# 1 LMSM: a modular approach for identifying lncRNA related

# 2 miRNA sponge modules in breast cancer

3 Junpeng Zhang[1,2,*,†], Taosheng Xu[3,†], Lin Liu[4], Wu Zhang[5], Chunwen Zhao[2], Sijing Li[2], Jiuyong Li[4], Nini Rao[1,*]

4 and Thuc Duy Le[4,*]

5 [1]Center for Informational Biology, School of Life Science and Technology, University of Electronic Science

6 and Technology of China, Chengdu, Sichuan 610054, China

7 [2]School of Engineering, Dali University, Dali, Yunnan 671003, China

8 [3]Institute of Intelligent Machines, Hefei Institutes of Physical Science, Chinese Academy of Sciences, Hefei,

9 Anhui 230031, China

10 [4]School of Information Technology and Mathematical Sciences, University of South Australia, Mawson Lakes,

11 SA 5095, Australia

12 [5]School of Agriculture and Biological Sciences, Dali University, Dali, Yunnan 671003, China

13 *To whom correspondence should be addressed. Email: zhangjunpeng_411@yahoo.com

14 Correspondence may also be addressed to Nini Rao. Email: raonn@uestc.edu.cn

15 Correspondence may also be addressed to Thuc Duy Le. Email: thuc.le@unisa.edu.au

16 [†]The authors wish it to be known that, in their opinion, the first two authors should be regarded as joint First

17 Authors.

## 18 Abstract

19 Until now, existing methods for identifying lncRNA related miRNA sponge modules mainly

20 rely on lncRNA related miRNA sponge interaction networks, which may not provide a full

21 picture of miRNA sponging activities in biological conditions. Hence there is a strong need

22    of new computational methods to identify lncRNA related miRNA sponge modules. In this

23    work, we propose a framework, LMSM, to identify LncRNA related MiRNA Sponge

24    Modules from heterogeneous data. To understand the miRNA sponging activities in

25    biological conditions, LMSM uses gene expression data to evaluate the influence of the

26    shared miRNAs on the clustered sponge lncRNAs and mRNAs. We have applied LMSM to

27    the human breast cancer (BRCA) dataset from The Cancer Genome Atlas (TCGA). As a

28    result, we have found that the majority of LMSM modules are significantly implicated in

29    BRCA and most of them are BRCA subtype-specific. Most of the mediating miRNAs act as

30    crosslinks across different LMSM modules, and all of LMSM modules are statistically

31    significant. Multi-label classification analysis shows that the performance of LMSM modules

32    is significantly higher than baseline's performance, indicating the biological meanings of

33    LMSM modules in classifying BRCA subtypes. The consistent results suggest that LMSM is

34    robust in identifying lncRNA related miRNA sponge modules. Moreover, LMSM can be

35    used to predict miRNA targets. Finally, LMSM outperforms a graph clustering-based strategy

36    in identifying BRCA-related modules. Altogether, our study shows that LMSM is a

37    promising method to investigate modular regulatory mechanism of sponge lncRNAs from

38    heterogeneous data.

39    **Author summary**

40    Previous studies have revealed that long non-coding RNAs (lncRNAs), as microRNA

41    (miRNA) sponges or competing endogenous RNAs (ceRNAs), can regulate the expression

42    levels of messenger RNAs (mRNAs) by decreasing the amount of miRNAs interacting with

43    mRNAs. In this work, we hypothesize that the "tug-of-war" between RNA transcripts for

44    attracting miRNAs is across groups or modules. Based on the hypothesis, we propose a

45    framework called LMSM, to identify LncRNA related MiRNA Sponge Modules. Based on

46    the two miRNA sponge modular competition principles, significant sharing of miRNAs and

47    high canonical correlation between the sponge lncRNAs and mRNAs, LMSM is also capable

48    of predicting miRNA targets. LMSM not only extends the ceRNA hypothesis, but also

49    provides a novel way to investigate the biological functions and modular mechanism of

50    lncRNAs in breast cancer.

## Introduction

51

52    Long non-coding RNAs (lncRNAs) are RNA transcripts with more than 200 nucleotides (nts)

53    in length [1]. More and more evidence has shown that lncRNAs play important functional

54    roles in many biological processes, including human cancers [2-4]. As a major class of non-

55    coding RNAs (ncRNAs), lncRNAs have attracted increasing interest from researchers in their

56    exploration of non-coding knowledge from the 'junk'.

57    Among the wide range of biological functions of lncRNAs, their role as competing

58    endogenous RNAs (ceRNAs) or miRNA sponges is in the limelight. As a family of small

59    ncRNAs (~18nts in length), miRNAs are important post-transcriptional regulators of gene

60    expression [5,6]. According to the ceRNA hypothesis [7], lncRNAs contain abundant miRNA

61    response elements (MREs) for competitively sequestering target mRNAs from miRNAs'

62    control. This regulation mechanism of lncRNAs when acting as miRNA sponges is highly

63    implicated in various human diseases [8], including breast cancer [9]. For example, lncRNA

64    *H19*, an imprinted gene is associated with breast cancer cell clonogenicity, migration and

65    mammosphere-forming ability. By sponging miRNA *let-7*, *H19* forms a *H19/let-7/LIN28*

66    reciprocal negative regulatory circuit to play a critical role in the breast cancer stem cell

67    maintenance [10].

68    To systematically investigate the functions of lncRNAs as miRNA sponges in human

69    cancer, a series of computational methods have been developed to infer lncRNA related

70    miRNA sponge interaction networks. The methods can be divided into three categories

71    according to the statistical or computational techniques employed: pair-wise correlation based

72    approach, partial association based approach, and mathematical modelling approach [11].

73    It is commonly known that to implement a specific biological function, genes tend to

74    cluster or connect in the form of modules or communities. Consequently, based on the

75    identified lncRNA related miRNA sponge interaction networks, several methods [12-17]

76    using graph clustering algorithms were developed to identify lncRNA related miRNA sponge

77    modules. For the identification of sponge lncRNA-mRNA pairs, most of existing methods

78    only consider pair-wise correlation of them. Since the lncRNA related miRNA sponge

79    interaction networks are created by simply putting together sponge lncRNA-mRNA pairs,

80    when the expression levels of each sponge lncRNA-mRNA pair are highly correlated, the

81    collective correlation between the set of sponge lncRNAs and the set of mRNAs in the same

82    identified module is not necessarily high. As we know, the pair-wise positive correlation

83    between the expression levels of a lncRNA and a mRNA pair is commonly used to identify

84    the sponge interactions between them. For the identification of lncRNA related miRNA

85    sponge modules, it is also necessary to investigate whether the clustered sponge lncRNAs and

86    mRNAs in a module have high collective positive correlation or not. Moreover, these

87    methods do not consider the influence of the shared miRNAs on the expression of the

88    clustered sponge lncRNAs and mRNAs. It is known that the "tug-of-war" between sponge

89    lncRNAs and mRNAs is mediated by miRNAs. Therefore, it is extremely important to

90    consider the influence of the shared miRNAs in identifying lncRNA related miRNA sponge

91    modules.

92      Recently, to study lncRNA, miRNA and mRNA-associated regulatory modules, Deng *et*

93      *al*. [18] and Xiao *et al*. [19] have proposed two types of joint matrix factorization methods to

94      identify mRNA-miRNA-lncRNA co-modules by integrating gene expression data and

95      putative miRNA-target interactions. However, it is still not clear how the shared miRNAs

96      influence the expression level of the sponge lncRNAs and mRNAs in a module.

97      To address the above issues, we firstly hypothesize that sponge lncRNAs form a group to

98      competitively release a group of target mRNAs from the control of the miRNAs shared by

99      the lncRNAs and mRNAs (details see Section Methods). We name this hypothesis the *miRNA*

100     *sponge modular competition hypothesis* in this paper. Then based on the hypothesis, we

101     propose a novel framework to identify LncRNA related MiRNA Sponge Modules (LMSM).

102     The framework firstly uses the WeiGhted Correlation Network Analysis (WGCNA) [20]

103     method to generate lncRNA-mRNA co-expression modules. Next, by incorporating matched

104     miRNA expression and putative miRNA-target interactions, LMSM applies three constraints

105     (see Section Methods) to obtain lncRNA related miRNA sponge modules (also called LMSM

106     modules in this paper). One of the constraints, high canonical correlation, is used to assess

107     whether the group of sponge lncRNAs and the group of mRNAs in the same module have a

108     high collective positive correlation or not. The other constraint, adequate sensitivity canonical

109     correlation conditioning on a group of miRNAs, is used to evaluate the influence of the

110     shared miRNAs on the clustered sponge lncRNAs and mRNAs.

111     To evaluate the LMSM approach, we apply it to matched miRNA, lncRNA and mRNA

112     expression data, and clinical information of breast cancer (BRCA) dataset from The Cancer

113     Genome Atlas (TCGA). The modular analysis results demonstrate that LMSM can help to

114     uncover modular regulatory mechanism of sponge lncRNAs in BRCA. LMSM is released

115    under the GPL-3.0 License, and is freely available through GitHub repository

116    (https://github.com/zhangjunpeng411/LMSM).


## Materials and methods

### A hypothesis on miRNA sponge modular competition

119    The ceRNA hypothesis [7] indicates that a pool of RNA transcripts (known as ceRNAs)

120    regulate each other's transcripts by competing for the shared miRNAs through MREs. Based

121    on this unifying hypothesis, a large-scale gene regulatory network including coding and non-

122    coding RNAs across the transcriptome can be formed, and it plays critical roles in human

123    physiological and pathological processes. However, by using MREs as letters of language,

124    the hypothesis only depicts the crosstalk between individual RNA transcript (e.g. coding

125    RNAs, lncRNAs, circRNAs or pseudogenes) and mRNA at the pair-wise interaction level

126    and the crosstalk between RNA transcripts and mRNAs at the module level is still an open

127    question.

128        There has been evidence showing that for the same transcriptional regulatory program,

129    biological process or signaling pathway, genes tend to form modules or communities to

130    coordinate biological functions [21]. These modules correspond to functional units in

131    complex biological systems, and they play important role in gene regulation. Based on these

132    findings, in this paper, we hypothesize that regarding miRNA sponging, the crosstalk

133    between different RNA transcripts is in the form of modular competition. We call the

134    hypothesis the *miRNA sponge modular competition hypothesis*.

135        As shown in Fig 1, based on our hypothesis, instead of having pair-wise competitions,

136    miRNA sponges form groups to compete at module level for common miRNAs. Here, a

6

137    miRNA sponge module consists of a competing group (other coding RNA group, pseudogene

138    group, circRNA group or lncRNA group) and a mRNA group. Here, other coding RNAs also

139    include mRNAs. From the perspective of modularity, the hypothesis at module level extends

140    the ceRNA hypothesis and provides a new channel to look into the functions and regulatory

141    mechanism of miRNA sponges or ceRNAs. Since the available resources of lncRNAs are

142    more abundant than those of other non-coding RNAs (e.g. circRNAs and pseudogenes), in

143    this paper, we focus on the competition between lncRNAs and mRNAs to validate and

144    demonstrate the proposed *miRNA sponge modular competition hypothesis*. Our goal is to

145    discover lncRNA related sponge modules, or LMSM modules. Here each LMSM module

146    contains a group of lncRNAs which compete collectively with a group of mRNAs for

147    sponging the same set of miRNAs.
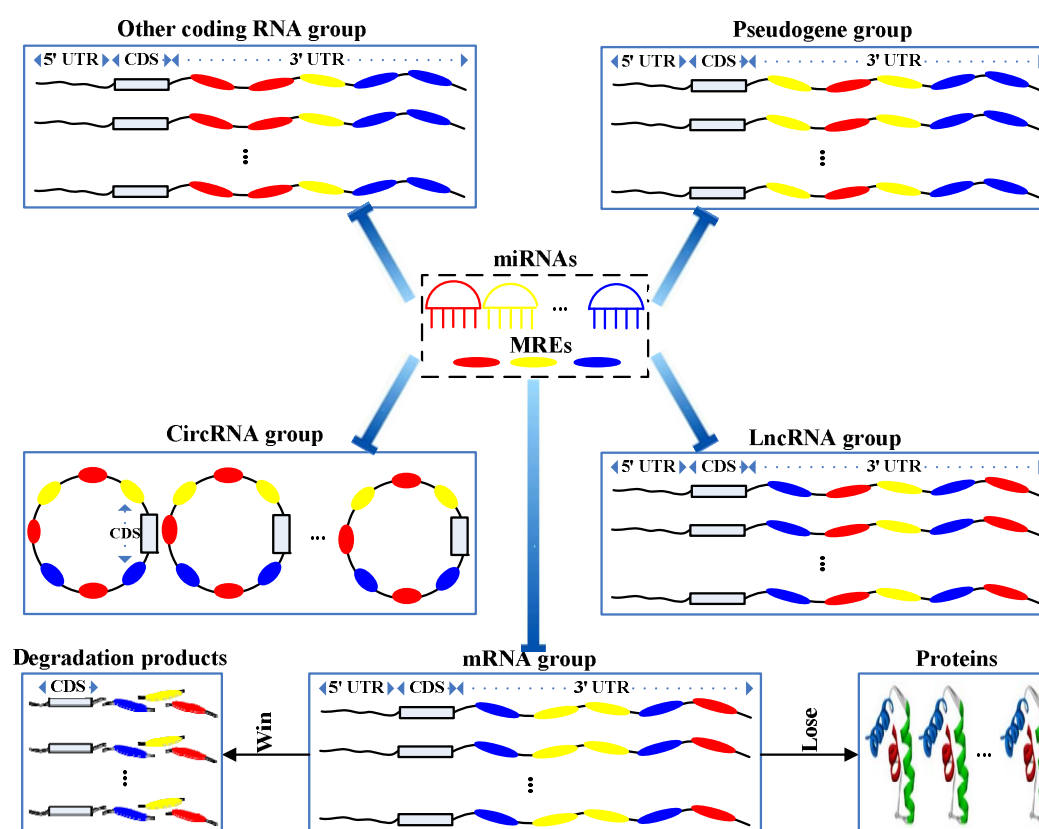


148
149    **Fig 1. An illustration of the miRNA sponge modular competition hypothesis.** The four
150    types of miRNA sponges (other coding RNAs, lncRNAs, circRNAs or pseudogenes),
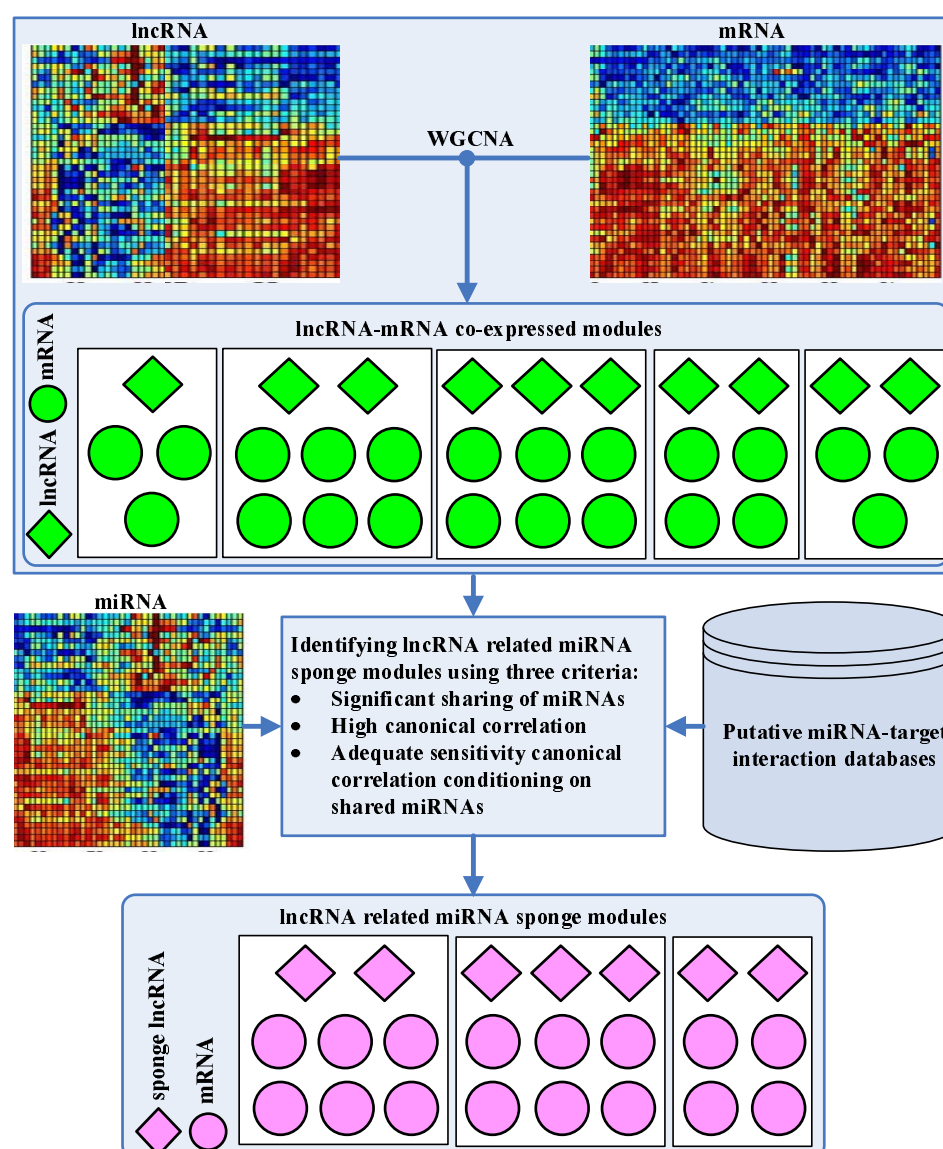151    miRNAs and their target mRNAs are shown. Each miRNA sponge module consists of a

7

152    group of the same type of miRNA sponges, e.g. a group of lncRNAs and a group of target
153    mRNAs. In the same module, the group of miRNA sponges competes with the group of
154    target mRNAs for binding with a set of miRNAs. If the miRNA sponges win the competition,
155    the group of target mRNAs will be released from repression and they will be translated into
156    proteins. If the miRNA sponges lose the competition, the group of target mRNAs will be
157    post-transcriptionally repressed and degraded.

158    **The LMSM framework**

159    *Overview of LMSM*. As shown in Fig 2, the proposed LMSM framework comprises two

160    stages. In stage 1, the WGCNA method [20] is used for finding lncRNA-mRNA co-

161    expression modules from matched lncRNA and mRNA expression data. Then in stage 2,

162    LMSM identifies LMSM modules from the lncRNA-mRNA co-expression modules using

163    three criteria. That is, a co-expression module is considered as a LMSM module if the group

164    of lncRNAs and the group of mRNAs in the co-expression module: (1) have significant

165    sharing of miRNAs, (2) have high canonical correlation between their expression levels, and

166    (3) have adequate sensitivity canonical correlation conditioning on their shared miRNAs.

167    LMSM checks the criteria one by one, and once a co-expression module does not meet a

168    criterion, it is discarded and will not be checked for the next criterion. In the following, we

169    will describe the two stages in detail.

170    *Identifying lncRNA-mRNA co-expression modules*. For identifying lncRNA-mRNA co-

171    expression modules, we use the WGCNA method. WGCNA is a popular method for

172    identifying co-expressed genes across samples and it can be used to identify clusters of highly

173    co-expressed lncRNAs and mRNAs. In our task, we use the matched lncRNA and mRNA

174    expression data as input to the *WGCNA* R package [20] to identify lncRNA-mRNA co-

175    expression modules. We use the scale-free topology criterion for soft thresholding. The

176    coefficient of determination $R^2$ (the range is from 0 to 1) is used to quantify the goodness of

177    scale-free topology, and larger $R^2$ values mean better scale-free topology. Normally, the $R^2$

8

178     value larger than 0.8 in power law curve fit is ranked as good-level in the WGCNA method.

179     Therefore, the desired minimum scale free topology fitting index $R^2$ is set as 0.8 in this work.



180
181     **Fig 2. Workflow of LMSM**. Firstly, we use the WGCNA method to infer lncRNA-mRNA
182     co-expression modules from the matched lncRNA and mRNA expression. Then by using
183     miRNA expression data and putative miRNA-target interactions, we infer lncRNA related
184     miRNA sponge modules (LMSM) by applying three criteria: significant sharing of miRNAs
185     by the group of lncRNAs and the group of target mRNAs in the same co-expression module,
186     high canonical correlation between the lncRNA group and the target mRNA group, and
187     adequate sensitivity canonical correlation between the lncRNA group and the target mRNA
188     group conditioning on shared miRNAs. Each LMSM module must contain at least two
189     sponge lncRNAs and two target mRNAs.

190 ***Inferring lncRNA related miRNA sponge modules***. To identify lncRNA related miRNA

191 sponge modules from the co-expression modules obtained in stage 1, we propose three

192 criteria (detailed below) by following the key tenet of our *miRNA sponge modular*

193 *competition hypothesis*. That is, a group of lncRNAs (acting as miRNA sponges) competes

194 with a group of mRNAs with respect to a set of miRNAs shared by the two groups.

195 The first criterion requires that the group of lncRNAs and the group of mRNAs in a

196 miRNA sponge module have a significant sharing of a set of miRNAs. LMSM uses a

197 hypergeometric test to assess the significance of the sharing of miRNAs between the group of

198 lncRNAs and the group of mRNAs in a co-expression module, based on putative miRNA-

199 target interactions. The *p*-value for the test is computed as:

200
$$p-value = 1 - \sum_{i_1=0}^{L_1-1} \frac{\binom{M_1}{i_1}\binom{N_1-M_1}{K_1-i_1}}{\binom{N_1}{K_1}} \qquad (1)$$

201 In the equation, $N_1$ is the number of all miRNAs in the dataset, $M_1$ and $K_1$ denote the total

202 numbers of miRNAs interacting with the group of lncRNAs and the group of mRNAs in the

203 co-expression module respectively, and $L_1$ (e.g. 3) is the number of miRNAs shared by the

204 group of lncRNAs and the group of mRNAs in the co-expression module.

205 The second criterion is to assure that the sponge modular competition between the group

206 of lncRNAs and the group of mRNAs in a miRNA sponge module is strong enough. In

207 existing work, to identify lncRNA related mRNA sponge interactions, a principle followed is

208 that the expression level of a lncRNA and the expression level of a mRNA need to be

209 strongly and positively correlated. Following the same principle on strong positive correlation

210 in expression levels while considering our modular competition hypothesis, LMSM requires

211    the *collective* correlation between the expression levels of the group of lncRNAs and the

212    group of target mRNAs in the same module to be strong and positive. To assess the *collective*

213    correlation, we perform canonical correlation analysis [22] to obtain the canonical correlation

214    between the group of lncRNAs and the group of mRNAs in a co-expression module. Let the

215    two column vectors $X = (x_1, x_2, ..., x_m)^T$ and $Y = (y_1, y_2, ..., y_n)^T$ represent the group of

216    lncRNAs and the group of mRNAs in a co-expression module respectively. $\Sigma_{XX}$, $\Sigma_{YY}$ and

217    $\Sigma_{XY}$ are the variance or cross-covariance matrices calculated from the expression data of $X$

218    and $Y$. The canonical correlation analysis seeks the canonical vectors $a$ ( $a \in \mathbb{R}^m$ ) and $b$

219    ( $b \in \mathbb{R}^n$ ) which maximize the correlation of $corr(a^T X, b^T Y)$. The canonical correlation

220    between the group of lncRNAs and the group of mRNAs, denoted as $CC_{lncR\text{-}mR}$, is then

221    calculated as follows with the found canonical vectors:

$$CC_{lncR-mR} = corr(a^T X, b^T Y) = \frac{a^T \sum_{XY} b}{\sqrt{a^T \sum_{XX} a}\sqrt{b^T \sum_{YY} b}} \tag{2}$$

223    In this work, we use the *PMA* R package [23] to compute canonical correlation.

224    Finally, the third criterion adapted from the sensitivity correlation [24] is employed to

225    assess if the miRNAs shared by the group lncRNAs and the group of mRNAs in a module

226    have large enough influence on the modular competition between the two groups of RNAs.

227    To check according to this criterion, we incorporate miRNA expression data, and compute

228    $SCC_{lncR\text{-}mR}$ the sensitivity canonical correlation between the group of lncRNAs and the group

229    of mRNAs in a co-expression module as follows:

$$SCC_{lncR-mR} = CC_{lncR-mR} - PCC_{lncR-mR} \tag{3}$$

11

231   where $PCC_{lncR\text{-}mR}$ is the partial canonical correlation between the group of lncRNAs and the

232   group of mRNAs, i.e. the canonical correlation conditioning on the expression of their shared

233   miRNAs in the co-expression module, or the canonical correlation between the two groups of

234   RNAs when the influence of the shared miRNAs is eliminated. Therefore, from Eq. (3), we

235   see that $SCC_{lncR\text{-}mR}$ implies the correlation between the two groups of RNAs under the

236   influence of their shared miRNAs.

237      $PCC_{lncR\text{-}mR}$ in Eq. (4) can be calculated as:

238
$$PCC_{lncR-mR} = \frac{CC_{lncR-mR} - CC_{miR-mR}CC_{miR-lncR}}{\sqrt{1 - CC_{miR-mR}^2}\sqrt{1 - CC_{miR-lncR}^2}} \tag{4}$$

239   where $CC_{miR\text{-}mR}$ ($CC_{miR\text{-}lncR}$) is the canonical correlation between the set of miRNAs in the co-

240   expression module and the group of mRNAs (lncRNAs) in the co-expression module.

241      In this study, empirically, a lncRNA-mRNA co-expressed module with $p$-value < 0.05 for

242   the hypergeometric test of miRNA sharing (criterion 1), $CC_{lncR\text{-}mR}$ > 0.8 for modular

243   competition strength assessment (criterion 2) and $SCC_{lncR\text{-}mR}$ > 0.1 for miRNA influence

244   (criterion 3) is regarded as a lncRNA related miRNA sponge module (a LMSM module).

245   **Evaluating statistical significance of LMSM modules**

246   To evaluate the statistical significance of LMSM modules, we adapt the null model method

247   proposed in [25]. The null model method hypothesizes that the shared miRNAs do not affect

248   the correlation between two genes, i.e. the sensitivity correlation (the difference between

249   correlation and partial correlation) between two genes is 0, and has been successfully applied

250   to evaluate statistical significance of ceRNA interactions. Similar to [25], LMSM is also

251   adapted from the Sensitivity Correlation (SC) method [24]. Therefore, the null model method

252   can be applied to evaluate the statistical significance of LMSM modules. In our null model,

12

253   the null hypothesis is that the group of the shared miRNAs does not influence the canonical

254   correlation between the group of lncRNAs and the group of mRNAs, i.e. $SCC_{lncR-mR} = 0$. For

255   each LMSM module, a group of lncRNAs or a group of mRNAs corresponds to a gene, and a

256   group of the shared miRNAs corresponds to a miRNA in the null model. For obtaining more

257   precise $p$-values, the number of datasets sampled is set to 1E+06 for the null model. Since the

258   sampling procedure is computationally intensive, we use the pre-computed sets of covariance

259   matrices in *SPONGE* R package [25] to build our null model. Based on the constructed null

260   model, we can infer adjusted $p$-values (adjusted by Benjamini and Hochberg method [26]) for

261   each LMSM module. A LMSM module with adjusted $p$-value less than 0.05 is regarded as a

262   statistically significant module.

263   **Application of LMSM in BRCA**

264   *BRCA enrichment analysis*. Instead of performing Gene Ontology (GO) and Kyoto

265   Encyclopedia of Genes and Genomes Pathway (KEGG) enrichment analysis, to investigate

266   whether an identified LMSM module is functionally associated with BRCA, we focus on

267   conducting BRCA enrichment analysis by using a hypergeometric test. For a LMSM module,

268   the $p$-value for the test is calculated as:

269
$$p-value = 1 - \sum_{i_2=0}^{L_2-1} \frac{\binom{M_2}{i_2}\binom{N_2-M_2}{K_2-i_2}}{\binom{N_2}{K_2}} \tag{5}$$

270   where $N_2$ is the number of genes (lncRNAs and mRNAs) in the dataset, $M_2$ denotes the

271   number of BRCA genes in the dataset, $K_2$ represents the number of genes in the LMSM

272   module, and $L_2$ is the number of BRCA genes in the LMSM module. A LMSM module with

273   $p$-value < 0.05 is regarded as a BRCA-related module.

274      *Module biomarker identification in BRCA*. The module survival analysis can imply whether

275      the identified LMSM modules are good biomarkers of the metastasis risks of cancer patients

276      or not, and it can give us the hint whether the LMSM modules may be related to and

277      potentially affect the metastasis or survival of cancer patients. For each BRCA sample, we fit

278      the multivariate Cox model (proportional hazards regression model) [27] using the genes

279      (lncRNAs and mRNAs) in LMSM modules to compute its risk score. All the BRCA samples

280      are equally divided into the high risk and the low risk groups according to their risk scores.

281      The Log-rank test is used to evaluate the difference of each LMSM module between the high

282      and the low risk BRCA groups. Moreover, we also calculate the proportional hazard ratio

283      (HR) between the high and the low risk BRCA groups. In this work, the *survival* R package

284      [28] is utilized, and a LMSM module with Log-rank $p$-value $< 0.05$ and HR $> 2$ is regarded

285      as a module biomarker in BRCA.

286      *Identification of BRCA subtype-specific modules*. As is known, BRCA is a heterogeneous

287      disease with several molecular subtypes, and the choice of chemotherapy for each BRCA

288      subtype is different. This diversity indicates that the genetic regulation of each BRCA

289      subtype is specific. To identify BRCA subtype-specific modules, we firstly identify BRCA

290      molecular subtypes using the PAM50 classifier [29]. By using a 50-gene subtype predictor,

291      the PAM50 classifier classifies a BRCA sample into one of the five "intrinsic" subtypes:

292      Luminal A (LumA), Luminal B (LumB), HER2-enriched (Her2), Basal-like (Basal) or

293      Normal-like (Normal). In this work, we use the *genefu* R package [30] to predict molecular

294      subtypes of each BRCA sample in the dataset used in our study.

295      To identify BRCA subtype-specific LMSM modules, we firstly need to estimate the

296      enrichment scores of LMSM modules in BRCA samples. To calculate the enrichment score

297      of each LMSM module in BRCA samples, the gene set variation analysis (GSVA) method

298    [31] is used. To calculate the enrichment score, the GSVA method uses the Kolmogorov-

299    Smirnov (KS) like random walk statistic as follows:

300
$$v_{jk}(1) = \frac{\sum_{i=1}^{1} |r_{ij}|^{\tau} I(g(i) \in \gamma_k)}{\sum_{i=1}^{p} |r_{ij}|^{\tau} I(g(i) \in \gamma_k)} - \frac{\sum_{i=1}^{1} I(g(i) \notin \gamma_k)}{p - |\gamma_k|}$$
(6)

301    where $\tau$ ($\tau$ =1 by default) is the weight of the tail in the random walk, $r_{ij}$ is the normalized

302    expression-level statistics of the $i$-th gene in the $j$-th sample as defined in [31], $\gamma_k$ is the $k$-th

303    LMSM module, $I(g(i) \in \gamma_k)$ is the indicator function on whether the $i$-th gene belongs to the

304    LMSM module $\gamma_k$, $|\gamma_k|$ is the number of genes in the $k$-th LMSM module, and $p$ is the

305    number of genes in the dataset.

306    To transform the KS like random walk statistic into an enrichment score (*ES*, also called

307    GSVA score), we calculate the maximum deviation from zero of the random walk of the $j$-th

308    sample with respect to the $k$-th LMSM module in the following:

309
$$ES_{jk}^{\max} = v_{jk}[\arg \max_{1=1,...,p} (abs(v_{jk}(1)))]$$
(7)

310    For each LMSM module $\gamma_k$, the formula generates a distribution of enrichment scores that is

311    bimodal (see the reference [31] for a more detailed description).

312    Based on the enrichment scores of LMSM modules in each BRCA sample, we further

313    identify two types of BRCA subtype-specific LMSM modules, up-regulated modules and

314    down-regulated modules. For one type of regulation pattern (up or down regulation), a

315    LMSM module is regarded to be specific to a BRCA subtype. For an up-regulated BRCA

316    subtype-specific LMSM module, the enrichment score of the LMSM module in the specific

317    BRCA subtype samples is significantly larger than the score in the other BRCA subtype

15

318    samples. For a down-regulated BRCA subtype-specific LMSM module, the enrichment score

319    of the LMSM module in the specific BRCA subtype samples is significantly smaller than the

320    score in the other BRCA subtype samples. For example, if a LMSM module $\gamma_k$ is up-

321    regulated Basal-like specific, the enrichment scores of the LMSM module in Basal-like

322    samples should be significantly larger than those in Luminal A, Luminal B, HER2-enriched

323    and Normal-like samples. In this work, for each LMSM module, we use Welch's $t$-test [32] to

324    calculate the significance $p$-value for the difference of the average enrichment scores between

325    any two BRCA subtype samples. Given a BRCA subtype, a LMSM module is considered as

326    an up-regulated (or down-regulated) module specific to this BRCA subtype if the module's

327    average enrichment score in samples of the given subtype is higher (or smaller) than the

328    average enrichment score in samples of any other subtype and the significance $p$-value of the

329    Welch's $t$-test between the samples of this subtype and any other subtype is less than 0.05.

330    ***Performance of LMSM modules in classifying BRCA subtypes***. In this section, to check the

331    biological relevance of the discovered LMSM modules, we conduct module classification of

332    BRCA subtypes. Here, classifying BRCA subtypes (LumA, LumB, Her2, Basal and Normal)

333    is a multi-class classification (also known as a special case of multi-label classification). To

334    understand the classification performance of the feature genes in each LMSM module, we

335    apply a state-of-the-art multi-label learning strategy called Binary Relevance (BR) [33]

336    implemented in the *utiml* R package [34] to conduct multi-label classification analysis. For

337    the BR strategy, we use the Support Vector Machine (SVM) classifier [35] with default

338    parameters implemented in *e1071* R package [36] as the base algorithm to build the multi-

339    label model. We select two commonly used multi-label classification measures: *Subset*

340    *accuracy* and *Hamming loss*, and conduct 10-fold cross-validation to evaluate the

341    performance of each LMSM module. In this work, *Subset accuracy* denotes the percentage of

342    correct predictions and *Hamming loss* is the fraction of wrong predictions to the total number

343    of predictions. Higher values of *Subset accuracy* and smaller values of *Hamming loss*

344    indicate better classification performance. In addition, for the evaluation, we use the baseline

345    method in [37], a commonly used multi-label classification method as the baseline for

346    comparison. The base algorithm of the baseline method is also the SVM classifier with

347    default parameters implemented in *e1071* R package [36].


348    **Results**


349    **Heterogeneous data sources**


350    We collect matched miRNA, lncRNA and mRNA expression data, and clinical data of BRCA

351    dataset from The Cancer Genome Atlas (TCGA, https://cancergenome.nih.gov/). A lncRNA

352    or mRNA without a corresponding gene symbol in the expression data of BRCA dataset is

353    removed. To obtain a unique expression value for replicates of miRNAs, lncRNAs or

354    mRNAs, we compute the average expression value of the replicates. As a result, we obtain

355    the matched expression data of 674 miRNAs, 12711 lncRNAs and 18344 mRNAs in 500

356    BRCA samples.


357       We retrieve putative miRNA-target interactions (including miRNA-lncRNA and miRNA-

358    mRNA interactions) from several high-confidence miRNA-target interaction databases and

359    use the combined database search results. Specifically, the putative miRNA-lncRNA

360    interactions are obtained from NPInter v3.0 [38] and the experimental module of DIANA-

361    LncBase v2.0 [39], and miRNA-mRNA interactions are from miRTarBase v8.0 [40],

362    TarBase v7.0 [41] and miRWalk v2.0 [42].

363   The BRCA related mRNAs are from DisGeNET v5.0 [43] and COSMIC v86 [44], and the

364   BRCA related lncRNAs are from LncRNADisease v2.0 [45], Lnc2Cancer v2.0 [46] and

365   MNDR v2.0 [47]. The ground truth of lncRNA related miRNA sponge interactions is

366   obtained by integrating the interactions from miRSponge [48], LncCeRBase [49] and

367   LncACTdb v2.0 [50].

368   **Most of the mediating miRNAs act as crosslinks across LMSM modules**

369   Following the steps shown in Fig 2, we have identified 17 LMSM modules (details can be

370   seen in S1 Data). The average size of the identified modules is 672.53 and the average

371   number of the shared miRNAs in a module is 232.82. In total, there are 549 unique miRNAs

372   mediating the 17 LMSM modules, and 90.16% (495 out of 549) miRNAs mediate at least

373   two LMSM modules (details can be seen in S2 Data). This result indicates that most of the

374   mediating miRNAs act as crosslinks across different LMSM modules.

375   **LMSM modules are all statistically significant**

376   In this section, by computing null-model-based $p$-values, we evaluate whether the identified

377   LMSM modules are statistically significant or not. As a result, the adjusted $p$-values for the

378   identified 17 LMSM modules (details can be seen in S3 Data) are all statistically significant

379   (adjusted $p$-value = 1.00E-06). This result demonstrates that LMSM modules are all

380   statistically significant.

381   **Most of LMSM modules are implicated in BRCA**

382   To investigate whether the identified LMSM modules are related to BRCA or not, we

383   conduct BRCA enrichment analysis and identify BRCA module biomarkers using the

384   methods described in Section Methods. For the BRCA enrichment analysis, we have

18

385     collected a list of 4819 BRCA genes (734 BRCA lncRNAs and 4085 BRCA mRNAs)

386     associated with the matched lncRNA and mRNA expression data (details in S4 Data). As

387     shown in Table 1, 10 out of 17 LMSM modules are functionally enriched in BRCA at a

388     significant level ($p$-value $< 0.05$). In Table 2, 15 out of 17 LMSM modules are regarded as

389     module biomarkers in BRCA at a significant level (Log-rank $p$-value $< 0.05$ and HR $> 2$).

390     Particularly, 90% (9 out of 10, excepting LMSM 14) of the BRCA-related LMSM modules

391     can act as module biomarker in BRCA. These results show that most of LMSM modules are

392     functionally implicated in BRCA.

393     **Table 1. BRCA-related LMSM modules**. $L_2$ is the number of BRCA genes in each LMSM
394     module, $K_2$ represents the number of genes in each LMSM module, the number of BRCA
395     genes in the dataset ($M_2$) is 4819, and the number of genes in the dataset ($N_2$) is 31055.

| Module ID | $L_2$ | $K_2$ | $p$-value |
|---|---|---|---|
| LMSM 2 | 327 | 1338 | 0 |
| LMSM 3 | 259 | 1340 | 7.34E-05 |
| LMSM 4 | 78 | 392 | 1.14E-02 |
| LMSM 5 | 89 | 449 | 8.07E-03 |
| LMSM 6 | 88 | 370 | 1.97E-05 |
| LMSM 8 | 275 | 880 | 0 |
| LMSM 12 | 24 | 110 | 4.95E-02 |
| LMSM 13 | 20 | 76 | 1.05E-02 |
| LMSM 14 | 252 | 1004 | 8.88E-16 |
| LMSM 16 | 48 | 182 | 1.11E-04 |

396     **Table 2. Survival analysis of LMSM modules in BRCA**. HRlow95 and HRup95 represent
397     the lower and upper of 95% confidence interval of HR, respectively.

| Module ID | Chi-square | $p$-value | HR | HRlow95 | HRup95 |
|---|---|---|---|---|---|
| LMSM 1 | 170.37 | 0 | 10.75 | 5.88 | 19.65 |
| LMSM 2 | 107.34 | 0 | 6.03 | 3.12 | 11.66 |
| LMSM 3 | 90.62 | 0 | 5.43 | 2.94 | 10.01 |
| LMSM 4 | 138.81 | 0 | 14.94 | 8.83 | 25.27 |
| LMSM 5 | 148.49 | 0 | 8.64 | 4.63 | 16.13 |
| LMSM 6 | 142.64 | 0 | 13.40 | 7.83 | 22.92 |
| LMSM 7 | 161.91 | 0 | 13.97 | 8.01 | 24.36 |
| LMSM 8 | 103.63 | 0 | 5.91 | 3.07 | 11.37 |
| LMSM 10 | 144.86 | 0 | 8.63 | 4.74 | 15.71 |
| LMSM 11 | 120.79 | 0 | 9.49 | 5.55 | 16.23 |
| LMSM 12 | 49.31 | 2.19E-12 | 5.46 | 3.38 | 8.80 |
| LMSM 13 | 60.08 | 9.10E-15 | 5.72 | 3.48 | 9.41 |

| | | | | | |
|---|---|---|---|---|---|
| LMSM 15 | 83.26 | 0 | 12.00 | 7.46 | 19.32 |
| LMSM 16 | 110.94 | 0 | 11.25 | 6.79 | 18.66 |
| LMSM 17 | 106.96 | 0 | 9.14 | 5.42 | 15.41 |

**LMSM modules are mostly BRCA subtype-specific**

In this section, we firstly divide the 500 BRCA samples into five "intrinsic" subtypes (Luminal A, Luminal B, HER2-enriched, Basal-like and Normal-like). The numbers of LumA, LumB, Her2, Basal and Normal samples are 190, 155, 52, 85 and 18, respectively. Then we calculate the enrichment scores of the identified 17 LMSM modules in the BRCA subtype samples respectively (details in S5 Data).

As illustrated in Fig 3, out of the 17 LMSM modules, 4 and 6 modules are identified as up-regulated and down-regulated BRCA subtype-specific LMSM modules, respectively. For the up-regulated BRCA subtype-specific LMSM modules, the numbers of Basal-specific, LumB-specific and Normal-specific modules are 1, 1 and 2, respectively. The numbers of Basal-specific, LumB-specific and Normal-specific modules are 3, 1 and 2 respectively among the down-regulated BRCA subtype-specific LMSM modules. In particular, only 1 module (LMSM 2) can act as both up-regulated and down-regulated BRCA subtype-specific LMSM module. In total, the unique number of BRCA subtype-specific LMSM modules is 9, indicating that most of LMSM modules are BRCA subtype-specific.

20

**A**

**Up-regulated BRCA subtype-specific LMSM modules**



**B**

**Down-regulated BRCA subtype-specific LMSM modules**



413
414 **Fig 3. Heatmap of the enrichment scores of BRCA subtype-specific LMSM modules in**
415 **five BRCA subtype samples**. (A) Up-regulated BRCA subtype-specific LMSM modules. (B)
416 Down-regulated BRCA subtype-specific LMSM modules.

417 **The performance of LMSM modules is significantly higher than baseline's performance**

418 **in classifying BRCA subtypes**

419 For the identified 17 LMSM modules, the average *Subset accuracy* and *Hamming loss* in

420 classifying BRCA subtypes is 0.7547 and 0.0892, respectively (details can be seen in S6

421 Data), The *Subset accuracy* and *Hamming loss* of the baseline are 0.3800 and 0.2480,

422 respectively. By using Welch's *t*-test method, the *Subset accuracy* achieved using the 17

423 LMSM modules is significantly larger (better) than the *Subset accuracy* of the baseline (*p*-

424 value < 2.20E-16), and the *Hamming loss* of the 17 LMSM modules is significantly smaller

425 (better) than the *Hamming loss* of the baseline (*p*-value < 2.20E-16). The better performance

426 than the baseline method indicates that LMSM modules are biological meaningful in

427 classifying BRCA subtypes.

21

428 **Several lncRNA-related miRNA sponge interactions are experimentally confirmed**

429 For the ground truth used in the validation, we have collected 581 experimentally validated

430 lncRNA-related miRNA sponge interactions associated with the matched lncRNA and

431 mRNA expression data (details in S4 Data). After we merge the sponge lncRNA-mRNA

432 pairs in the identified 17 LMSM modules, we have predicted 1471664 unique lncRNA-

433 related miRNA sponge interactions (details at https://github.com/zhangjunpeng411/LMSM).

434 For each LMSM module, the number of shared miRNAs, lncRNAs, mRNAs, predicted

435 lncRNA-related miRNA sponge interactions can be seen in S7 Data.

436     As shown in Table 3, there are 4 LMSM modules (LMSM 2, LMSM 3, LMSM 5 and

437 LMSM 8) containing 14 experimentally validated lncRNA-related miRNA sponge

438 interactions in total. It is noted that all the lncRNAs and mRNAs in these confirmed lncRNA-

439 related miRNA sponge interactions are BRCA-related genes, indicating they may have

440 potentially involved in BRCA.

441 **Table 3. Validated lncRNA-related miRNA sponge interactions**.

| Module ID | Validated lncRNA-related miRNA sponge interactions |
|---|---|
| LMSM 2 | *H19*: *HMGA2, H19:IGF2, H19:ITGB1, H19*: *TGFB1, H19*: *VIM, H19:RUNX1, H19:CDH13, H19:KLF4, H19:TGFBI, H19:VDR* |
| LMSM 3 | *LINC00152*: *MCL1* |
| LMSM 5 | *NEAT1*: *EMP2* |
| LMSM 8 | *LINC00324*: *BTG2, DLEU2*: *CCNE1* |

442 **LMSM is capable of predicting miRNA targets**

443 LMSM use high-confidence miRNA-target interactions as seeds to predict miRNA-target

444 interactions. A miRNA-mRNA or miRNA-lncRNA pair in a LMSM module has the potential

445 to be a miRNA-target pair for the following reasons. Firstly, at sequence level, the sponge

446 lncRNAs and mRNAs in each LMSM module have a significant sharing of miRNAs.

447 Secondly, at expression level, the sponge lncRNAs and mRNAs in each LMSM module are

22

448    highly correlated. As a result, the sponge lncRNAs and mRNAs of each LMSM module have

449    a high chance to be target genes of the shared miRNAs. Thus, based on the identified LMSM

450    modules, we have predicted 2820524 unique miRNA-target interactions (including 2023304

451    miRNA-lncRNA and 797220 miRNA-mRNA interactions) (details at

452    https://github.com/zhangjunpeng411/LMSM). For each LMSM module, the numbers of

453    predicted miRNA-lncRNA interactions and miRNA-mRNA interactions can be seen in S7

454    Data.

455    In addition, we investigate the intersection of the miRNA-target interactions predicted by

456    LMSM with the other well-cited miRNA-target prediction methods. In terms of miRNA-

457    mRNA interactions, we select TargetScan v7.2 [51], DIANA-microT-CDS v5.0 [52],

458    starBase v3.0 [53] and miRWalk v3.0 [54] for investigation. We choose starBase v3.0 [53]

459    and DIANA-LncBase v2.0 [39] for investigation in terms of miRNA-lncRNA interactions.

460    As shown in the UpSet plot [55] of Fig 4A, the number of miRNA-mRNA interactions

461    identified by all the five methods is only 21842. However, the percentage of overlap between

462    LMSM and each of the other four methods achieves ~63.74% (1289620 out of 2023304). As

463    shown in Fig 4B, the number of miRNA-lncRNA interactions identified by all the three

464    methods is only 1160. Since the miRNA-lncRNA interactions are still limited, most of the

465    miRNA-lncRNA interactions (~93.90%, 748609 out of 797220) are individually predicted by

466    LMSM.

**A**



**B**



**Fig 4. Overlaps and differences between predicted miRNA-target interactions by LMSM and other methods**. (A) Predicted miRNA-mRNA interactions between LMSM and TargetScan, DIANA_microT_CDS, starBase, miRWalk. (B) Predicted miRNA-lncRNA interactions between LMSM and starBase, DIANA_LncBase. Each column corresponds to an exclusive intersection that includes the elements of the sets denoted by the dark or red circles, but not of the others. The overlap size between different methods denotes exclusive overlaps, i.e. the overlap set not in a subset of any other overlap set.

**Comparison with graph clustering-based strategy**

24

476    Graph clustering-based strategy [12-17] is an alternative approach to identifying lncRNA

477    related miRNA sponge modules. As there is no graph clustering-based strategy specifically

478    designed for finding lncRNA related miRNA sponge modules, so we create a baseline Graph

479    Clustering-based method (called GC in this paper) which uses well-known network

480    construction and graph clustering methods as described in the following. The GS method

481    includes two steps: i) identifying lncRNA related miRNA sponge interaction network, and ii)

482    identifying lncRNA related miRNA sponge modules from the identified network. In step 1,

483    we adapt the well-cited Sensitivity Correlation (SC) method [24] implemented in the

484    *miRspongeR* R package [56] to infer lncRNA related miRNA sponge interaction network. A

485    lncRNA-mRNA pair is considered as an interacting pair in the network if they have

486    significant sharing of the miRNAs, significant correlation and adequate sensitivity correlation.

487    We require that the pairs must share at least 3 miRNAs and their sensitivity correlation (the

488    difference between correlation and partial correlation) must be larger than 0.1. The

489    statistically significance of the miRNA sharing and positive correlations are tested using

490    hypergeometric test and Welch's *t*-test respectively, with a significant level at 0.05. In step 2,

491    we use the well-cited Markov cluster (MCL) algorithm [57] to infer lncRNA related miRNA

492    sponge modules. Here, each obtained cluster corresponds to a module. Each module should

493    contain at least 2 sponge lncRNAs and 2 target mRNAs. In total, by using the GC method, we

494    have obtained 108 lncRNA related miRNA sponge modules.

495    We compare LMSM and GC in terms of the percentage of BRCA-related modules, the

496    percentage of module biomarkers in BRCA, the classification performance (mean *Subset*

497    *accuracy* and mean *Hamming loss*) in classifying BRCA subtypes, and the number of

498    validated lncRNA-related miRNA sponge interactions. As shown in Table 4, the comparison

499    result indicates that LMSM always performs better than the GC method. The detailed results

500    of the GC method can be seen in S8 Data.

25

501 **Table 4. Comparison results between LMSM and GC**.

| Method | %BRCA-related modules | %Module biomarkers | Mean *Subset accuracy* | Mean *Hamming loss* | #Validated interactions |
|---|---|---|---|---|---|
| LMSM | **58.82%** | **88.24%** | **0.7547** | **0.0892** | **14** |
| GC | 32.41% | 66.67% | 0.6586 | 0.1319 | 2 |

502 **LMSM is robust**

503 To demonstrate the robustness of the LMSM workflow, we use the sparse group factor

504 analysis (SGFA) method [58], instead of the WGCNA method to identify lncRNA-mRNA

505 co-expression modules. The SGFA method is extended from the group factor analysis (GFA)

506 method [59-61], and it can reliably infer biclusters (modules) from multiple data sources, and

507 provide predictive and interpretable structure existing in any subset of the data sources. Given

508 *B* biclusters to be identified, the SGFA method assigns each column (lncRNA or mRNA) or

509 row (sample) a grade of membership (association) belonging to these biclusters. The range of

510 the values of the associations is [-1, 1]. We use the absolute value of association (*AVA*) to

511 evaluate the strength of lncRNAs and mRNAs belonging to a bicluster, and the cutoff of *AVA*

512 is also set to 0.8. Specifically, we use the *GFA* R package [58] to identify lncRNA-mRNA

513 co-expression modules. The parameter settings for inferring lncRNA-related miRNA sponge

514 modules are the same.

515     By using the SGFA method, we have identified 51 LMSM modules (details can be seen in

516 S1 Data). The average size of these LMSM modules is 277.63 and the average number of the

517 shared miRNAs is 135.65. There are 490 unique miRNAs mediating the 51 LMSM modules,

518 and 84.90% (416 out of 490) miRNAs mediate at least two LMSM modules (details can be

519 seen in S2 Data). As the result obtained using the WGCNA method, the result with the SGFA

520 method also implies that the mediating miRNAs mostly act as crosslinks across different

521 LMSM modules. In addition, by using a null-model-based *p*-value computation method, the

522    identified 51 LMSM modules are also all statistically significant with adjusted $p$-value $\leq$

523    5.00E-06 (details can be seen in S3 Data).

524    As shown in S1 Table of S1 File, 3 out of the 51 LMSM modules are functionally

525    enriched in BRCA at a significant level ($p$-value < 0.05). Moreover, 49 out of the 51 LMSM

526    modules are regarded as module biomarkers in BRCA (see in S2 Table of S1 File). The

527    results indicate that most of LMSM modules are related to BRCA.

528    We also compute the enrichment scores of the identified 51 LMSM modules in the BRCA

529    subtype samples (details in S5 Data). As illustrated in S1 Fig of S1 File, out of the 51 LMSM

530    modules, 33 and 24 modules are regarded as up-regulated and down-regulated BRCA

531    subtype-specific LMSM modules, respectively. For the up-regulated BRCA subtype-specific

532    LMSM modules, the numbers of Basal-specific, Her2-specific, LumB-specific and Normal-

533    specific modules are 27, 2, 2 and 2, respectively. The numbers of Basal-specific, Her2-

534    specific, LumA-specific, LumB-specific and Normal-specific modules are 2, 3, 15, 3 and 1

535    respectively for the down-regulated BRCA subtype-specific LMSM modules. Particularly, 16

536    modules can act as both up-regulated and down-regulated BRCA subtype-specific LMSM

537    module. Overall, the unique number of BRCA subtype-specific LMSM modules is 41. This

538    result also indicates that the identified LMSM modules are mostly BRCA subtype-specific.

539    The average value of *Subset accuracy* and *Hamming loss* of the identified 51 LMSM

540    modules in classifying BRCA subtypes is 0.6921 and 0.1135, respectively (details can be

541    seen in S6 Data). In classifying BRCA subtypes, the baseline value of *Subset accuracy* and

542    *Hamming loss* is 0.3800 and 0.2480, respectively. By using Welch's $t$-test method, the value

543    of *Subset accuracy* for 51 LMSM modules is significantly larger (better) than the baseline

544    value of *Subset accuracy* ($p$-value < 2.20E-16), and the value of *Hamming loss* for 51 LMSM

545    modules is significantly smaller (better) than the baseline value of *Hamming loss* ($p$-value <

27

546    2.20E-16). The better performance than the baseline method also indicates that LMSM

547    modules are biological meaningful in classifying BRCA subtypes.

548    Moreover, we have predicted 605456 unique lncRNA-related miRNA sponge interactions

549    in the identified 51 LMSM modules (details at https://github.com/zhangjunpeng411/LMSM).

550    The number of the shared miRNAs, lncRNAs, mRNAs, predicted lncRNA-related miRNA

551    sponge interactions of each LMSM module can be seen in S7 Data. Since the experimentally

552    validated lncRNA-related miRNA sponge interactions are still limited, only 4 LMSM

553    modules containing 4 lncRNA-related miRNA sponge interactions (see S3 Table of S1 File)

554    are experimentally validated. All lncRNAs and mRNAs in the confirmed lncRNA-related

555    miRNA sponge interactions are also BRCA-related genes.

556    LMSM also has identified a large number of potential miRNA-target interactions

557    (1646449 in total, including 435345 miRNA-mRNA and 1211104 miRNA-lncRNA

558    interactions, details at https://github.com/zhangjunpeng411/LMSM). The number of

559    predicted miRNA-lncRNA interactions, predicted miRNA-mRNA interactions, putative

560    miRNA-lncRNA interactions and putative miRNA-mRNA interactions can be seen in S7

561    Data. As illustrated in S2 Fig of S1 File, the number of the miRNA-mRNA interactions

562    identified by all the five methods is 4897 and the number of the miRNA-lncRNA interactions

563    identified by all the three methods is 1149. Most of the identified miRNA-mRNA interactions

564    by LMSM (~58.55%, 254910 out of 435345) are also predicted by one of the other four

565    methods. In terms of the predicted miRNA-lncRNA interactions, ~94.23% (1141232 out of

566    1211104) miRNA-lncRNA interactions are also individually predicted by LMSM.

567    Finally, in terms of the percentage of BRCA-related modules, the percentage of module

568    biomarkers in BRCA, the classification performance (mean *Subset accuracy* and mean

569    *Hamming loss*) in classifying BRCA subtypes, and the number of validated lncRNA-related

28

570     miRNA sponge interactions, LMSM also generally performs better than the GC method (see

571     S4 Table of S1 File).

572     Altogether, the above results are consistent with those obtained using the WGCNA

573     method, indicating that our LMSM workflow is robust for studying lncRNA-related miRNA

574     sponge modules.

## Discussion

576     The crosstalk between different RNA transcripts in a miRNA-dependent manner forms a

577     complex miRNA sponge interaction network and depicts a novel layer of gene expression

578     regulation. Until now, several types of RNA transcripts, e.g. lncRNAs, pseudogenes,

579     circRNAs and mRNAs, have been confirmed to act as miRNA sponges. Since lncRNAs are a

580     large class of ncRNAs and function in many aspects of cell biology, including human

581     cancers, we focus on identifying lncRNA related miRNA sponge modules in this work.

582     By integrating multiple data sources, previous studies mainly investigate the identification

583     of lncRNA related miRNA sponge interaction network. Based on the identified lncRNA

584     related miRNA sponge interaction network, they use graph clustering algorithms to further

585     infer lncRNA related miRNA sponge modules. Different from existing computational

586     methods on lncRNA related miRNA sponge modules, in this work, we propose a novel

587     method named LMSM to directly identify lncRNA related miRNA sponge modules from

588     heterogeneous data. It is noted that the LMSM method depends on our presented hypothesis

589     of miRNA sponge modular competition. In the hypothesis, miRNA sponges tend to form a

590     group to compete with a group of target mRNAs for binding with miRNAs.

591      We have applied the LMSM method to the BRCA dataset from TCGA. For the putative

592      miRNA-target interactions, we integrate high-confidence miRNA-target interactions from

593      several databases. The analysis results demonstrate that our LMSM method is useful in

594      identifying lncRNA related miRNA sponge modules, and it can help with understanding

595      regulatory mechanism of lncRNAs.

596      LMSM is a flexible method to investigate miRNA sponge modules in human cancer.

597      Firstly, any biclustering or clustering algorithm (e.g. the joint non-negative matrix

598      factorization methods presented by Deng *et al*. [18] and Xiao *et al*. [19]) can be plugged in

599      stage 1 of LMSM to identify lncRNA-mRNA co-expression modules. The only condition for

600      using these algorithms is that they can be used to identify biclusters or clusters from high-

601      dimensional expression data. Secondly, LMSM is a parametric model, and the parameter

602      settings of LMSM can be replaced according to the practical requirements of researchers. For

603      example, the threshold of the three metrics in stage 2 for identifying lncRNA related miRNA

604      sponge modules can be looser or stricter. Thirdly, LMSM can also be extended to study other

605      ncRNA (e.g. circRNA and pseudogene) related miRNA sponge modules. For instance, if we

606      change the matched lncRNA expression data and the miRNA-lncRNA interactions to

607      matched circRNA expression data and the miRNA-circRNA interactions respectively, the

608      pipeline of LMSM is to identify circRNA related miRNA sponge modules.

609      It is noted that each LMSM module contains many sponge lncRNAs and mRNAs, so it is

610      hard to experimentally validate such a module by follow-up wet-lab experiments. This is a

611      common issue of existing computational methods, including LMSM. We suggest that

612      biologists can select some sponge lncRNAs and mRNAs of interest in each LMSM module,

613      and then validate the modular competition between the selected sponge lncRNAs and target

614      mRNAs. We believe that LMSM is still useful in shortlisting high-confidence sponge

615    lncRNAs and mRNAs for experimental validation. For example, previous study [62] has

616    shown that lncRNA *MIR22HG* is functionally complementary to lncRNA *H19*. In the

617    identified LMSM module no. 2 (LMSM 2), lncRNA *H19* is experimentally validated to

618    compete with 10 target mRNAs (*HMGA2*, *IGF2*, *ITGB1*, *TGFB1*, *VIM*, *RUNX1*, *CDH13*,

619    *KLF4*, *TGFBI* and *VDR*). Thus, biologists can select 2 lncRNAs (*H19* and *MIR22HG*) and 10

620    target mRNAs (*HMGA2*, *IGF2*, *ITGB1*, *TGFB1*, *VIM*, *RUNX1*, *CDH13*, *KLF4*, *TGFBI* and

621    *VDR*) in LMSM 2 to validate the modular competition between them.

622    Taken together, based on the hypothesis of miRNA sponge modular competition, we

623    propose a new approach to identifying lncRNA related miRNA sponge modules by

624    integrating expression data and miRNA-target binding information. Our method not only

625    extends the ceRNA hypothesis, but also provides a novel way to investigate the biological

626    functions and modular mechanism of lncRNAs in BRCA. We believe that our method can be

627    also applied to other human cancer datasets assists in human cancer research.

628    ## Supporting information

629    **S1 Data. The identified LMSM modules.**

630    **S2 Data. The distribution of the shared miRNAs in LMSM modules.**

631    **S3 Data. Statistically significant analysis results of LMSM modules.**

632    **S4 Data. BRCA-related genes and experimentally validated lncRNA related miRNA sponge interactions.**

633    **S5 Data. The enrichment scores of the identified LMSM modules in the BRCA subtype samples.**

634    **S6 Data. Classification analysis results of LMSM modules in classifying BRCA subtypes.**

635    **S7 Data. The number of shared miRNAs, lncRNAs, mRNAs, predicted interactions for each LMSM**

636    **module.**

637    **S8 Data. The results of a graph clustering-based strategy.**

638    **S1 File. Supporting file. Supplementary file.**

# References

640    1.  Derrien T, Johnson R, Bussotti G, Tanzer A, Djebali S, Tilgner H, et al. The GENCODE v7

641        catalog of human long noncoding RNAs: analysis of their gene structure, evolution, and

642        expression. Genome Res 2012; 22(9):1775-1789. doi: 10.1101/gr.132159.111.

643    2.  Fang Y, Fullwood MJ. Roles, functions, and mechanisms of long non-coding RNAs in cancer.

644        Genomics Proteomics Bioinformatics 2016; 14(1):42-54. doi: 10.1016/j.gpb.2015.09.006.

645    3.  Bhan A, Soleimani M, Mandal SS. Long noncoding RNA and cancer: a new paradigm.

646        Cancer Res 2017; 77(15):3965-81. doi: 10.1158/0008-5472.CAN-16-2634.

647    4.  Kopp F, Mendell JT. (2018) Functional classification and experimental dissection of long

648        noncoding RNAs. Cell 2018; 172(3):393-407. doi: 10.1016/j.cell.2018.01.011.

649    5.  Ambros V. The functions of animal microRNAs. Nature 2004; 431(7006):350-5. doi:

650        10.1038/nature02871.

651    6.  Bartel DP. MicroRNAs: genomics, biogenesis, mechanism, and function. Cell 2004;

652        116(2):281-97. doi: 10.1016/s0092-8674(04)00045-5.

653    7.  Salmena L, Poliseno L, Tay Y, Kats L, Pandolfi PP. A ceRNA hypothesis: the Rosetta Stone

654        of a hidden RNA language? Cell 2011; 146(3):353-8. doi: 10.1016/j.cell.2011.07.014.

655    8.  Tay Y, Rinn J, Pandolfi PP. The multilayered complexity of ceRNA crosstalk and competition.

656        Nature 2014; 505(7483):344-52. doi: 10.1038/nature12986.

657    9.  Zhou S, He Y, Yang S, Hu J, Zhang Q, Chen W, et al. The regulatory roles of lncRNAs in the

658        process of breast cancer invasion and metastasis. Biosci Rep 2018; 38(5):BSR20180772. doi:

659        10.1042/BSR20180772.

660    10. Peng F, Li TT, Wang KL, Xiao GQ, Wang JH, Zhao HD, et al. H19/let-7/LIN28 reciprocal

661        negative regulatory circuit promotes breast cancer stem cell maintenance. Cell Death Dis

662        2017; 8(1):e2569. doi: 10.1038/cddis.2016.438.

663    11. Le TD, Zhang J, Liu L, Li J. Computational methods for identifying miRNA sponge

664        interactions. Brief Bioinform 2017; 18(4):577-590. doi: 10.1093/bib/bbw042.

665    12. Shao T, Wu A, Chen J, Chen H, Lu J, Bai J, et al. Identification of module biomarkers from

666        the dysregulated ceRNA-ceRNA interaction network in lung adenocarcinoma. Mol Biosyst

667        2015; 11(11):3048-58. doi: 10.1039/c5mb00364d.

668    13. Zhang Y, Xu Y, Feng L, Li F, Sun Z, Wu T, et al. Comprehensive characterization of

669        lncRNA-mRNA related ceRNA network across 12 major cancers. Oncotarget 2016;
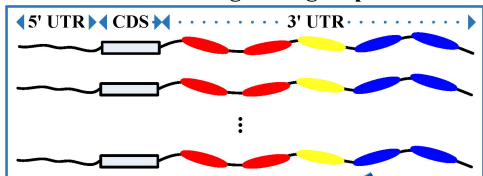
670        7(39):64148-67. doi: 10.18632/oncotarget.11637.

14. Zhang J, Le TD, Liu L, Li J. Inferring miRNA sponge co-regulation of protein-protein interactions in human breast cancer. BMC Bioinformatics 2017; 18(1):243. doi: 10.1186/s12859-017-1672-2.

15. Wang H, Xu D, Huang H, Cui Y, Li C, Zhang C, et al. Detection of dysregulated competing endogenous RNA modules associated with clear cell kidney carcinoma. Mol Med Rep 2018; 18(2):1963-72. doi: 10.3892/mmr.2018.9189.

16. Do D, Bozdag S. Cancerin: A computational pipeline to infer cancer-associated ceRNA interaction networks. PLoS Comput Biol 2018; 14(7):e1006318. doi: 10.1371/journal.pcbi.1006318.

17. Zhang J, Liu L, Li J, Le TD. LncmiRSRN: identification and analysis of long non-coding RNA related miRNA sponge regulatory network in human cancer. Bioinformatics 2018; 34(24):4232-40. doi: 10.1093/bioinformatics/bty525.

18. Deng J, Kong W, Wang S, Mou X, Zeng W. Prior knowledge driven joint NMF algorithm for ceRNA co-module identification. Int J Biol Sci 2018; 14(13):1822-1833. doi: 10.7150/ijbs.27555.

19. Xiao Q, Luo J, Liang C, Cai J, Li G, Cao B. CeModule: an integrative framework for discovering regulatory patterns from genomic data in cancer. BMC Bioinformatics 2019; 20(1):67. doi: 10.1186/s12859-019-2654-3.

20. Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. BMC Bioinformatics 2008; 9:559. doi: 10.1186/1471-2105-9-559.

21. Segal E, Shapira M, Regev A, Pe'er D, Botstein D, Koller D, et al. Module networks: identifying regulatory modules and their condition-specific regulators from gene expression data. Nat Genet 2003; 34(2):166-76. doi: 10.1038/ng1165.

22. Hotelling H. Relations between two sets of variates. Biometrika 1936; 28(3/4):321-377. doi: 10.2307/2333955.

23. Witten DM, Tibshirani R, Hastie T. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. Biostatistics 2009; 10(3):515-34. doi: 10.1093/biostatistics/kxp008.

24. Paci P, Colombo T, Farina L. Computational analysis identifies a sponge interaction network between long non-coding RNAs and messenger RNAs in human breast cancer. BMC Syst Biol 2014; 8:83. doi: 10.1186/1752-0509-8-83.

25. List M, Dehghani Amirabad A, Kostka D, Schulz MH. Large-scale inference of competing endogenous RNA networks with sparse partial correlation. Bioinformatics 2019; 35(14):i596-i604. doi:10.1093/bioinformatics/btz314.

705    26. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful
706        approach to multiple testing. J R Stat Soc Ser B (Methodological) 1995; 57(1): 289-300. doi:
707        10.2307/2346101.

708    27. Andersen P, Gill R. Cox's regression model for counting processes, a large sample study. Ann
709        Stat 1982; 10(4):1100-20. doi:10.1214/aos/1176345976.

710    28. Therneau TM, Grambsch PM. Modeling survival data: extending the Cox model. Springer-
711        Verlag, New York, 2000.

712    29. Parker JS, Mullins M, Cheang MC, Leung S, Voduc D, Vickery T, et al. Supervised risk
713        predictor of breast cancer based on intrinsic subtypes. J Clin Oncol 2009; 27(8):1160-7. doi:
714        10.1200/JCO.2008.18.1370.

715    30. Gendoo DM, Ratanasirigulchai N, Schröder MS, Paré L, Parker JS, Prat A, et al. Genefu: an
716        R/Bioconductor package for computation of gene expression-based signatures in breast cancer.
717        Bioinformatics 2016; 32(7):1097-9. doi: 10.1093/bioinformatics/btv693.

718    31. Hänzelmann S. GSVA: gene set variation analysis for microarray and RNA-seq data. BMC
719        Bioinformatics 2013; 14:7. doi: 10.1186/1471-2105-14-7.

720    32. Welch BL. The generalisation of student's problems when several different population
721        variances are involved. Biometrika 1947; 34(1-2):28-35. doi: 10.1093/biomet/34.1-2.28.

722    33. Tsoumakas G, Katakis I, Vlahavas I. Mining Multi-Label Data. In O. Maimon and L. Rokach,
723        editors, Data Mining and Knowledge Discovery Handbook, chapter 34, pages 667–685.
724        Springer-Verlag, 2 edition, 2010. ISBN 0387244352. doi:10.1007/978-0-387-09823-4_34.

725    34. Rivolli A, de Carvalho AC. The utiml package: Multi-label classification in R. The R Journal
726        2018; 10(2): 24-37. doi:10.32614/RJ-2018-041.

727    35. Chang CC, Lin CJ. LIBSVM: a library for support vector machines. ACM T Intel Syst Tec.
728        2011; 2(3):1-27. doi:10.1145/1961189.1961199.

729    36. Meyer D, Dimitriadou E, Hornik K, Weingessel A, Leisch F. e1071: Misc functions of the
730        department of statistics, probability theory group (Formerly: E1071), TU Wien, R package
731        version, 1.7-3. 2019. https://CRAN.R-project.org/package=e1071.

732    37. Metz J, de Abreu LF, Cherman EA, Monard MC. On the estimation of predictive evaluation
733        measure baselines for multi-label learning. In 13th Ibero-American Conference on AI, pages
734        189–198, Cartagena de Indias, Colombia, 2012. doi:10.1007/978-3-642-34654-5_20.

735    38. Hao Y, Wu W, Li H, Yuan J, Luo J, Zhao Y, et al. NPInter v3.0: an upgraded database of
736        noncoding    RNA-associated    interactions.    Database    (Oxford)    2016.    doi:
737        10.1093/database/baw057.

738    39. Paraskevopoulou MD, Vlachos IS, Karagkouni D, Georgakilas G, Kanellos I, Vergoulis T, et
739        al. DIANA-LncBase v2: indexing microRNA targets on non-coding transcripts. Nucleic Acids
740        Res 2016; 44(Database issue):D231-D238. doi: 10.1093/nar/gkv1270.
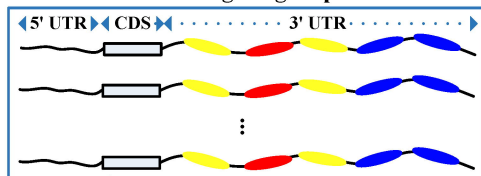
741    40. Huang HY, Lin YC, Li J, Huang KY, Shrestha S, Hong HC, et al. miRTarBase 2020: updates
742          to the experimentally validated microRNA-target interaction database. Nucleic Acids Res
743          2019. doi: 10.1093/nar/gkz896.

744    41. Vlachos IS, Paraskevopoulou MD, Karagkouni D, Georgakilas G, Vergoulis T, Kanellos I, et
745          al. DIANA-TarBase v7.0: indexing more than half a million experimentally supported
746          miRNA: mRNA interactions. Nucleic Acids Res 2015; 43(Database issue):D153-D159. doi:
747          10.1093/nar/gku1215.

748    42. Dweep H, Gretz N. miRWalk2.0: a comprehensive atlas of microRNA-target interactions. Nat
749          Methods 2015; 12(8):697. doi: 10.1038/nmeth.3485.

750    43. Piñero J, Bravo À, Queralt-Rosinach N, Gutiérrez-Sacristán A, Deu-Pons J, Centeno E, et al.
751          DisGeNET: a comprehensive platform integrating information on human disease-associated
752          genes and variants. Nucleic Acids Res 2017; 45(Database issue):D833-D839. doi:
753          10.1093/nar/gkw943.

754    44. Forbes SA, Beare D, Boutselakis H, Bamford S, Bindal N, Tate J, et al. COSMIC: somatic
755          cancer genetics at high-resolution. Nucleic Acids Res 2017; 45(Database issue):D777-D783.
756          doi: 10.1093/nar/gkw1121.

757    45. Bao Z, Yang Z, Huang Z, Zhou Y, Cui Q, Dong D. LncRNADisease 2.0: an updated database
758          of long non-coding RNA-associated diseases. Nucleic Acids Res 2019; 47(Database
759          issue):D1034-D1037. doi: 10.1093/nar/gky905.

760    46. Gao Y, Wang P, Wang Y, Ma X, Zhi H, Zhou D, et al. Lnc2Cancer v2.0: updated database of
761          experimentally supported long non-coding RNAs in human cancers. Nucleic Acids Res 2019;
762          47(Database issue):D1028-D1033. doi: 10.1093/nar/gky1096.

763    47. Cui T, Zhang L, Huang Y, Yi Y, Tan P, Zhao Y, et al. MNDR v2.0: an updated resource of
764          ncRNA-disease associations in mammals. Nucleic Acids Res 2018; 46(Database issue):D371-
765          D374. doi: 10.1093/nar/gkx1025.

766    48. Wang P, Zhi H, Zhang Y, Liu Y, Zhang J, Gao Y, et al. miRSponge: a manually curated
767          database for experimentally supported miRNA sponges and ceRNAs. Database (Oxford) 2015.
768          doi: 10.1093/database/bav098.

769    49. Pian C, Zhang G, Tu T, Ma X, Li F. LncCeRBase: a database of experimentally validated
770          human competing endogenous long non-coding RNAs. Database (Oxford) 2018. doi:
771          10.1093/database/bay061.

772    50. Wang P, Li X, Gao Y, Guo Q, Wang Y, Fang Y, et al. LncACTdb 2.0: an updated database of
773          experimentally supported ceRNA interactions curated from low- and high-throughput
774          experiments. Nucleic Acids Res 2019; 47(Database issue):D121-D127. doi:
775          10.1093/nar/gky1144.

776   51. Agarwal V, Bell GW, Nam JW, Bartel DP. Predicting effective microRNA target sites in
777       mammalian mRNAs. Elife 2015; 4. doi: 10.7554/eLife.05005.

778   52. Paraskevopoulou MD, Georgakilas G, Kostoulas N, Vlachos IS, Vergoulis T, Reczko M, et al.
779       DIANA-microT web server v5.0: service integration into miRNA functional analysis
780       workflows. Nucleic Acids Res 2013; 41(Web Server issue):W169-73. doi:
781       10.1093/nar/gkt393.

782   53. Li JH, Liu S, Zhou H, Qu LH, Yang JH. starBase v2.0: decoding miRNA-ceRNA, miRNA-
783       ncRNA and protein-RNA interaction networks from large-scale CLIP-Seq data. Nucleic
784       Acids Res 2014; 42(Database issue):D92-7. doi: 10.1093/nar/gkt1248.

785   54. Sticht C, De La Torre C, Parveen A, Gretz N. miRWalk: An online resource for prediction of
786       microRNA    binding    sites.    PLoS    One    2018;    13(10):e0206239.    doi:
787       10.1371/journal.pone.0206239.

788   55. Conway JR, Lex A, Gehlenborg N. UpSetR: an R package for the visualization of intersecting
789       sets    and    their    properties.    Bioinformatics    2017;    33(18):2938-2940.    doi:
790       10.1093/bioinformatics/btx364.

791   56. Zhang J, Liu L, Xu T, Xie Y, Zhao C, Li J, et al. miRspongeR: an R/Bioconductor package
792       for the identification and analysis of miRNA sponge interaction networks and modules. BMC
793       Bioinformatics 2019; 20(1):235. doi:10.1186/s12859-019-2861-y.

794   57. Enright AJ, Van Dongen S, Ouzounis CA. An efficient algorithm for large-scale detection of
795       protein families. Nucleic Acids Res. 2002; 30(7):1575–1584. doi:10.1093/nar/30.7.1575.

796   58. Bunte K, Leppäaho E, Saarinen I, Kaski S. Sparse group factor analysis for biclustering of
797       multiple    data    sources.    Bioinformatics    2016;    32(16):2457-63.    doi:
798       10.1093/bioinformatics/btw207.

799   59. Klami A, Virtanen S, Leppäaho E, Kaski S. Group factor analysis. IEEE Trans Neural Netw
800       Learn Syst 2015; 26(9):2136-2147. doi: 10.1109/TNNLS.2014.2376974.

801   60. Suvitaival T, Parkkinen JA, Virtanen S, Kaski S. Cross-organism toxicogenomics with group
802       factor analysis. Syst Biomed 2014; 2(4):71-80. doi: 10.4161/sysb.29291.

803   61. Virtanen S, Klami A, Khan S, Kaski S. Bayesian group factor analysis. In: Lawrence,N. and
804       Girolami,M. (eds), Proc. of the 15th International Conference on Artificial Intelligence and
805       Statistics, 2012; pp. 1269-1277.

806   62. Wang P, Ning S, Zhang Y, Li R, Ye J, Zhao Z, et al. Identification of lncRNA-associated
807       competing triplets reveals global patterns and prognostic markers for cancer. Nucleic Acids
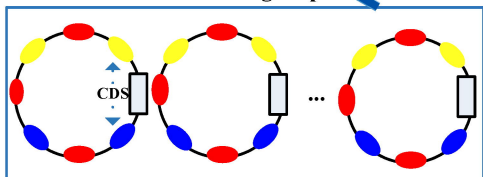808       Res 2015; 43(7):3478-89. doi: 10.1093/nar/gkv233.
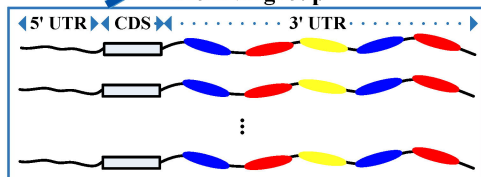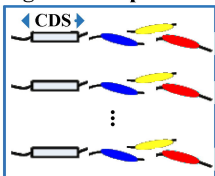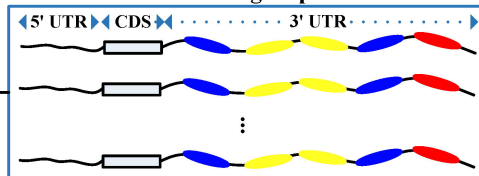
Other coding RNA group

Pseudogene group

5' UTR  CDS  3' UTR

CircRNA group

CDS

LncRNA group

5' UTR  CDS  3' UTR

miRNAs

MREs

Degradation products

CDS

mRNA group

5' UTR  CDS  3' UTR

Proteins

Win
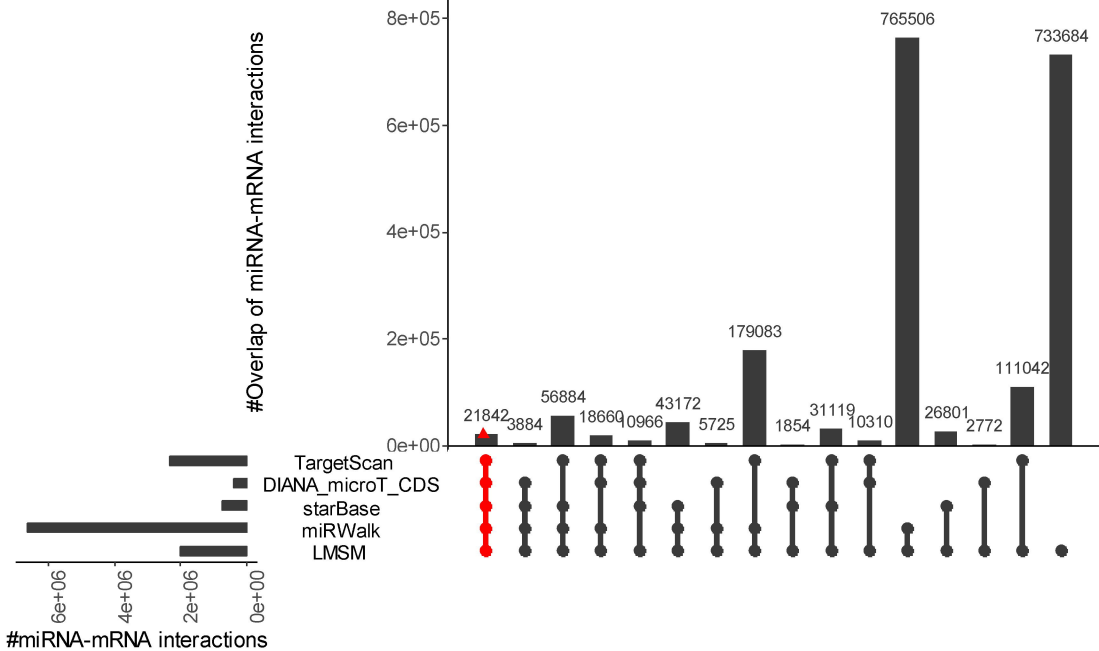
Lose

**A**

**Up-regulated BRCA subtype-specific LMSM modules**

**B**

**Down-regulated BRCA subtype-specific LMSM modules**

**A**



**B**