to both be personally accountable for the author's
402 contributions and ensure to answer any questions related to the accuracy or integrity of
403 any part of the work. All authors read and approved the final manuscript.

404 **Funding**

405 This work was supported by National Natural Science Foundation of China [32070672],
406 Postdoctoral Scientific Research Development Fund, and Yu Weihai Outstanding
407 Youth Training Fund of Harbin Medical University.

408 **Competing interests**

409 The authors declare no conflicts of interest regarding this work.

410 **Acknowledgements**

411 Not applicable

412 **Availability of data and material**

413 Not applicable

414 **Consent for publication**

415 Not applicable

416

417 **References**

- 418 1. Siegel RL, Miller KD, Jemal A. Cancer statistics, 2020. *CA: a cancer journal for clinicians.*
419 2020;70(1):7-30.
- 420 2. Hanahan D, Weinberg RA. Hallmarks of cancer: the next generation. *Cell.* 2011;144(5):646-74.
- 421 3. Vogelstein B, Papadopoulos N, Velculescu VE, Zhou S, Diaz LA, Jr., Kinzler KW. Cancer genome
422 landscapes. *Science.* 2013;339(6127):1546-58.
- 423 4. Cancer Genome Atlas Research N, Weinstein JN, Collisson EA, Mills GB, Shaw KR, Ozenberger BA,
424 et al. The Cancer Genome Atlas Pan-Cancer analysis project. *Nature genetics.* 2013;45(10):1113-20.
- 425 5. Cheng F, Zhao J, Zhao Z. Advances in computational approaches for prioritizing driver mutations
426 and significantly mutated genes in cancer genomes. *Briefings in bioinformatics.* 2016;17(4):642-56.
- 427 6. Gutschner T, Hammerle M, Eissmann M, Hsu J, Kim Y, Hung G, et al. The noncoding RNA MALAT1
428 is a critical regulator of the metastasis phenotype of lung cancer cells. *Cancer Res.* 2013;73(3):1180-9.
- 429 7. Niu Y, Ma F, Huang W, Fang S, Li M, Wei T, et al. Long non-coding RNA TUG1 is involved in cell
430 growth and chemoresistance of small cell lung cancer by regulating LIMK2b via EZH2. *Molecular cancer.*
431 2017;16(1):5.
- 432 8. Chakravarty D, Sboner A, Nair SS, Giannopoulou E, Li R, Hennig S, et al. The oestrogen receptor
433 alpha-regulated lncRNA NEAT1 is a critical modulator of prostate cancer. *Nat Commun.* 2014;5:5383.
- 434 9. Guo X, Gao L, Liao Q, Xiao H, Ma X, Yang X, et al. Long non-coding RNAs function annotation: a
435 global prediction method based on bi-colored networks. *Nucleic acids research.* 2013;41(2):e35.
- 436 10. Li J, Han L, Roebuck P, Diao L, Liu L, Yuan Y, et al. TANRIC: An Interactive Open Platform to Explore
437 the Function of lncRNAs in Cancer. *Cancer Res.* 2015;75(18):3728-37.
- 438 11. Li JH, Liu S, Zhou H, Qu LH, Yang JH. starBase v2.0: decoding miRNA-ceRNA, miRNA-ncRNA and
439 protein-RNA interaction networks from large-scale CLIP-Seq data. *Nucleic acids research.*
440 2014;42(Database issue):D92-7.
- 441 12. Yi Y, Zhao Y, Li C, Zhang L, Huang H, Li Y, et al. RAID v2.0: an updated resource of RNA-associated
442 interactions across organisms. *Nucleic acids research.* 2017;45(D1):D115-D8.
- 443 13. Hao Y, Wu W, Li H, Yuan J, Luo J, Zhao Y, et al. NPInter v3.0: an upgraded database of noncoding
444 RNA-associated interactions. *Database : the journal of biological databases and curation.* 2016;2016.
- 445 14. Chen G, Wang Z, Wang D, Qiu C, Liu M, Chen X, et al. lncRNADisease: a database for long-non-
446 coding RNA-associated diseases. *Nucleic acids research.* 2013;41(Database issue):D983-6.
- 447 15. Ning S, Zhang J, Wang P, Zhi H, Wang J, Liu Y, et al. lnc2Cancer: a manually curated database of
448 experimentally supported lncRNAs associated with various human cancers. *Nucleic acids research.*
449 2016;44(D1):D980-5.
- 450 16. Futreal PA, Coin L, Marshall M, Down T, Hubbard T, Wooster R, et al. A census of human cancer
451 genes. *Nature reviews Cancer.* 2004;4(3):177-83.
- 452 17. Mora A, Donaldson IM. iRefR: an R package to manipulate the iRefIndex consolidated protein
453 interaction database. *BMC bioinformatics.* 2011;12:455.
- 454 18. Schmitt AM, Chang HY. Long Noncoding RNAs in Cancer Pathways. *Cancer cell.* 2016;29(4):452-63.
- 455 19. Ying L, Chen Q, Wang Y, Zhou Z, Huang Y, Qiu F. Upregulated MALAT-1 contributes to bladder cancer
456 cell migration by inducing epithelial-to-mesenchymal transition. *Molecular bioSystems.* 2012;8(9):2289-
457 94.
- 458 20. Ji Q, Zhang L, Liu X, Zhou L, Wang W, Han Z, et al. Long non-coding RNA MALAT1 promotes tumour
459 growth and metastasis in colorectal cancer through binding to SFPQ and releasing oncogene PTBP2 from
460 SFPQ/PTBP2 complex. *British journal of cancer.* 2014;111(4):736-48.

- 461 21. Xu TP, Huang MD, Xia R, Liu XX, Sun M, Yin L, et al. Decreased expression of the long non-coding
462 RNA FENDRR is associated with poor prognosis in gastric cancer and FENDRR regulates gastric cancer
463 cell metastasis by affecting fibronectin1 expression. *Journal of hematology & oncology*. 2014;7:63.
- 464 22. Grote P, Wittler L, Hendrix D, Koch F, Wahrisch S, Beisaw A, et al. The tissue-specific lncRNA Fendrr
465 is an essential regulator of heart and body wall development in the mouse. *Developmental cell*.
466 2013;24(2):206-14.
- 467 23. Kuleshov MV, Jones MR, Rouillard AD, Fernandez NF, Duan Q, Wang Z, et al. Enrichr: a
468 comprehensive gene set enrichment analysis web server 2016 update. *Nucleic acids research*.
469 2016;44(W1):W90-7.
- 470 24. Li Y, Qiu C, Tu J, Geng B, Yang J, Jiang T, et al. HMDD v2.0: a database for experimentally supported
471 human microRNA and disease associations. *Nucleic acids research*. 2014;42(Database issue):D1070-4.
- 472 25. Zhang YC, Liao JY, Li ZY, Yu Y, Zhang JP, Li QF, et al. Genome-wide screening and functional analysis
473 identify a large number of long noncoding RNAs involved in the sexual reproduction of rice. *Genome
474 biology*. 2014;15(12):512.
- 475 26. Johnsson P, Lipovich L, Grander D, Morris KV. Evolutionary conservation of long non-coding RNAs;
476 sequence, structure, function. *Biochimica et biophysica acta*. 2014;1840(3):1063-71.
- 477 27. Li H, Yu B, Li J, Su L, Yan M, Zhu Z, et al. Overexpression of lncRNA H19 enhances carcinogenesis
478 and metastasis of gastric cancer. *Oncotarget*. 2014;5(8):2318-29.
- 479 28. Liu Y, Zhao J, Zhang W, Gan J, Hu C, Huang G, et al. lncRNA GAS5 enhances G1 cell cycle arrest via
480 binding to YBX1 to regulate p21 expression in stomach cancer. *Scientific reports*. 2015;5:10159.
- 481 29. Kim P, Cheng F, Zhao J, Zhao Z. ccmGDB: a database for cancer cell metabolism genes. *Nucleic acids
482 research*. 2016;44(D1):D959-68.
- 483 30. Hu B, Cai H, Zheng R, Yang S, Zhou Z, Tu J. Long non-coding RNA 657 suppresses hepatocellular
484 carcinoma cell growth by acting as a molecular sponge of miR-106a-5p to regulate PTEN expression.
485 *The international journal of biochemistry & cell biology*. 2017;92:34-42.
- 486 31. Liu H, Li J, Koirala P, Ding X, Chen B, Wang Y, et al. Long non-coding RNAs as prognostic markers in
487 human breast cancer. *Oncotarget*. 2016;7(15):20584-96.
- 488 32. Ding X, Zhang Y, Yang H, Mao W, Chen B, Yang S, et al. Long non-coding RNAs may serve as
489 biomarkers in breast cancer combined with primary lung cancer. *Oncotarget*. 2017;8(35):58210-21.
- 490 33. Wang W, Zhao Z, Yang F, Wang H, Wu F, Liang T, et al. An immune-related lncRNA signature for
491 patients with anaplastic gliomas. *Journal of neuro-oncology*. 2018;136(2):263-71.
- 492 34. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al. Cytoscape: a software
493 environment for integrated models of biomolecular interaction networks. *Genome research*.
494 2003;13(11):2498-504.
- 495 35. Quinn JJ, Chang HY. Unique features of long non-coding RNA biogenesis and function. *Nat Rev
496 Genet*. 2016;17(1):47-62.
- 497 36. Yang G, Lu X, Yuan L. lncRNA: a link between RNA and cancer. *Biochimica et biophysica acta*.
498 2014;1839(11):1097-109.
- 499 37. Malek E, Jagannathan S, Driscoll JJ. Correlation of long non-coding RNA expression with metastasis,
500 drug resistance and clinical outcome in cancer. *Oncotarget*. 2014;5(18):8027-38.

501

502

503 **Figure legends**

504 **Fig.1** The workflow for prioritizing cancer lncRNA modulators. Step 1: Score the
505 mutated gene. Step 2: Construction of lncRNA-mRNA functional network. Step 3:
506 Score the lncRNA modulators by combining gene's score and network.

507 **Fig.2** lncRNA-mRNA functional network and properties. (A) A global network
508 consisting of 19883 coding genes and 12150 lncRNAs. (B) The degree of lncRNA in
509 the network. (C) The betweenness of lncRNA in the network.

510 **Fig.3** Assessment of predictive power for our method. (A) The performance of our
511 method for cancer genes. (B) The performance of four networks for lncRNA modulators.
512 (C) The performance of three methods for lncRNA modulators. (D) The performance
513 of our method with part of edges deleted.

514 **Fig.4** Functional Analysis of high-risk cancer lncRNA modulators. (A) The heatmap of
515 80 lncRNA scores for 9 cancers. (B) Functional annotation for common character
516 lncRNA modulators. (C) Functional annotation for cancer-specific lncRNA modulators.
517 (D) Venn diagram showing the number of intersections between neighbor miRNAs and
518 genes and cancer miRNAs and genes. (E) Schematic representation of lncRNA-linked
519 cancer miRNAs and genes in the network.

520 **Fig.5** Comparison of biological characteristics of different levels of lncRNA
521 modulators, including(A) number of exons, (B)number of miRNA binding sites, (C)
522 conservation score and (D)lncRNA expression.

523 **Fig.6** Analysis of top20 stomach cancer lncRNA modulators. (A) 4 lncRNAs were
524 significantly associated with survival. (B) Functional annotation for top 20 lncRNA
525 modulators. (C) The number of intersections between neighbor miRNAs and genes and
526 cancers miRNA and genes for 4 lncRNAs. (D) The subnet of the cancer genes linked to
527 these lncRNAs in the network. (E) The subnet of the cancer miRNAs linked to these
528 lncRNAs in the network.

529

Figures

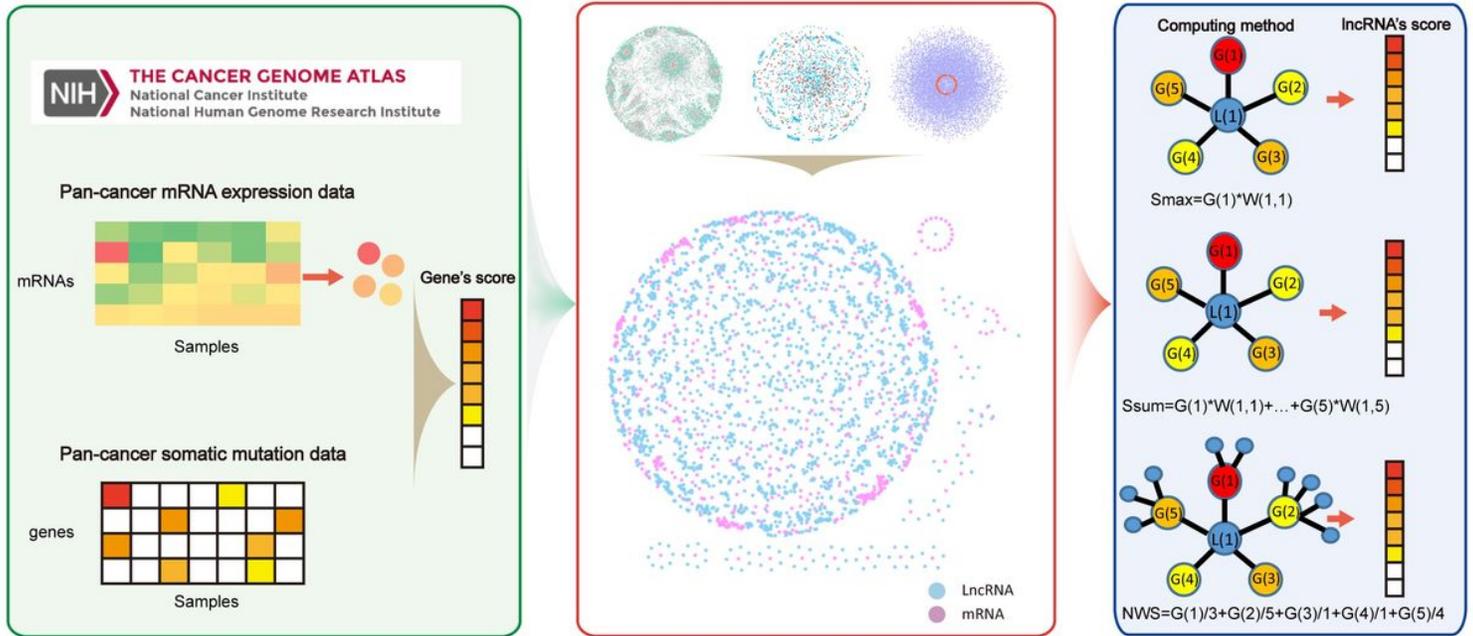


Figure 1

The workflow for prioritizing cancer lncRNA modulators. Step 1: Score the mutated gene. Step 2: Construction of lncRNA-mRNA functional network. Step 3: Score the lncRNA modulators by combining gene's score and network.

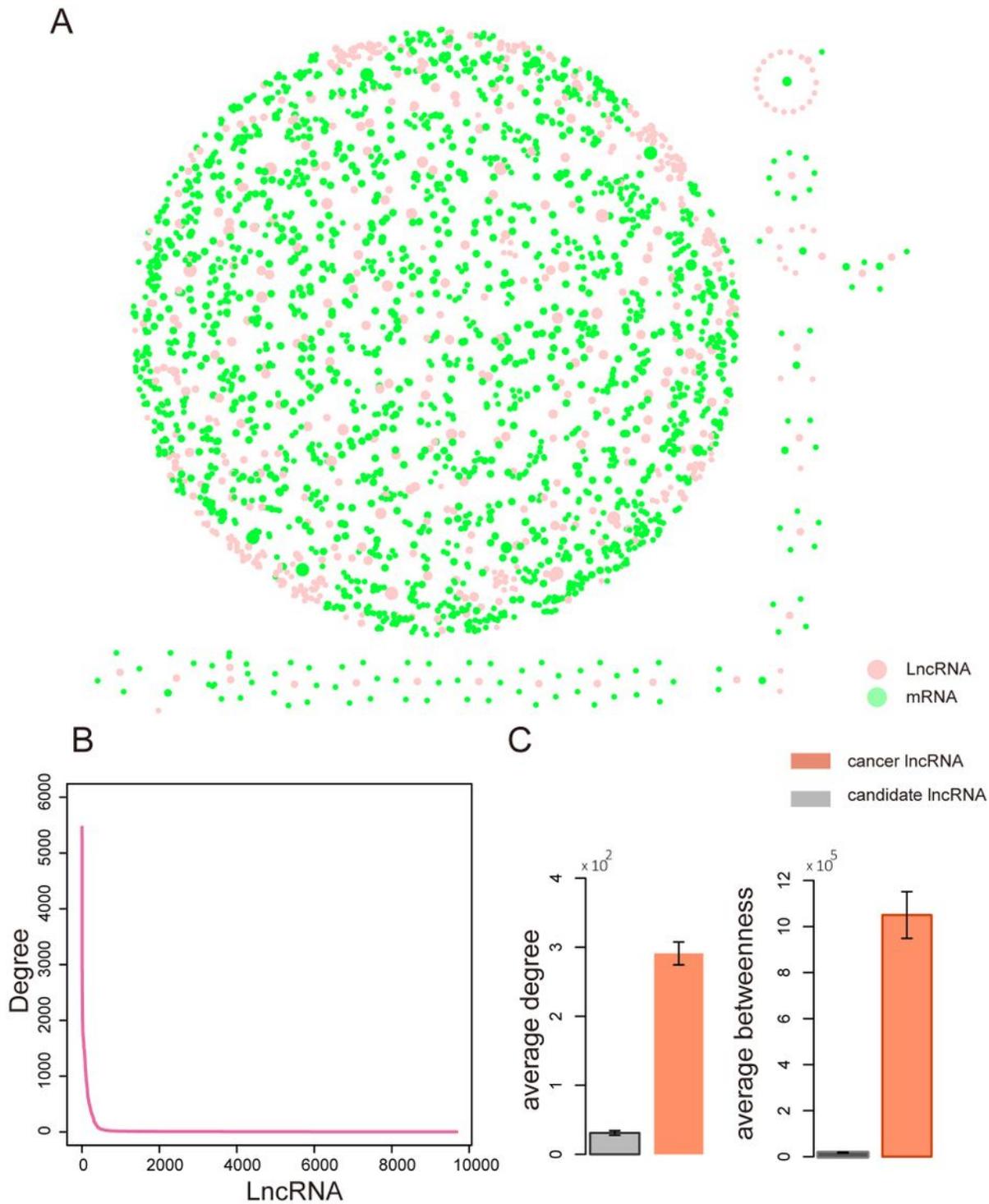


Figure 2

lncRNA-mRNA functional network and properties. (A) A global network consisting of 19883 coding genes and 12150 lncRNAs. (B) The degree of lncRNA in the network. (C) The betweenness of lncRNA in the network.

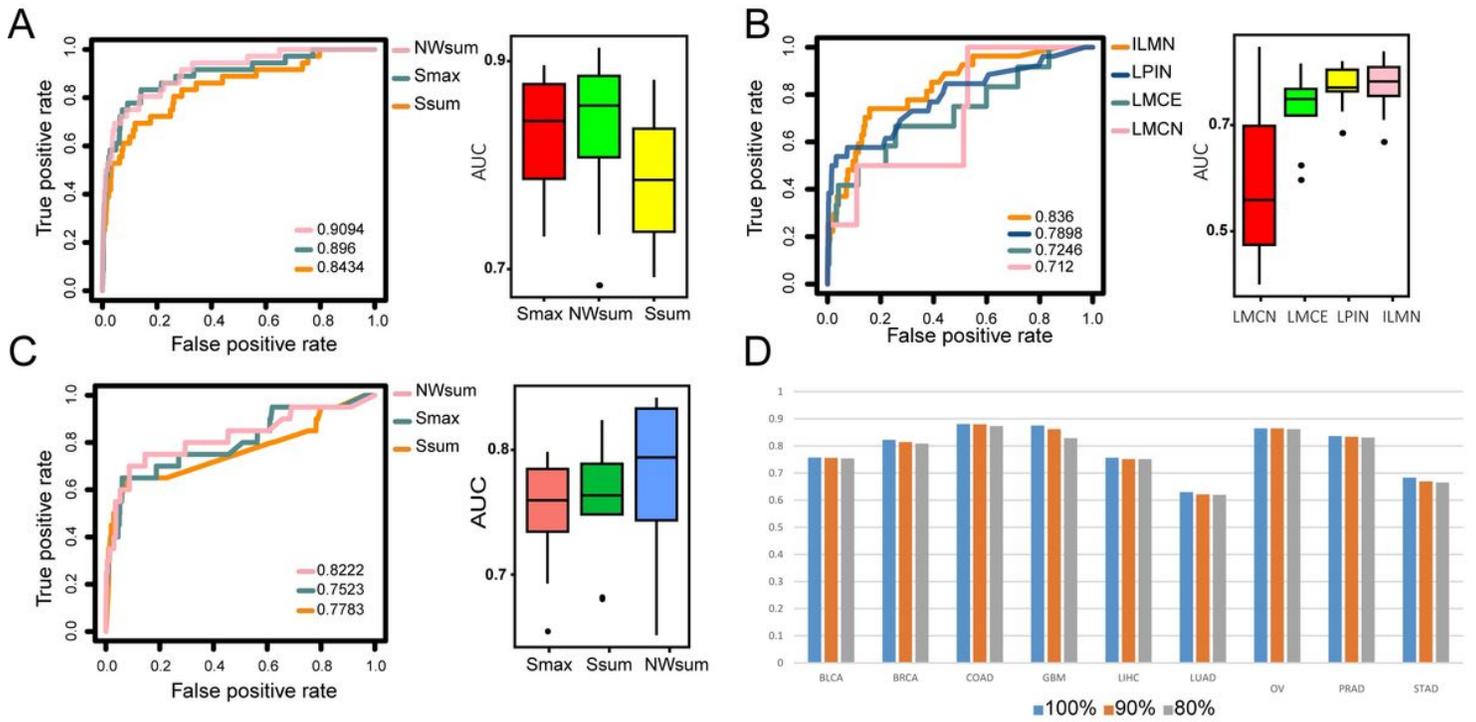


Figure 3

Assessment of predictive power for our method. (A) The performance of our method for cancer genes. (B) The performance of four networks for lncRNA modulators. (C) The performance of three methods for lncRNA modulators. (D) The performance of our method with part of edges deleted.

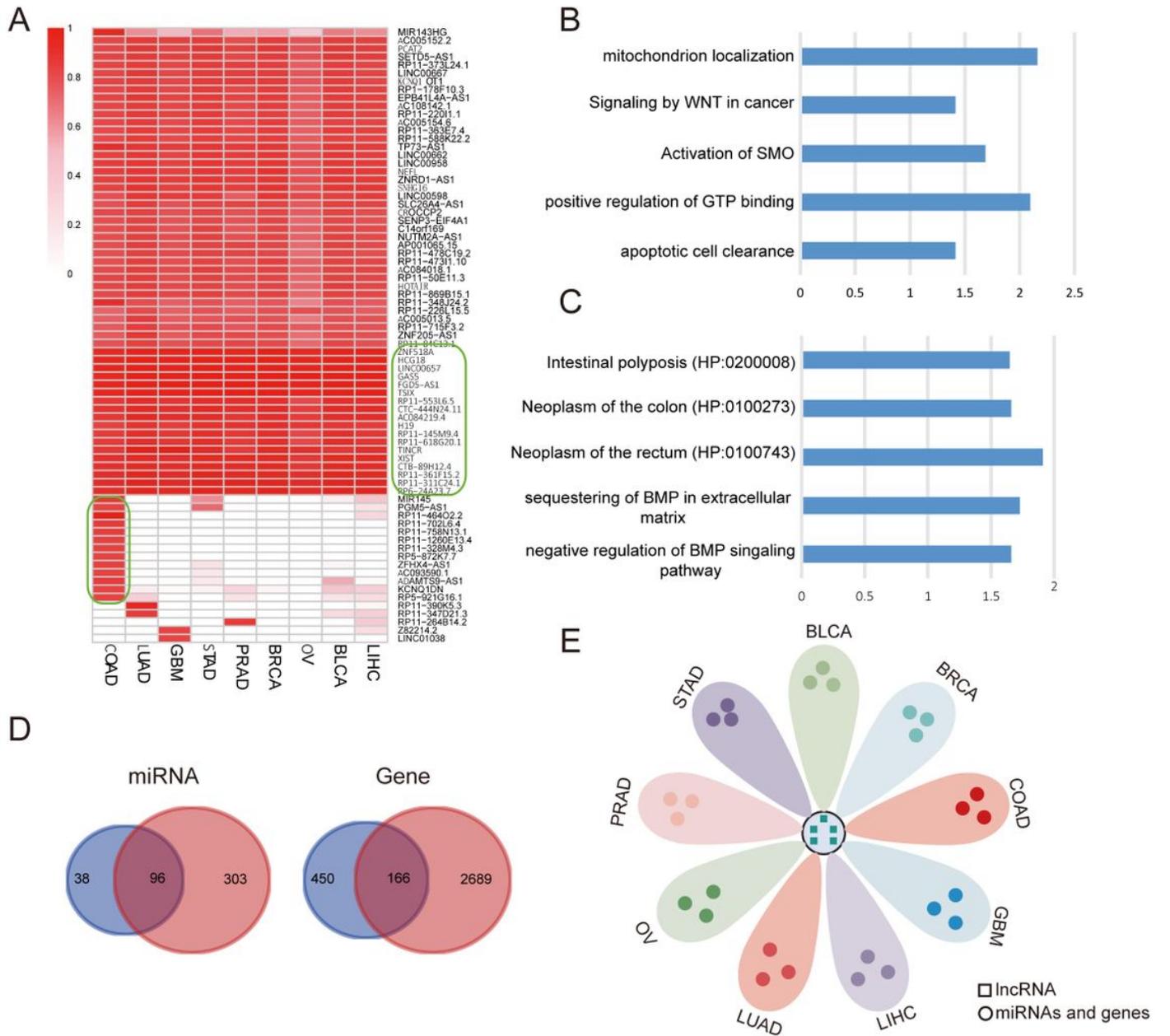


Figure 4

Functional Analysis of high-risk cancer lncRNA modulators. (A) The heatmap of 80 lncRNA scores for 9 cancers. (B) Functional annotation for common character lncRNA modulators. (C) Functional annotation for cancer-specific lncRNA modulators. (D) Venn diagram showing the number of intersections between neighbor miRNAs and genes and cancer miRNAs and genes. (E) Schematic representation of lncRNA-linked cancer miRNAs and genes in the network.

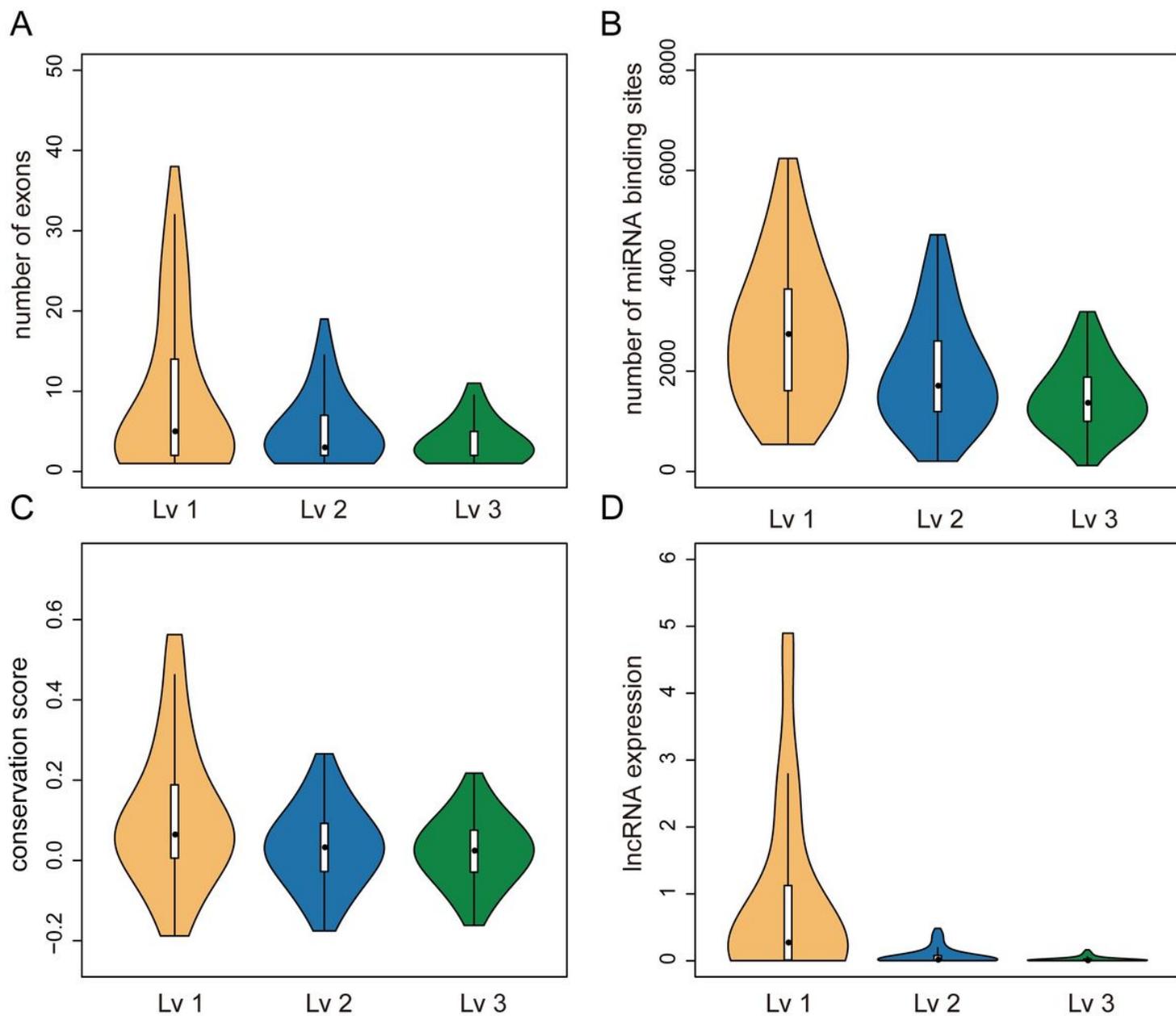


Figure 5

Comparison of biological characteristics of different levels of lncRNA modulators, including (A) number of exons, (B) number of miRNA binding sites, (C) conservation score and (D) lncRNA expression.

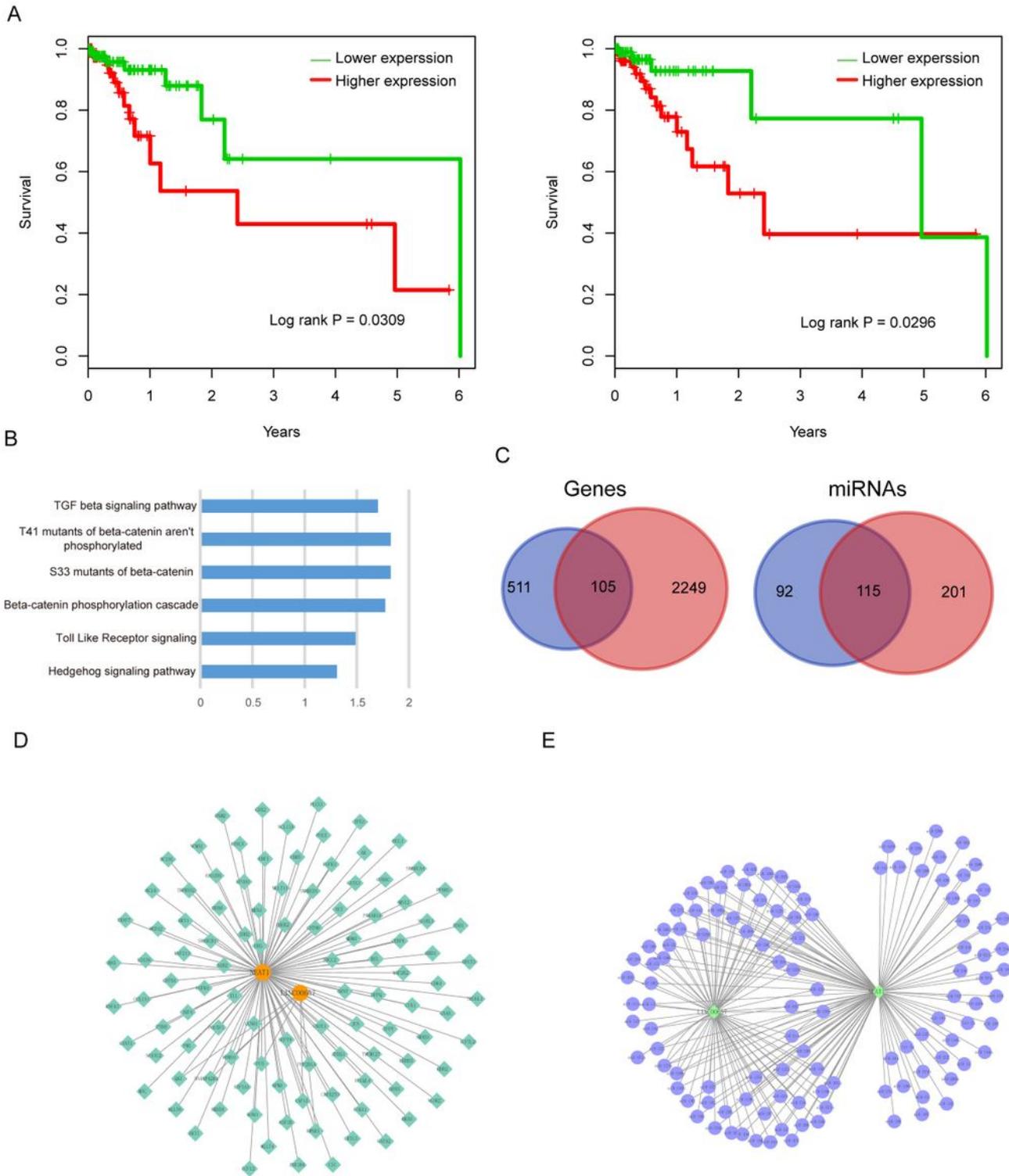


Figure 6

Analysis of top20 stomach cancer lncRNA modulators. (A) 4 lncRNAs were significantly associated with survival. (B) Functional annotation for top 20 lncRNA modulators. (C) The number of intersections between neighbor miRNAs and genes and cancers miRNA and genes for 4 lncRNAs. (D) The subnet of the cancer genes linked to these lncRNAs in the network. (E) The subnet of the cancer miRNAs linked to these lncRNAs in the network.

Supplementary Files

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

- [SupplementaryTablesS1.xlsx](#)
- [SupplementaryTablesS2.xlsx](#)