

Identifying Network Biomarkers of Cancer By Sample-Specific Differential Network

Yu Zhang

Shandong University

Xiao Chang

Anhui University of Finance & Economics

Jie Xia

Chinese Academy of Science

Yanhong Huang

Shandong University

Shaoyan Sun

Ludong University

Luonan Chen

University of Chinese Academy of Sciences

Xiaoping Liu (✉ xpliu@sdu.edu.cn)

University of Chinese Academy of Sciences

Research Article

Keywords: single-sample-differential-network gene-expression-data cancer driver-gene enrichment analysis

Posted Date: August 31st, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-677372/v1>

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

[Read Full License](#)

1 **Identifying network biomarkers of cancer by sample-specific differential**
2 **network**

3 **Yu Zhang^{1,2,3,#}, Xiao Chang^{4,#,*}, Jie Xia^{5,#}, Yanhong Huang³, Shaoyan Sun⁶, Luonan Chen^{1,}**
4 **^{2,5,7}, Xiaoping Liu^{1,2,3,*}**

5 ¹Key Laboratory of Systems Biology, Hangzhou Institute for Advanced Study, University of
6 Chinese Academy of Sciences, Hangzhou 310024, China

7 ²Key Laboratory of Systems Health Science of Zhejiang Province

8 ³School of Mathematics and Statistics, Shandong University, Weihai, Shandong, 264209,
9 China

10 ⁴Institute of Statistics and Applied Mathematics, Anhui University of Finance & Economics,
11 Bengbu, 233030, China.

12 ⁵Center for Excellence in Molecular Cell Science, Shanghai Institute of Biochemistry and Cell
13 Biology, Chinese Academy of Science, Shanghai 200031, China.

14 ⁶School of Mathematics and Statistics, Ludong University, Yantai 264025, China.

15 ⁷School of Life Science and Technology, ShanghaiTech University, Shanghai 201210, China.

16 ***To whom correspondence should be addressed. Xiao Chang, E-mail: chxlaugh@163.com;**

17 **Xiaoping Liu, E-mail: xpliu@sdu.edu.cn**

18 **#These authors contributed equally to this paper as the first authors.**

19 Number of Tables: 3

20 Number of Figures: 7

21

22 **ABSTRACT**

23 Abundant datasets generated from various big science projects on diseases have presented great
24 challenges and opportunities, which are contributed to unfold the complexity of diseases. The
25 discovery of disease- associated molecular networks for each individual plays an important

26 role in personalized therapy and precision treatment of cancer based on the reference networks.
27 However, there are no effective ways to distinguish the consistency of different reference
28 networks. In this study, we developed a statistical method, i.e. a sample-specific differential
29 network (SSDN), to construct and analyze such networks based on gene expression of a single
30 sample against a reference dataset. We proved that the SSDN is structurally consistent even
31 with different reference datasets if the reference dataset can follow certain conditions. The
32 SSDN also can be used to identify patient-specific disease modules or network biomarkers as
33 well as predict the potential driver genes of a tumor sample.

34 **Keywords:** single-sample-differential-network gene-expression-data cancer driver-gene
35 enrichment analysis

36

37 **Introduction**

38 With the rapid advance of deep sequencing technology for cancer genomes, several large-scale
39 projects, i.e. The Cancer Genome Atlas (TCGA)^{1, 2} and International Cancer Genome
40 Consortium (ICGC)^{3, 4}, were performed to provide the opportunities for the comprehensive
41 understanding of molecular mechanisms and pathogenesis underlying cancer^{5, 6}. One crucial
42 challenge for cancer omics data sets is to get insight into the mechanism of tumor progression⁷⁻
43 ⁹. The studies have shown that the molecular mechanisms of most complex diseases were due
44 to the dysfunction of relevant systems/networks instead of the malfunction of single
45 molecules¹⁰⁻¹². Therefore, constructing a network to analyze molecular mechanism has become
46 an effective method for studying complex diseases. The dynamic interactions and regulations
47 between molecules¹³⁻¹⁹ can detect the causal disease genes/module biomarkers at a single
48 sample level. The edge biomarker's method^{20, 21} calculates the difference between normal
49 network and disease network, and discovers a set of differentially correlated gene pairs. Besides,
50 network biomarkers²² or subnetwork markers^{14, 23} can accurately characterize disease states. In

51 actuality, these methods construct the individual-specific network based on the reference
52 network from a group of reference samples. However, it is unclear how can a different reference
53 sample set or the reference network affect the structure of single-sample network, or if or not
54 a different reference sample set can result in a different network structure. In other words, these
55 methods cannot distinguish the consistency of single-sample network with different reference
56 sample sets.

57 The SSN method²⁴ estimates the perturbations of Pearson's correlation coefficient (PCC) for
58 each pair of genes in a single sample, and it can be used to construct the individual-specific
59 network for disease samples and control samples, which called the Disease network and
60 Control network. Compared with the previously described SSN, we establish the SSDN by
61 comparing and finding the difference between the Disease network and Control network. A
62 reference sample set is required to construct a reference network in SSDN, and the consistency
63 of single-sample-Pearson correlation coefficients (*s-PCC*) needs to be considered in the SSDN.
64 For this consideration, we analyzed the conditions of consistency of *s-PCC* based on different
65 reference networks in this work, and proved that the *s-PCC* based on different reference
66 networks are consistent in the following two cases: the number of reference samples is
67 sufficiently large; the reference sample sets follow the same distribution. In other words,
68 provided that if one of these two conditions is satisfied for the reference samples, we have the
69 same SSDN structure, which is independent of the choice of the reference samples. This result
70 provides a theoretical foundation on determining the reference network in the construction of
71 SSDN.

72 In this work, we first gave the theoretical result on the conditions of the reference samples
73 to construct a consistent SSDN, and then validated the consistency of *s-PCC* based on different
74 reference networks both by simulated data and by three gastric cancer datasets from GEO
75 datasets from TCGA database. For clarifying the sample-specific characteristics of SSDN, we

76 established a disease-specific sample network (DSSN), which is similar to SSDN but is based
77 on non-paired sample data, to identify potential sample-specific driver genes to assess clinic
78 prognosis information, which is strongly correlated with individual somatic mutation genes
79 and validated by the enrichment analysis. The results of survival analysis for the potential
80 sample-specific driver genes demonstrate that the networks with those genes can be used as
81 effective module biomarkers to predict the prognosis for patients.

82

83 **MATERIAL AND METHODS**

84 *Data processing*

85 The gene expression profiles for gastric cancer were from the GEO database
86 (<http://www.ncbi.nlm.nih.gov/geo/>) including datasets GSE27342, GSE63089, and GSE33335.

87 The three datasets contain 80 pairs, 45 pairs and 25 pairs from gastric cancer tissues and
88 matched adjacent tumor-adjacent tissues from 150 cancer patients. All profiles were
89 normalized by the RMA (robust multi-array averaging) methods, and the probe sets were
90 mapped to their corresponding gene symbols. The expression values of replicated probe sets
91 were averaged to one gene. As a result, 17,325 genes were gotten for the following study. In
92 addition, four tumor datasets, which were Breast invasive carcinoma (BRCA), Lung
93 adenocarcinoma (LUAD), Lung squamous cell carcinoma (LUSC) and Liver Hepatocellular
94 Carcinoma (LIHC), were gotten from the TCGA data portal (<http://cancergenome.nih.gov>).

95 There were 1102 tumor and 113 tumor-adjacent samples in BRCA, 533 tumor and 59 tumor-
96 adjacent samples in LUAD, 502 tumor and 49 tumor-adjacent samples in LUSC, 371 tumor
97 and 50 tumor-adjacent samples in LIHC and the clinic information of these samples were also
98 downloaded from TCGA. Then, 24,991 mRNAs/genes were obtained for each sample in TCGA
99 and the data of tumor-adjacent tissue was considered as the normal samples for further study.

100 Finally, we obtained 50 tumor samples for BRCA from the ICGC database (International

101 Cancer Genome Consortium, <https://icgc.org/>) as a follow-up verification.

102

103 ***Functional enrichment for the individual specific network***

104 The existing cancer genes were gathered from the Cancer Gene Census database²⁵ (CGC,
105 <https://cancer.sanger.ac.uk/census/>) and a hypergeometric test was used to calculate the
106 functional enrichment of genes in the SSDN. The formula of the hypergeometric test is:

$$107 \quad P(X \geq k) = 1 - \sum_{i=0}^{k-1} \frac{\binom{K}{i} \binom{N-K}{n-i}}{\binom{N}{n}},$$

108 where N is the number of genes of the gene expression profiles, K is the number of genes
109 existing cancer-related genes in the CGC database, n is the number of genes in the SSDN of a
110 single sample and i is the number of overlapped genes between K and n . P is the statistical
111 significance of the hypergeometric test. If $P < 0.05$, then we regarded that the enrichment for
112 CGC database considered statistically significant. In addition, the enrichment analysis of genes
113 in the SSDN was conducted using DAVID Bioinformatics Tool (version 6.8,
114 <https://david.ncifcrf.gov/home.jsp>)²⁶ in the cancer pathway from the KEGG (Kyoto
115 Encyclopedia of Genes and Genomes).

116

117 ***Survival analysis for the individual specific network***

118 In order to confirm whether the genes from SSDN are related to disease, we used them as a
119 network biomarker to observe the effect between gene expression and survival rate in samples.

120 Here we defined the hub gene, which is a gene that is highly connected with others, or a gene
121 with a high degree. First, we computed the top m highest degree genes for SSDN of all samples
122 composing of the hub genes in one cancer. Second, for a single sample, if the top n highest
123 degree genes of this sample included half of the hub genes, then the gene was chosen into high-

124 risk group, on the contrary, it would be taken into low-risk group. Survival analysis was
 125 performed on the disease samples based on the hub genes. Furthermore, the log-rank test (with
 126 $p < 0.05$ considered significantly) in R/Bioconductor²⁷ was used to evaluate the statistically
 127 significant in the survival curves between the high and low-risk groups. An independent data
 128 from ICGC database were used to validate our results.

129

130 *The theoretical foundation of SSN based on different reference networks*

131 Assume that $X = [x_1, \dots, x_n]$ and $Y = [y_1, \dots, y_n]$ are two expression vectors for gene X and Y
 132 in reference samples (reference I) with length n , where x_i is the expression of gene X for the
 133 i th ($1 \leq i \leq n$) sample in reference samples, and y_j is the expression of gene Y for the j th
 134 ($1 \leq j \leq n$) sample in reference samples. Here, n can be considered as the number of the
 135 reference samples, x_i and y_j represent the expression values of two gene X and Y in
 136 reference samples. And the PCC for gene X and Y can be calculated as follows.

$$137 \quad R_n = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)}{\sqrt{\sum_{i=1}^n (x_i - \bar{x}_n)^2 \sum_{i=1}^n (y_i - \bar{y}_n)^2}}.$$

138 There were two new samples S_a with expression (x_a, y_a) and S_b with expression
 139 (x_b, y_b) for gene X and Y . The two samples were added into reference samples to form new
 140 vector pairs $[(X, x_a), (Y, y_a)]$ and $[(X, x_b), (Y, y_b)]$. Then the PCC s between vectors
 141 (X, x_a) and (Y, y_a) , between vectors (X, x_b) and (Y, y_b) with the length $(n+1)$ were
 142 calculated as R_{na} and R_{nb} . The differences of PCC s between before and after adding the new
 143 samples were $\Delta_{na} = R_{na} - R_n$ and $\Delta_{nb} = R_{nb} - R_n$.

144 Then we have another two reference vectors $X' = [x'_1, x'_2, \dots, x'_m]$ and $Y' = [y'_1, y'_2, \dots, y'_m]$

145 with length m for another reference samples (reference2), where x'_i is the i th ($1 \leq i \leq m$)
 146 element of gene X' and y'_j is the j th ($1 \leq j \leq m$) element of gene Y' in reference2. Here,
 147 m can be considered to be the number of the reference samples, x'_i and y'_j represent the
 148 expression levels of two molecules X' and Y' respectively. When the same two new
 149 samples S_a and S_b added to X' and Y' , the differences of PCC are $\Delta_{ma} = R_{ma} - R_m$ and
 150 $\Delta_{mb} = R_{mb} - R_m$.

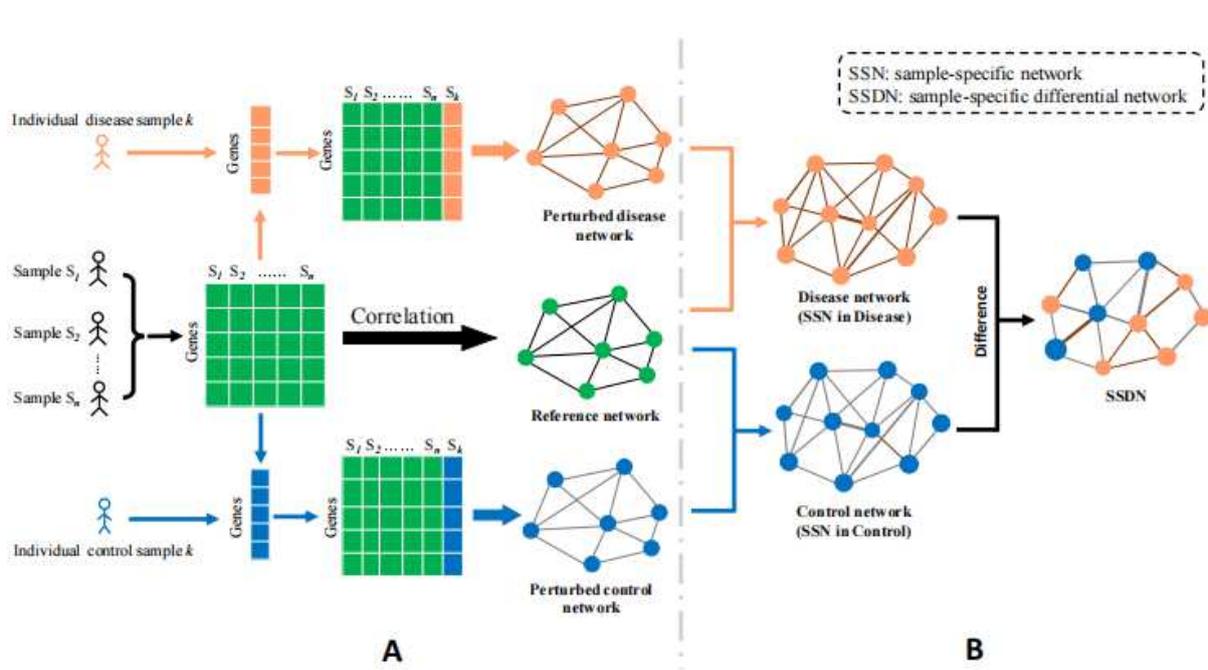
151 Derived from our mathematical theory (Supplementary Data Notes S1), if given the
 152 relationship of Δ_{na} and Δ_{nb} , then the relationship of Δ_{ma} and Δ_{mb} got two conclusions:
 153 One, assuming $\Delta_{na} > \Delta_{nb}$ based on reference1, if $n, m \rightarrow \infty$, we can get $\Delta_{ma} > \Delta_{mb}$ based on
 154 reference2, vice versa. Another, if vector X and vector X' belong to one independent
 155 identically distributed random variables $\{S_n\}$, Y and vector Y' belong to one independent
 156 identically distributed random variables $\{W_n\}$, and $\Delta_{na} > \Delta_{nb}$, then we got $\Delta_{ma} > \Delta_{mb}$. The
 157 details of mathematical explanations of the two conclusions for a single sample are given in
 158 Supplementary Data Notes S1. For the convenience, if $\Delta_{na} > \Delta_{nb}$ (or $\Delta_{na} < \Delta_{nb}$) in reference1,
 159 and $\Delta_{ma} > \Delta_{mb}$ (or $\Delta_{ma} < \Delta_{mb}$) in reference2, that means, Δ_{na} and Δ_{nb} , Δ_{ma} and Δ_{mb}
 160 have the same relationship, we defined it as single-sample-Pearson correlation coefficients,
 161 which implies as s - PCC in the following paper.

162

163 ***Constructing an individual-specific differential network***

164 The sample-specific network for an individual patient is constructed based on the statistical
 165 perturbation analysis of this sample against a group of given control samples. So, we required
 166 expression profiles for a group of normal samples, which served as the reference/control
 167 samples. We construct a reference network by Pearson correlation coefficients (PCC) using the

168 reference samples (Figure 1A). We calculate the *PCC* of each pairs of genes as an edge with or
169 without a background²⁴. Then, a disease sample k obtained from cancer tissues of a patient was
170 added to the reference samples and construct a Perturbed disease network by *PCC* (Figure 1A).
171 After that, a Disease network for disease sample k can be obtained by calculating the different
172 edges between the Perturbed disease network and Reference network, and a Disease network
173 is an SSN for the disease sample k in disease status (Figure 1B). At the same time, we also add
174 a control sample k , which was obtained from normal tissue of the same patient to the reference
175 network to construct the perturbed control network (Figure 1A) and Control network (Figure
176 1B) through the same procedure. The Control network is an SSN from a control sample. The
177 disease sample k was from the tumor tissue of patient k , and the control sample k was from the
178 tumor-adjacent tissue of patient k . The difference between the Disease network and Control
179 network is probably due to cancer-related genes. If the changes between two networks in terms
180 of the network structure are obvious, the genes that caused the changes are highly possible to
181 be cancer-related. On the contrary, if genes are insignificantly changing in the structure
182 between two networks, these are likely not to be the cancer-related. Thus, this new network
183 was called sample-specific differential network (SSDN) for sample k by obtaining the
184 differences between above two networks (Figure 1B), i.e. for an edge in the Disease network,
185 if it is not in the Control network, then the edge was kept in final SSDN, and vice versa.



186

187 **Fig 1. Construction of individual specific different network.** (A) A reference network is
 188 established through a set of control samples by *PCC*. A patient sample *k* includes a paired tumor
 189 sample (disease sample) and an adjacent normal tissues sample (control sample). A disease
 190 sample *k* is added to the reference network and constructed a Perturbed disease network (SSN
 191 in Disease). A Perturbed control network (SSN in Control) is obtained by a control sample *k* in
 192 the same way. The SSN is sample-specific network. (B) Based on the theory of predecessors,
 193 the Disease network and Control network is constituted by the significant of each edge. The
 194 Sample-Specific Differential Network (SSDN) is constructed by quantifying above two
 195 networks' differences.

196

197 In this study, a protein-protein interaction network (PPIN) from HPRD database (Human
 198 Protein Reference Database, <http://www.hprd.org/>) was used as the background network to
 199 filter the potential false-positive edges from the correlation networks. If there is an edge in
 200 PPIN for a pair of genes, the *PCC* of the gene pair would be calculated for Reference network,
 201 Perturbed disease network, and Perturbed normal network. If there is no connection in PPIN
 202 for a pair of genes, we ignored the calculation of *PCC* for the gene pair. We used the background

203 network in order to reduce the amount of calculation of the *PCC* and reduce the existence of
204 false positive gene pairs. (Supplement Data Notes S2)

205

206 **RESULTS**

207 *Numerical simulation of s-PCC based on different reference networks*

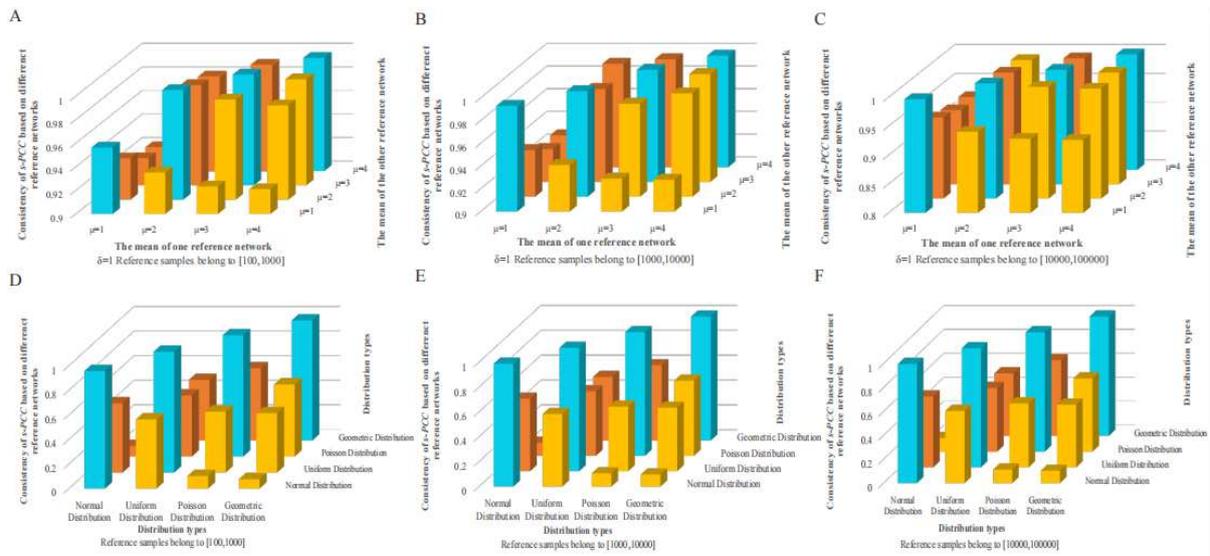
208 In order to verify our conclusion, a numerical simulation was done for a single sample
209 correlation. In the above paper, we have proved that *PCC* is consistent in two cases, and defined
210 the single-sample-Pearson correlation coefficients as *s-PCC*. Firstly, two simulated sample sets
211 were generated based on normal distributions from different mean, variance, and sample size
212 as two reference datasets, and the two reference datasets were called *reference1* and *reference2*.
213 Secondly, the *s-PCC* between two genes can be calculated based on the two reference datasets
214 (*reference1* and *reference2*), to obtain *s-PCC1* and *s-PCC2* in a simulated single sample.
215 Assuming the gene-pair (x, y) and (x', y') in the single sample, the correlations of the two
216 gene-pairs based on *reference1* are $s-PCC1_{(x, y)}$ and $s-PCC1_{(x', y')}$, the two gene-pairs based on
217 *reference2* are $s-PCC2_{(x, y)}$, and $s-PCC2_{(x', y')}$. If $s-PCC1_{(x, y)} > s-PCC1_{(x', y')}$ and $s-PCC2_{(x, y)} > s-$
218 $PCC2_{(x', y')}$, or $s-PCC1_{(x, y)} < s-PCC1_{(x', y')}$ and $s-PCC2_{(x, y)} < s-PCC2_{(x', y')}$, means the tendency of
219 the *s-PCC* for the two gene-pairs is consistent based on the two reference datasets. We regarded
220 the two gene-pairs as consistent gene-pairs. Then the consistency of *s-PCC* for two reference
221 datasets was defined as the percentage of consistent gene-pairs between the two reference
222 datasets. Finally, we evaluate the consistency of *s-PCC* among different reference datasets.

223 The two reference datasets were respectively generated from the normal distribution with
224 the mean value ($\mu=1, 2, 3, 4$) and the variance ($\delta=1$), and the sample size of the two reference
225 datasets were same and randomly obtained from a range. The consistency of *s-PCC* based on
226 different sample size of reference dataset was shown in Figure 2A-C (random value from range
227 100 to 1000, range 1000 to 10000, and range 10000 to 100000). If two reference datasets

228 generated from same distribution (same mean and variance), the consistency of *s-PCC* would
229 be higher than in other situations that the two reference datasets came from different
230 distributions (different mean) (Figure 2A-C). For example, the two reference datasets range
231 from 100 to 1000, and generated from same normal distribution with same mean ($\mu=1$) and
232 variance ($\delta=1$), the consistency of *s-PCC* is 95.64% (Figure 2A). When the two reference
233 datasets generated from different distributions with different mean ($\mu=1$ and $\mu=2$) and same
234 variance ($\delta=1$), the consistency of *s-PCC* is 93.51% (Figure 2A). When the two reference
235 datasets generated from the distributions $\mu=1$ and $\mu=3$ and same variance ($\delta=1$), the consistency
236 of *s-PCC* is 92.33% (Figure 2A). When the two reference datasets generated from the
237 distributions $\mu=1$ and $\mu=4$ and same variance ($\delta=1$), the consistency of *s-PCC* is 92.1% (Figure
238 2A). The more different the distributions of the two reference datasets generate from, the lower
239 consistency of *s-PCC* for the two reference datasets is. The same tendency was also shown in
240 Figure 2B and 2C. And with the increase of sample size of reference datasets, the consistency
241 of *s-PCC* was also raised from range 100 to 1000, range 1000 to 10000, and range 10000 to
242 100000 (Figure 2A-C). If two reference datasets generated from the different distributions
243 (Normal Distribution, Uniform Distribution, Poisson Distribution, Geometric Distribution), a
244 similar tendency was also shown in Figure 2D-F. For example, the two reference datasets range
245 from 100 to 1000, when the reference datasets both generated from Normal Distributions, the
246 consistency of *s-PCC* is 95.58% (Figure 2D). If one reference dataset generated from Normal
247 Distribution, the other generated from Uniform Distribution, the consistency of *s-PCC* is 56.23%
248 (Figure 2D). If one reference dataset generated from Normal Distribution, the other generated
249 from Poisson Distribution, the consistency of *s-PCC* is 10.23% (Figure 2D). If one reference
250 datasets generated from Normal Distribution, the other generated from Geometric Distribution,
251 the consistency of *s-PCC* is 7.65% (Figure 2D). The same tendency was also shown in Figure
252 2E-F. If two reference datasets both generated from Normal Distributions, with the increase of

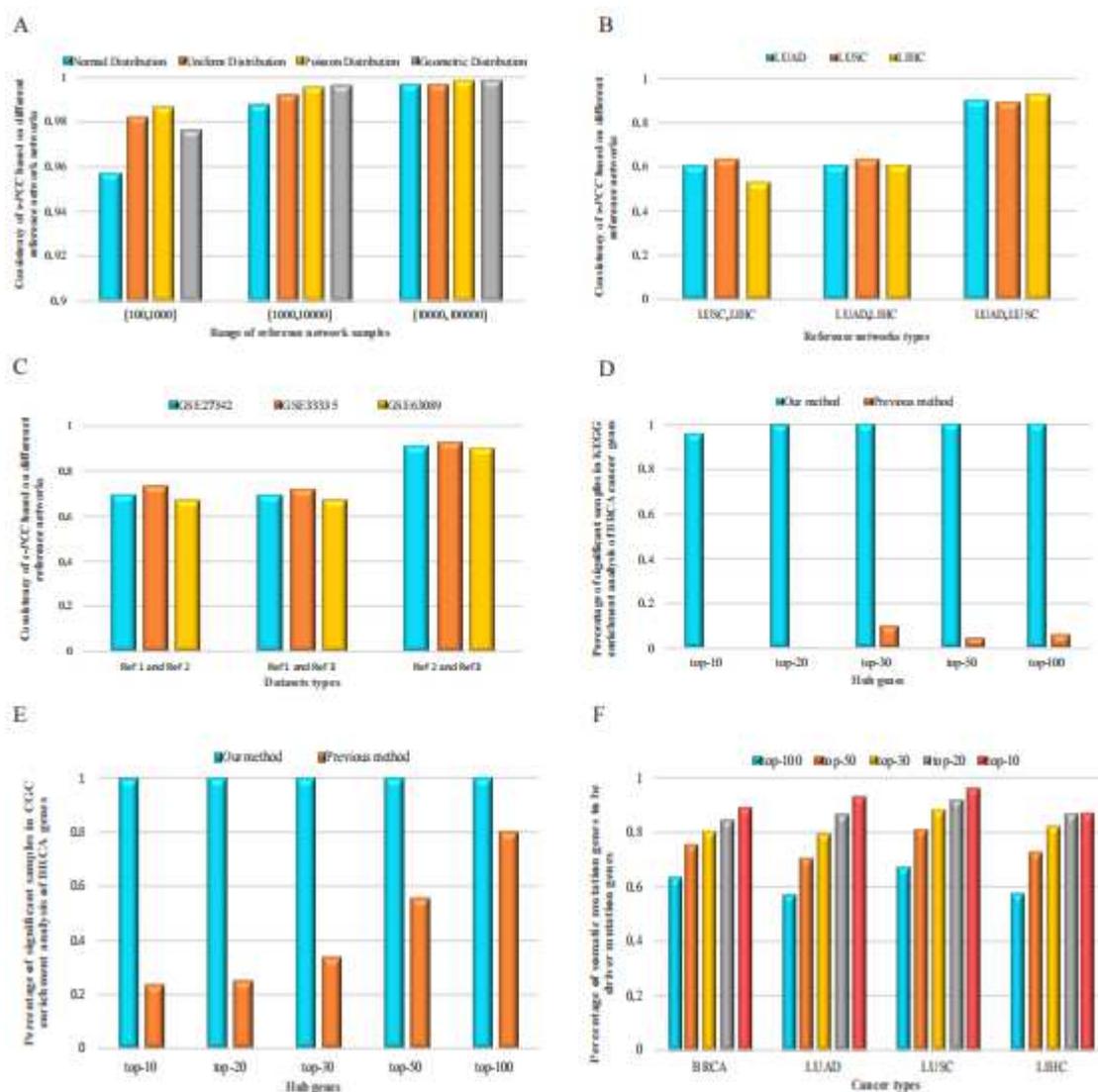
253 sample size of reference datasets, the consistency of s -PCC was also raised from range 100 to
 254 1000, range 1000 to 10000, and range 10000 to 100000 (Figure 3A). The results of numerical
 255 simulation showed that the consistency of s -PCC would reduce with the different distribution
 256 of the reference datasets, and raise with the increase of sample size of the reference datasets. It
 257 is also consistent with the theoretical analysis in the last section.

258



259

260 **Fig 2. Simulation data verify our conclusions.** (A) The variance of two normal distributions
 261 is one. Two reference datasets ranging from 100 to 1000. (B) The variance of two normal
 262 distributions is one. Two reference datasets ranging from 1000 to 10000. (C) The variance of
 263 two normal distributions is one. Two reference datasets ranging from 10000 to 100000. (D)
 264 Randomly generated four distributions to verify our conclusion. Two reference datasets ranging
 265 from 100 to 1000. (E) Randomly generated four distributions to verify our conclusion. Two
 266 reference datasets ranging from 1000 to 10000. (F) Randomly generated four distributions to
 267 verify our conclusion. Two reference datasets ranging from 10000 to 100000.



268

269

270

271

272

273

274

275

276

277

278

Fig 3. Validating sample-individual differential networks and predicting driver genes in cancer. (A) Four distribution samples ranging from 100 to 1000, 1000 to 10000 and 10000 to 100000. (B) Cross validation of three cancer. (C) Cross validation of three gastric cancer databases. (D) The proportion of significant samples in the enrichment analysis of top 100, 50, 30, 20 and 10 highest degree genes for BRCA DSSN in the KEGG pathway and compare with previous method (SSN). (E) The proportion of significant samples in the enrichment analysis of top 100, 50, 30, 20 and 10 highest degree genes for BRCA DSSN in the CGC database and compare with previous method (SSN). (F) The proportion of somatic mutation genes to be driver mutation genes in top 100, 50, 30, 20 and 10 highest degree genes for each DSSN.

279 ***Real data validation for the consistency of s -PCC in different reference sets***

280 In addition to the simulated data, three tumor datasets (LUAD, LUSC, and LIHC) were
281 obtained from TCGA (<https://www.cancer.gov/>) database to validate the results. The
282 control/normal samples (more than ten samples) were randomly selected from the three
283 datasets to form three reference sample sets, and used to construct three reference networks.
284 Each tumor sample in the three tumor datasets constructed Perturbed disease networks based
285 on the three reference networks. The average consistency of s -PCC was calculated based on
286 different reference networks, and the random selection from Perturbed disease networks for
287 different gene pairs was repeated 10^5 times. We regarded that the normal samples from the
288 same tissue follow the same distributions, so the normal samples from Lung adenocarcinoma
289 (LUAD) and Lung squamous cell carcinoma (LUSC) follow the same distribution, and Liver
290 Hepatocellular Carcinoma (LIHC) follows other distribution compare with LUAD and LUSC.
291 The results showed that if we only used the LUAD dataset as reference sample sets to construct
292 reference networks, the average consistency of s -PCC is 91.8%. If we only used the LUSC
293 dataset as reference sample sets to construct reference networks, the average consistency of s -
294 PCC is 92%. If we only used the LIHC dataset as reference sample sets to construct reference
295 networks, the average consistency of s -PCC is 92.6%. If we used the reference sample sets
296 from LUAD and LUSC, the average consistency is 90.4% (Figure 3B). While changing the
297 reference sample sets to LUAD and LIHC, the average consistency is changed to 61.61%
298 (Figure 3B). The average consistency is 60.66% when the reference sample sets were LUSC
299 and LIHC (Figure 3B). The results showed when the reference networks obey the same
300 distribution, the consistency of s -PCC will be higher than the reference networks with a
301 different distribution (Figure 3B).

302 Here we also used three gastric cancer databases from the GEO database
303 (<https://www.ncbi.nlm.nih.gov/geo/>) as an example, these datasets are GSE33335, GSE63089,

304 and GSE27342. We selected the control/normal samples of the three datasets as the reference
305 samples to construct three reference networks, and each tumor sample was used to construct
306 Perturbed disease networks based on the three reference networks. For convenience, we noted
307 the reference network from GSE33335 as *Ref1*, from GSE63089 as *Ref2* and from GSE27342
308 as *Ref3*. The number of reference samples in GSE33335 was 24, GSE63089 was 45, GSE27342
309 was 80. The average consistency of *s-PCC* was calculated based on different reference
310 networks and the random selection from Perturbed disease networks for different gene pairs
311 was repeated 10^5 times. The results were shown that if taken *Ref1* and *Ref2* as reference
312 networks, constructed the Perturbed disease networks for GSE27342, the consistency of *s-PCC*
313 was 69.34%. constructed the Perturbed disease networks for GSE 33335, the consistency of *s-*
314 *PCC* was 73.28%, constructed the Perturbed disease networks for GSE63089, the consistency
315 of *s-PCC* was 67.18% (Table 1 and Figure 3C). When the Perturbed disease networks were
316 constructed taken *Ref1* and *Ref3* as reference network, the consistency of the three datasets was
317 similar to the consistency based on *Ref1* and *Ref2* (Table 1 and Figure 3C). And when the
318 Perturbed disease networks were constructed taken *Ref2* and *Ref3* as reference network, the
319 consistency would be rapidly increased by over 90% (Table 1 and Figure 3C). It is an agreement
320 with the theoretical derivation that the consistency of *s-PCC* would be raised with the increase
321 number of reference samples.

322

323 **Table 1: the comparison for the consistency of different reference samples**

Dataset	<i>Ref1</i> and <i>Ref2</i>	<i>Ref1</i> and <i>Ref3</i>	<i>Ref2</i> and <i>Ref3</i>
GSE27342	69.34%	69.23%	90.71%
GSE33335	73.28%	71.80%	92.55%
GSE63089	67.18%	67.35%	90.20%

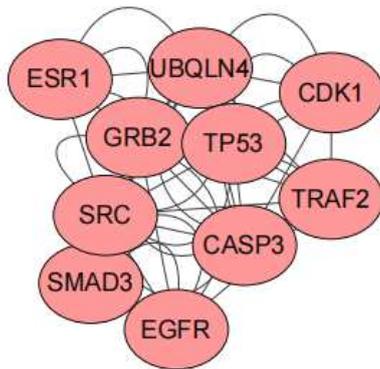
324

325 ***DSSN reveal individual features by pathway and disease gene enrichment***

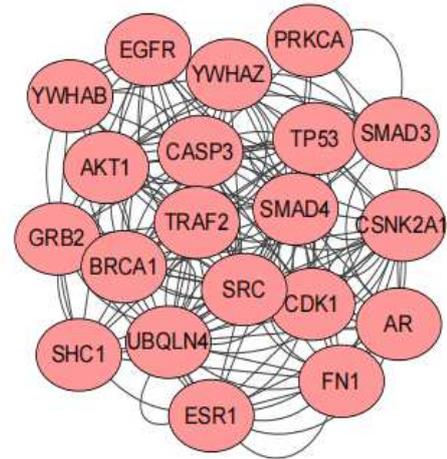
326 For revealing the disease modules of non-paired tumor samples, a common normal network
327 was constructed by collecting the common edges of Control networks, the edges existed in
328 more than 1/3 Control networks for a type of cancer. The common normal network was used
329 as the Control network to deduce the patient-specific disease modules. A disease-specific
330 sample network (DSSN) was established by identifying the differential edges between the
331 Disease network and the Control network (Figure 1). That is, for an edge in the Disease network,
332 if it is not in the control network, then the edge was kept in the final DSSN. The hub nodes of
333 DSSN were the potential cause modules of this tumor sample, and then the top- 100, 50, 30,
334 20 and 10 hub genes with a high degree in DSSN were respectively selected as potential disease
335 modules for every tumor sample in TCGA.

336 We chose the Breast invasive carcinoma (BRCA) and LIHC to draw the disease modules
337 networks. Disease modules reflect different extent of aggregation in different networks. In
338 BRCA reference network, we selected the top- 10 and 20 hub genes as potential disease
339 modules, and calculated the *PCC* between these genes. The gene pairs with p-value less than
340 0.01 form the edges of the network (Figure 4). In BRCA Control network, the top- 10 disease
341 modules cannot be aggregated, while are scattered by several modules (Figure 5A). Because
342 our modules are selected from Disease network hub genes, so it cannot be significantly
343 aggregated in the Control network. In BRCA Disease network, we selected three individual
344 samples to view the network (Figure 5B-D). The top- 10 disease modules are significantly
345 aggregated in the Disease network. The top 20- disease modules in Disease network are also
346 shown in Figure S1. In the same way, the top- 10 and 20 disease modules in LIHC are drawn
347 in Figure S2-S4. It also implies that hub genes existed in the form of modules in the Disease
348 network.

A



B

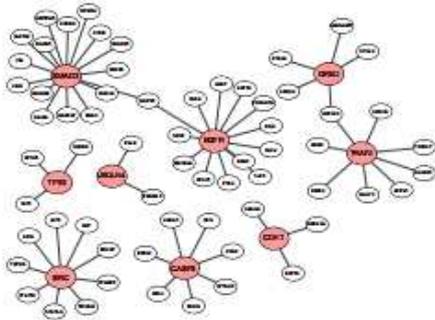


349

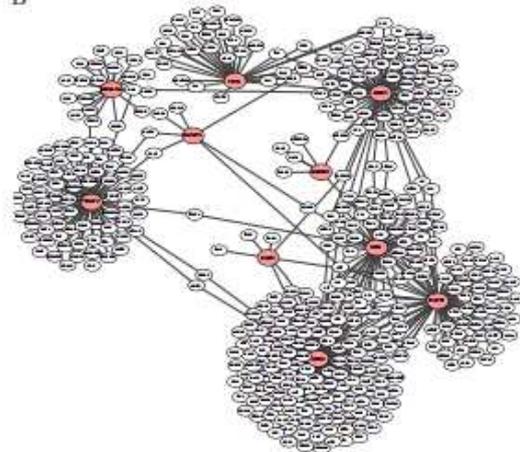
350 **Fig 4. The potential disease modules in BRCA reference network. (A) The network modules**

351 **among the top- 10 hub genes. (B) The network modules among the top- 20 hub genes.**

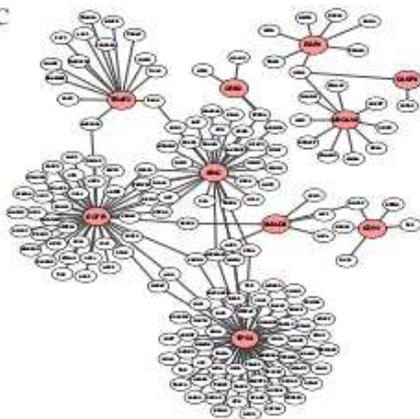
A



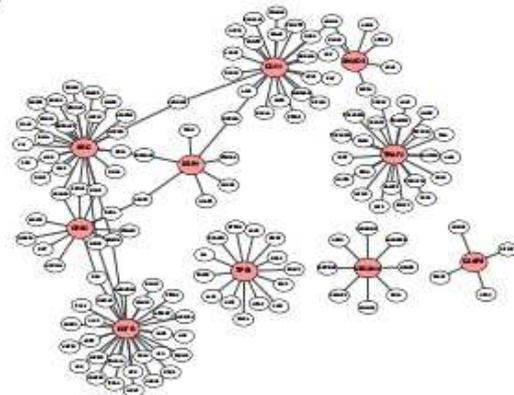
B



C



D



352

353 **Fig 5. The potential disease modules in BRCA Control network and Disease network. (A)**

354 **The network modules among the top- 10 hub gene in Control network. (B) The network**

355 **modules among the top- 10 hub gene in Disease network in sample BRCA_AAAK. (C) The**

356 network modules among the top- 10 hub gene in Disease network in sample BRCA_A0CZ. (D)
 357 The network modules among the top- 10 hub gene in Disease network in sample BRCA_A440.

358

359 The potential disease modules of each sample were enriched to the corresponding pathway
 360 in KEGG and enriched the disease genes in CGC database. The hypergeometric test was used
 361 to test the significant level of the enrichment analysis and the percentage of significant samples.
 362 The results showed that the top- 100, 50, 30, 20 and 10 hub genes of more than 95% samples
 363 in BRCA were significantly enriched to Pathways in cancer and Breast cancer pathway in
 364 KEGG (Figure 3D, Table 2). All potential disease modules (top- 100, 50, 30,20 and 10 hub
 365 genes) were significantly enriched to tumor genes in CGC database (Figure 3E, Table 3). There
 366 are similar results for the potential disease modules in LUAD, LUSC, and LIHC with BRCA
 367 (Figure S5). For example, there are more than 95% LIHC samples to be significantly enriched
 368 to Pathways in cancer and Hepatocellular carcinoma pathway in KEGG and all samples to be
 369 significantly enriched to existing tumor genes in CGC by top-100, 50, 30, 20 and 10 genes of
 370 DSSN (Figure S5C and F). For LUAD and LUSC, the potential disease modules of almost all
 371 tumor samples were significantly enriched corresponding Pathways in cancer, Non-small cell
 372 lung cancer pathway in KEGG and existing tumor genes in CGC (Figure S5A, B, D, and E).
 373 Compared with the SSN method²⁴, our method can obtain higher accuracy of significant
 374 samples in the enrichment analysis of disease pathways and cancer-related genes by the
 375 potential disease modules from DSSN.

376

377 **Table 2 The enrichment in KEGG pathway compared with our method and SSN method**

BRCA	top-10	top-20	top-30	top-50	top-100
Our method	95.64%	99.64%	99.99%	100%	100%
SSN method	0.00%	0.18%	10.00%	4.63%	6.17%
LUAD	top-10	top-20	top-30	top-50	top-100
Our method	96.81%	95.68%	99.81%	100%	100%

SSN method	0.00%	0.00%	0.38%	3.57%	15.57%
LUSC	top-10	top-20	top-30	top-50	top-100
Our method	98.41%	97.21%	99.20%	100%	100%
SSN method	0.20%	0.00%	0.80%	4.38%	12.25%
LIHC	top-10	top-20	top-30	top-50	top-100
Our method	100.00%	100%	100%	100%	100%
SSN method	0.00%	0.00%	0.54%	3.78%	7.55%

378

379 **Table 3 The enrichment in CGC database compared with our method and SSN method**

BRCA	top-10	top-20	top-30	top-50	top-100
Our method	99.82%	100.00%	100.00%	100.00%	100.00%
SSN method	23.96%	25.41%	34.03%	55.63%	80.03%
LUAD	top-10	top-20	top-30	top-50	top-100
Our method	99.62%	100.00%	100.00%	100.00%	100.00%
SSN method	16.89%	22.14%	24.39%	50.09%	76.55%
LUSC	top-10	top-20	top-30	top-50	top-100
Our method	100.00%	100.00%	100.00%	100.00%	100.00%
SSN method	7.37%	10.76%	12.95%	32.07%	69.52%
LIHC	top-10	top-20	top-30	top-50	top-100
Our method	98.92%	99.73%	100.00%	100.00%	100.00%
SSN method	10.78%	10.51%	15.36%	30.19%	60.91%

380

381

382 *Predicting individual driver mutation by DSSN*

383 Somatic mutation genes of a tumor sample can provide individual-specific information for this
384 sample²⁸ and can be used to verify the potential driver genes of the sample. There are 125
385 existing driver mutation genes to have been determined for cancer in reference²⁹. As we have
386 referred, a hub gene in SSDN is a crucial gene from normal to tumor state. If a hub gene of
387 SSDN was mutated, the gene may impact more genes than the non-hub gene and would be the
388 potential driver mutation gene for this sample. Based on such an assumption, DSSN was
389 involved in the network change between normal and tumor, and the hub genes in DSSN are
390 more likely to associate with disease. So, the probability/proportion, which a mutated gene is
391 an existing driver mutation gene, was respectively calculated for the top-100, 50, 30, 20 and
392 10 hub genes of each DSSN in each tumor and the average probability of each tumor was

393 shown in Figure 3F. The results showed that the probability/proportion was monotonically
394 increased from the top- 100 to 10 hub genes of DSSN (Figure 3F). As an example, if a gene in
395 the top 100 hub genes of DSSN of a BRCA sample was mutated, the probability of this gene
396 being a driver mutation gene is 63.04% (Figure 3F). If a gene was mutated in the top 10 hub
397 genes of DSSN of a BRCA sample, the probability would rise to 88.96% (Figure 3F). It means
398 if a hub gene of DSSN is mutated, the gene is a high probability to be a driver mutation gene
399 in BRCA. Similar results were shown in LUAD, LUSC, and LIHC (Figure 3F). Therefore, the
400 hub genes of DSSN are strongly related to the disease cause for one sample, and high-degree
401 genes are more likely to be carcinogenic factors.

402

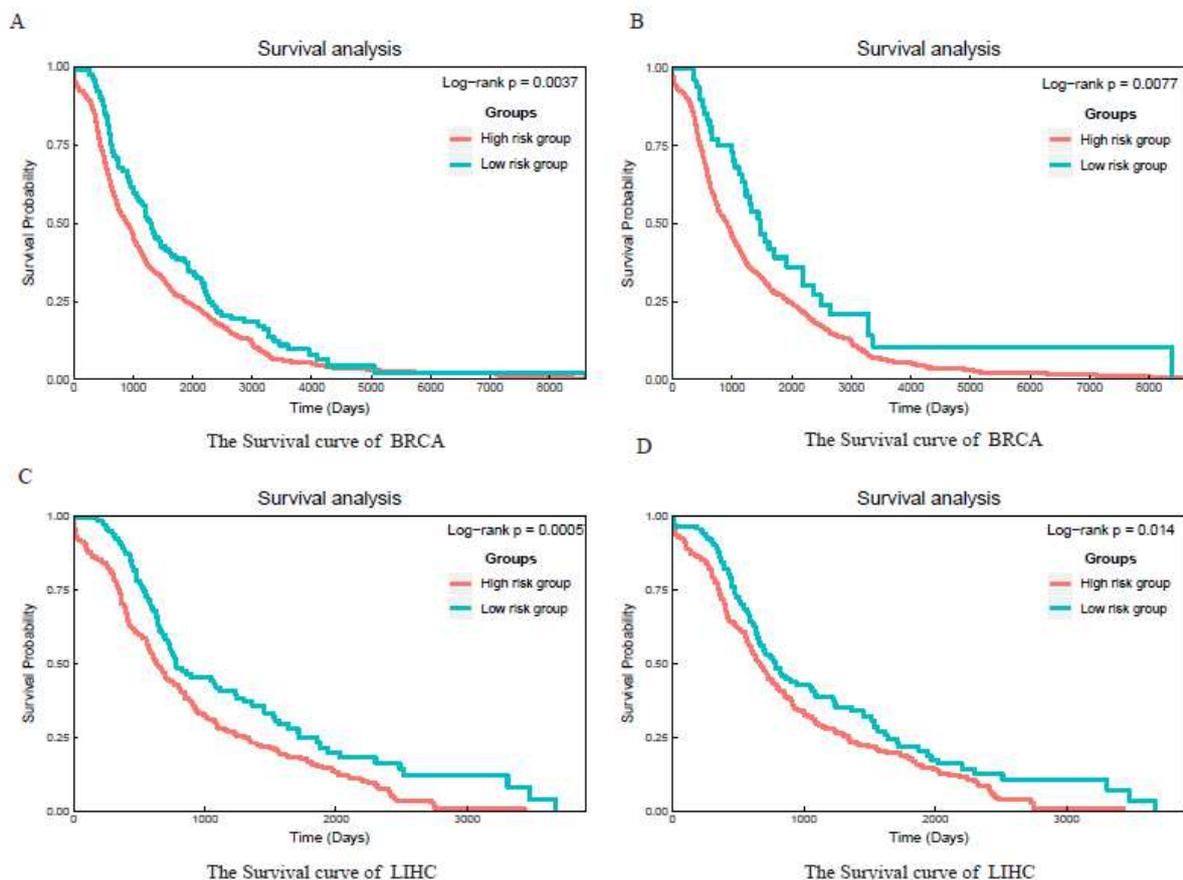
403 ***Prognosis analysis for tumor samples***

404 The clinic follow-up information was collected for each sample in TCGA, and the survival
405 times (unit is days) and vital status (alive or dead) were filtered out. The samples that missed
406 survival times or vital status were ignored. For BRCA and LIHC, the repetition hub genes were
407 identified based on the top 10 hub genes of each DSSN, and the most frequent 10 repetition
408 hub genes were used to survival analysis for tumor samples. We used the 10 repetition hub
409 genes to divide tumor samples into two groups, one included the samples that there were at
410 least 4 repetition hub genes to be in the top 10 hub genes of this sample; another included the
411 samples that had less than 4 repetition hub genes to be in the top 10 hub genes of this sample.
412 The repetition hub genes can also be identified based on the top 20 hub genes of each DSSN,
413 and the most frequent 20 repetition hub genes were used to survival analysis for tumor samples.
414 In the same way, the 20 repetition hub genes to divide tumor samples into two groups, one
415 included the samples that there were at least 9 repetition hub genes to be in the top 20 hub
416 genes of this sample; another included the samples that had less than 9 repetition hub genes to
417 be in the top 20 hub genes of this sample. A log-rank test was employed to test the significance

418 of the survival time.

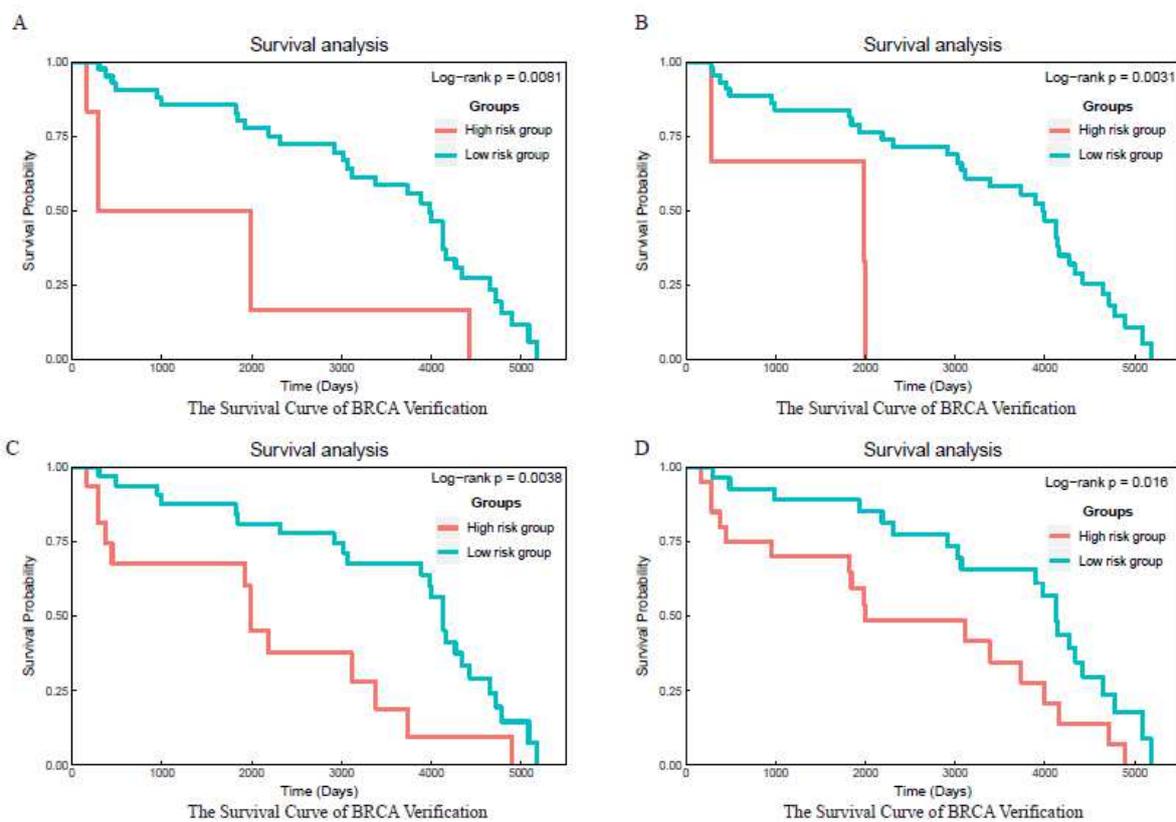
419 The most frequent 10 and 20 repetition hub genes can significantly distinguish the samples
420 with different survival time. For BRCA, the most frequent 10 repetition hub genes samples in
421 the two groups (high and low risk) can be significantly distinguished with p -value 0.0037 of
422 log-rank test (Figure 6A); the most frequent 20 repetition hub genes samples in the two groups
423 also can be significantly distinguished with p -value 0.0077 of log-rank test (Figure 6B). The p -
424 value of the most frequent 10 and 20 repetition hub genes of LIHC samples for survival analysis
425 are respectively 0.0005 and 0.014 (Figure 6C and 6D). Including LUSC samples, the repetition
426 hub genes can be used to prognosis analysis in BRCA, LUAD, and LIHC (Figure 6, Figure S6).
427 An independent dataset from ICGC database for BRCA was used to validate the effectiveness
428 of the DSSN in survival analysis, and the significant results were shown in Figure 7. The
429 repetition hub genes can be considered as the potential biomarkers for prognosis of the tumors.

430



431

432 **Fig 6. Survival curve for BRCA and LIHC.** (A) Survival curve for BRCA survival analysis
 433 when using the most frequent 10 repetition hub genes to divide tumor samples into two groups.
 434 (B) Survival curve for BRCA survival analysis when using the most frequent 20 repetition hub
 435 genes to divide tumor samples into two groups. (C) Survival curve for LIHC survival analysis
 436 when using the most frequent 10 repetition hub genes to divide tumor samples into two groups.
 437 (D) Survival curve for LIHC survival analysis when using the most frequent 20 repetition hub
 438 genes to divide tumor samples into two groups.



439
 440 **Fig 7. Examine the statistical significance of hub genes using BRCA test samples.** (A)
 441 Survival curve for BRCA verification survival analysis when using the most frequent 10
 442 repetition hub genes to divide tumor samples into two groups. (B) Survival curve for BRCA
 443 verification survival analysis when using the most frequent 20 repetition hub genes to divide
 444 tumor samples into two groups. (C) Survival curve for BRCA verification survival analysis
 445 when using the most frequent 30 repetition hub genes to divide tumor samples into two groups.

446 (D) Survival curve for BRCA verification survival analysis when using the most frequent 50
447 repetition hub genes to divide tumor samples into two groups.

448

449 **CONCLUSION**

450 In this work, we proposed the SSDN method and overcame the shortcomings of the previous
451 method. We compared the difference between the SSN in Disease and SSN in Control,
452 constructed the SSDN and DSSN. And we verified the consistency of *s-PCC* based on different
453 reference networks from both theoretical and practical. SSDN and DSSN can also be used to
454 select disease modules (hub genes) to evaluate individual features. The enrichment analysis of
455 KEGG pathway and CGC database indicated the disease modules play an important role in
456 cancer-related function by using the TCGA datasets for various cancer. The results of survival
457 analysis demonstrated these genes can be used independently as individual module biomarkers.
458 We expect SSDN and DSSN to generalize well to further studies, including application to
459 pathological diagnosis and drug therapy for patient-specific cancer care.

460

461 **DISCUSSION**

462 The SSDN method is an improvement of SSN method. The SSN method uses the reference
463 network, but it does not explain whether the choice of different reference networks will affect
464 the construction of Disease/Control network. We used the mathematical explanations to prove
465 that in two cases, one is the size of samples is sufficiently large, other is the samples come from
466 the same distribution. We have proved that the *PCC* in these two cases are consistent. And
467 verified that the different reference networks have no effect on the construction of the network
468 by using the real cancer data. Another improvement is the SSN method have constructed the
469 individual-specific network for disease sample and is called the Disease network, the
470 individual-specific network for control sample is called the Control network. We compared the

471 two networks, and found the difference between the Disease network and Control network.
472 Because if we only rely on the *PCC* to judge the significance of a gene pair/edge, there will
473 have a high false positive. Comparing the difference between the two networks, we can select
474 the gene pairs that only exist in the Disease network, or select the gene pairs that only exist in
475 the Control network. In this way, we can reduce the false positive gene pairs, find the real
476 difference between the two networks and remove the disturbance. Moreover, by using the real
477 cancer data, we verified the potential disease modules for individual sample in SSDN method
478 were enriched to the corresponding pathway in KEGG and enriched to the disease genes in
479 CGC database (Table 2, Table 3).

480 Several limitations of the current study should be considered, based on the limitations, we
481 construct two kinds of different networks, SSDN and DSSN. For constructing SSDN, a sample
482 pair are necessary for every patient. It means that there must be a tumor and control sample
483 from the same patient tissue, respectively, only then we can construct a Disease network and a
484 Control network for each patient. While if the normal samples cannot correspond to tumor
485 samples, the situation has changed slightly. We create a new network, common normal network,
486 for which the edge includes more than 1/3 Control networks' edges for a same cancer. The
487 common normal network was called the Control network, and the rest steps of constructing
488 DSSN are the same as SSDN. For an edge in the Disease network, if it is not in the Control
489 network, then the edge was kept in the DSSN. The DSSN is a disease-specific sample network
490 because the common normal network is constructed based on the same cancer, the subsequent
491 conclusions are also concluded on the same cancer. Therefore, SSDN is constructed by
492 identifying the different edges between the SSN in disease and SSN in control, DSSN is
493 constructed by identifying the different edges between the SSN in disease and common normal
494 network. Here, in order to construct SSDN, we used three gastric cancer datasets from the GEO
495 database because the samples with the same cancer obey the same distribution. These three

496 datasets have 25, 45, and 80 paired samples, which include tumor and control samples from
497 the same patient tissue. Since paired samples are two samples from the same patient, we can
498 construct a Disease network and a Control network for individual patient. However, by
499 analyzing four cancer datasets from TCGA database, due to the normal samples that cannot
500 correspond to tumor samples, we have to use DSSN to solve the problem. Although DSSN
501 cannot completely consistent with our previous methods, it indeed the best way to solve the
502 problem of missing data.

503 We used three gastric cancer datasets from the GEO database and these datasets are
504 GSE33335, GSE63089, and GSE27342. For convenience, we noted the reference network
505 from GSE33335 as *Ref1*, from GSE63089 as *Ref2* and from GSE27342 as *Ref3*. We used the
506 *Ref2* and *Ref3* as two reference networks. For example, we add a normal sample of GSE27342
507 in a reference network and obtain the *s-PCC* between each gene pair. If the relationship between
508 a pair of genes is significant, this pair of genes can be linked to an edge in the Control network.
509 In the same way, we add a disease sample of GSE27342 in reference networks and obtain a
510 Disease network. Hence, there is a Disease network and a Control network for a individual
511 sample. For the Disease network and Control network, there are six ways to obtain SSDN. The
512 first SSDN is constructed based on specific genes in Disease and Control networks, the second
513 SSDN is constructed based on common genes in Disease and Control network, and the rest are
514 constructed based on the specific genes in Control networks, the specific genes in Disease
515 networks, the genes only in Control networks, and the genes only in Disease networks. We
516 selected the most frequent 300, 100, 50, 30, 20 and 10 repetition hub genes, and calculated the
517 percentage of common hub genes based on *Ref1 Ref2, Ref3*. For dataset GSE27342 and the
518 most frequent 10 hub genes, the percentage of common hub genes in the first SSDN is 50%,
519 the percentage of common hub genes in the second SSDN is 54.77%. The results of other
520 percentage hub genes were shown in Table S1. According to our theory and the real data

521 validation, when the sample size is sufficiently large or the distribution is the same, the
522 consistency of s -PCC will be very high. While the sample size of the above three gastric cancer
523 datasets is just 25, 45, and 80, the sample size is too small to get a high consistency. The SSDN
524 established based on *Ref2* and *Ref3* with a relatively large sample size is already the best result
525 at present. Furthermore, we hope to develop a method to construct SSDN independent of the
526 sample number.

527 For the survival analysis, we also used other parameters (the number of top hub genes and
528 the number of repetitions hub genes) to test the robustness of the DSSN. For BRCA, the most
529 frequent 30 repetition hub genes samples in the two groups (high and low risk) can be
530 significantly distinguished with p -value 0.013 of the log-rank test (Figure S2A); the most
531 frequent 50 repetition hub genes samples in the two groups also can be significantly
532 distinguished with p -value 0.024 of the log-rank test (Figure S2A). The results showed that
533 different parameters can obtain similar results for prognosis (Figure S2).

534 Compared with the prior work¹⁵, the previous work of Liu, X et al. We took the normal
535 samples of the same cancer as a whole, and calculate the correlation by Liu's method, and took
536 the disease samples of the same cancer as a whole, calculate the correlation by Liu's method.
537 Then we calculate the CGC enrichment and KEGG enrichment of the whole cancer. The results
538 show that SSDN method is prior than the Liu's method. We have proved, although the SSDN
539 method is used for single sample difference analysis, it turns out if treat the cancer as a whole,
540 the enrichment analysis results are still good (Supplement Table S2, Table S3).

541 And we compared with the SSN method and our method. For BRCA, we used the SSDN
542 method, and selected the gene pairs that specific existing in the disease network, then identified
543 the top 6 genes with the highest degree in these gene pairs. For one sample, we have constructed
544 the SSN, and calculated the top 10 genes with the highest degree in SSN. If the top 10 genes
545 have less than 2 genes of the top 6 genes, the sample will be regarded as the normal sample.

546 Otherwise, the sample can be regarded as the disease sample. The percentage accuracy of the
547 SSDN classification is 88.56%. On the other hand, we used the SSN method to classify the
548 control sample and disease sample. We have selected the gene pairs that existing in the control
549 network, and calculated the top 6 genes with the highest degree in these gene pairs. Then we
550 have constructed the SSN in control and SSN in disease, and calculated the top 10 genes with
551 the highest degree in SSN. If the top 10 genes have less than 2 genes of the top 6 genes, the
552 sample will be regarded as the control sample, otherwise it's a disease sample. Finally, we have
553 selected the gene pair that existing in the disease network, and calculated the top 6 genes with
554 the highest degree in these gene pairs. Based on the top 6 genes, we classify again. As a result,
555 for BRCA, the accuracy of SSN in control classification is 54.40%, the accuracy of SSN in
556 disease classification is 85.27%. The accuracy of classification for the other three cancer are
557 also shown in the Supplement Table S4. The result shows that the SSDN method select is
558 indeed the specific gene in the disease network, and is better than SSN method. Because the
559 SSN method can only reflect the information of a single individual, and our method can
560 maximize the use of normal sample and disease sample.

561 As for survival analysis, we also compared with the SSN method. For BRCA and LIHC,
562 the repetition hub genes were identified based on the top 10 hub genes of each DSSN, and the
563 most frequent 10 repetition hub genes were used to survival analysis for tumor samples. We
564 calculated the top 6 genes with the highest degree in SSN in control, and used the 10 repetition
565 hub genes to divide tumor samples into two groups, one included the samples that there were
566 at least 2 repetition hub genes to be in the top 6 hub genes of this sample; another included the
567 samples that had less than 2 repetition hub genes to be in the top 10 hub genes of this sample.
568 In the same way, we calculate the top 6 genes with the highest degree in SSN in disease, and
569 repeat the same classification process above. The *p*-value in DSSN method is below
570 0.05(Figure 6), is better than SSN (Supplement Table S5). It further illustrates that

571 SSDN/DSSN method can find the difference between disease and control samples.

572

573 **ACKNOWLEDGE**

574 National Key R&D Program of China (No. 2017YFA0505500), National Natural Science
575 Foundation of China (NSFC) [61403363,11901272], Key Project of Natural Science of Anhui
576 Provincial Education Department (No. KJ2020A0018, KJ2016A002), Key project of teaching
577 and research of Anhui Finance and Economics University (No. acjyzd201606), and Shanghai
578 Municipal Science and Technology Major Project (No. 2017SHZDZX01).

579

580 **Ethics approval and consent to participate**

581 The authors declare that they have no competing interests.

582

583 **Consent for publication**

584 Not Applicable.

585

586 **REFERENCES**

- 587 1. Network, T. C. G. A. R., Corrigendum: Comprehensive genomic characterization defines
588 human glioblastoma genes and core pathways. *Nature* **2013**, 494, 506.
- 589 2. Uhm, J., Comprehensive genomic characterization defines human glioblastoma genes and
590 core pathways TCGA Research Network Nature 2008 455 1061 8 2671642 18772890. 2009; Vol.
591 2009, p 117-118.
- 592 3. International Cancer Genome, C.; Hudson, T. J.; Anderson, W.; Artez, A.; Barker,
593 A. D.; Bell, C., et al., International network of cancer genome projects. *Nature* **2010**, 464
594 (7291), 993-998.
- 595 4. Hudson, T. J.; Anderson, W.; Aretz, A.; Barker, A. D.; Bell, C.; Bernabé, R. R.,
596 et al., International network of cancer genome projects. *Nature* **2010**, 464, 993-998.
- 597 5. Zhang, J.; Zhang, S., Discovery of cancer common and specific driver gene sets. *Nucleic*
598 *Acids Research* **2017**, 45 (10), e86-e86.
- 599 6. Zhang, S.; Liu, C.-C.; Li, W.; Shen, H.; Laird, P. W.; Zhou, X. J., Discovery of
600 multi-dimensional modules by integrative analysis of cancer genomic data. *Nucleic acids*
601 *research* **2012**, 40 (19), 9379-9391.
- 602 7. Chin, L.; Andersen, J. N.; Futreal, P. A., Cancer genomics: from discovery science to
603 personalized medicine. *Nature Medicine* **2011**, 17, 297.
- 604 8. Schilsky, R. L., Personalized medicine in oncology: the future is now. *Nature Reviews*
605 *Drug Discovery* **2010**, 9, 363.

606 9. Kiesel, A.; Chia, B. K. H.; Bertrand, D.; Chng, K. R.; Nagarajan, N.; Hillmer, A.,
607 et al., Patient-specific driver gene prediction and risk assessment through integrated
608 network analysis of cancer omics profiles. *Nucleic Acids Research* **2015**, 43 (7), e44-e44.
609 10. Chen, L.; Wang, R. S.; Zhang, X. S., *Biomolecular Networks: Methods and Applications*
610 in Systems Biology. 2009.
611 11. Hood, L.; Flores, M., A personal view on systems medicine and the emergence of proactive
612 P4 medicine: predictive, preventive, personalized and participatory. *New Biotechnology* **2012**,
613 29 (6), 613-624.
614 12. Chen, L.; Wang, R.; Li, C.; Aihara, K., *Modeling Biomolecular Networks in Cells:*
615 *Structures and Dynamics*. Springer-Verlag. London: 2010.
616 13. Liu, X.; Chang, X.; Liu, R.; Yu, X.; Chen, L.; Aihara, K., Quantifying critical
617 states of complex diseases using single-sample dynamic network biomarkers. *Plos*
618 *Computational Biology* **2017**, 13 (7), e1005633.
619 14. Liu, X.; Liu, Z. P.; Zhao, X. M.; Chen, L., Identifying disease genes and module
620 biomarkers by differential interactions. *Journal of the American Medical Informatics*
621 *Association Jamia* **2016**, 19 (2), 241-248.
622 15. Liu, X.; Chang, X., Identifying module biomarkers from gastric cancer by differential
623 correlation network. *Oncotargets & Therapy* **2016**, 9, 5701-5711.
624 16. Gov, E.; Arga, K. Y., Differential co-expression analysis reveals a novel prognostic
625 gene module in ovarian cancer. *Sci Rep-Uk* **2017**, 7.
626 17. Langfelder, P.; Horvath, S., WGCNA: an R package for weighted correlation network
627 analysis. *Bmc Bioinformatics* **2008**, 9.
628 18. Gill, R.; Datta, S.; Datta, S., A statistical framework for differential network
629 analysis from microarray data. *Bmc Bioinformatics* **2010**, 11.
630 19. Hu, B.; Chang, X.; Liu, X., Predicting Functional Modules of Liver Cancer Based on
631 Differential Network Analysis. *Interdisciplinary Sciences Computational Life Sciences* **2019**,
632 (8), 1-9.
633 20. Zhang, W. W.; Zeng, T.; Liu, X. P.; Chen, L. N., Diagnosing phenotypes of single-
634 sample individuals by edge biomarkers. *J Mol Cell Biol* **2015**, 7 (3), 231-241.
635 21. Zhang, W. W.; Zeng, T.; Chen, L. N., Edge Marker: Identifying differentially correlated
636 molecule pairs as edge-biomarkers. *J Theor Biol* **2014**, 362, 35-43.
637 22. Rui, L.; Xiangdong, W.; Kazuyuki, A.; Luonan, C., Early diagnosis of complex diseases
638 by molecular biomarkers, network biomarkers, and dynamical network biomarkers. *Medicinal*
639 *Research Reviews* **2014**, 34 (3), 455-478.
640 23. Ideker, T.; Krogan, N. J., Differential network biology. *Mol Syst Biol* **2012**, 8.
641 24. Liu, X.; Wang, Y.; Ji, H.; Aihara, K.; Chen, L., Personalized characterization of
642 diseases using sample-specific networks. *Nucleic Acids Research* **2016**, 44 (22), e164-e164.
643 25. P Andrew, F.; Lachlan, C.; Mhairi, M.; Thomas, D.; Timothy, H.; Richard, W., et
644 al., A census of human cancer genes. *Nature Reviews Cancer* **2004**, 4 (3), 177-183.
645 26. Sherman, B. T.; Huang, D. W.; Lempicki, R. A., Bioinformatics enrichment tools: paths
646 toward the comprehensive functional analysis of large gene lists. *Nucleic Acids Research*
647 **2008**, 37 (1), 1-13.
648 27. R Team, C., R: A language and environment for statistical computing. [http://www.R-](http://www.R-project.org)
649 [project.org](http://www.R-project.org) **2013**.
650 28. Blau, C. A.; Liakopoulou, E., Can we deconstruct cancer, one patient at a time? *Trends*
651 *in Genetics* **2013**, 29 (1), 6-10.
652 29. Bert, V.; Nickolas, P.; Velculescu, V. E.; Shibin, Z.; Diaz, L. A.; Kinzler, K. W.,
653 *Cancer genome landscapes*. *Science* **2013**, 339, 1546-1558.

654

655 **Availability of data and materials**

656 **LIHC cancer gene, the raw data of LIHC cancer gene.**

657 LIHC normal gene, the raw data of LIHC normal gene.
658 LUSC cancer gene, the raw data of LUSC cancer gene.
659 LUSC normal gene, the raw data of LUSC normal gene.
660 LUAD cancer gene, the raw data of LUAD cancer gene.
661 LUAD normal gene, the raw data of LUAD normal gene.
662 BRCA cancer gene, the raw data of BRCA cancer gene.
663 BRCA normal gene, the raw data of BRCA normal gene.
664 All the datasets generated and analysed during the current study are available in the [TCGA]
665 repository, [[https://www.cancer.gov/about-nci/organization/ccg/research/structural-](https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga)
666 [genomics/tcga](https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga)].
667 Gastric Cancer Sample, GSE27342
668 [<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE27342>]
669 GSE63089 [<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE63089>]
670 GSE33335 [<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE33335>]

671 **Authors' Contributions**

672 Yu Zhang and Xiaoping Liu proposed the sample-specific differential network, and were the
673 major contributor in writing the manuscript. Yu Zhang and Xiao Chang prepared Figure 1, and
674 construct the network. Jie Xia prepared Figure 2, Shaoyan Sun prepared Figure 3-7, and
675 examine the cancer and normal database. All authors reviewed the manuscript. All authors read
676 and approved the final manuscript.

677 **Supplementary Figure Legends**

678 **Figure S1. The network modules with the potential disease modules in BRCA Control**
679 **network and Disease network.** (A) The network modules among the top- 20 hub gene in
680 Control network. (B) The network modules among the top- 20 hub gene in Disease network in
681 sample BRCA_A0T6. (C) The network modules among the top- 20 hub gene in Disease

682 network in sample BRCA_A4RY. (D) The network modules among the top- 20 hub gene in
683 Disease network in sample BRCA_A1IX.

684 **Figure S2. The network modules with the potential disease modules in LIHC reference**
685 **network.** (A) The network modules among the top- 10 hub genes. (B) The network modules
686 among the top- 20 hub genes.

687 **Figure S3. The network modules with the potential disease modules in LIHC Control**
688 **network and Disease network.** (A) The network modules among the top- 10 hub gene in
689 Control network. (B) The network modules among the top- 10 hub gene in Disease network in
690 sample LIHC_A9H1. (C) The network modules among the top- 10 hub gene in Disease
691 network in sample LIHC _A69I. (D) The network modules among the top- 10 hub gene in
692 Disease network in sample LIHC _AAC9.

693 **Figure S4. The network modules with the potential disease modules in LIHC Control**
694 **network and Disease network.** (A) The network modules among the top- 20 hub gene in
695 Control network. (B) The network modules among the top- 20 hub gene in Disease network in
696 sample LIHC_A110. (C) The network modules among the top- 20 hub gene in Disease network
697 in sample LIHC _A520. (D) The network modules among the top- 20 hub gene in Disease
698 network in sample LIHC _AA0V.

699 **Figure S5. The enrichment in KEGG pathway and CGC database compared with our**
700 **method and SSN method.** (A)The proportion of significant samples in the enrichment analysis
701 of top- 100, 50, 30, 20 and 10 highest degree genes for LUAD DSSN in the KEGG pathway
702 and compare with the previous method. The x-axis is the hub genes and the y-axis is the
703 percentage of significant samples in KEGG enrichment analysis. (B) The proportion of
704 significant samples in the enrichment analysis of top- 100, 50, 30, 20 and 10 highest degree
705 genes for LUSC DSSN in the KEGG pathway and compare with the previous method. The x-
706 axis is the hub genes and the y-axis is the percentage of significant samples in KEGG

707 enrichment analysis. (C) The proportion of significant samples in the enrichment analysis of
708 top- 100, 50, 30, 20 and 10 highest degree genes for LIHC DSSN in the KEGG pathway and
709 compare with the previous method. The x-axis is the hub genes and the y-axis is the percentage
710 of significant samples in KEGG enrichment analysis. (D) The proportion of significant samples
711 in the enrichment analysis of top- 100, 50, 30, 20 and 10 highest degree genes for LUAD DSSN
712 in the CGC database and compare with the previous method. The x-axis is the hub genes of
713 cancer and the y-axis is the percentage of significant samples in CGC enrichment analysis. (E)
714 The proportion of significant samples in the enrichment analysis of top- 100, 50, 30, 20 and 10
715 highest degree genes for LUSC DSSN in the CGC database and compare with the previous
716 method. The x-axis is the hub genes of cancer and the y-axis is the percentage of significant
717 samples in CGC enrichment analysis. (F) The proportion of significant samples in the
718 enrichment analysis of top- 100, 50, 30, 20 and 10 highest degree genes for LIHC DSSN in the
719 CGC database and compare with the previous method. The x-axis is the hub genes of cancer
720 and the y-axis is the percentage of significant samples in CGC enrichment analysis.

721 **Figure S6. Survival curve for BRCA and LIHC.** (A) Survival curve for BRCA survival
722 analysis when using the most frequent 30 repetition hub genes to divide tumor samples into
723 two groups. (B) Survival curve for BRCA survival analysis when using the most frequent 50
724 repetition hub genes to divide tumor samples into two groups. (C) Survival curve for LIHC
725 survival analysis when using the most frequent 30 repetition hub genes to divide tumor samples
726 into two groups. (D) Survival curve for LIHC survival analysis when using the most frequent
727 50 repetition hub genes to divide tumor samples into two groups. (E) Survival curve for LUAD
728 survival analysis when using the most frequent 30 repetition hub genes to divide tumor samples
729 into two groups. (F) Survival curve for LUSC survival analysis when using the most frequent
730 20 repetition hub genes to divide tumor samples into two groups.

731

Supplementary Files

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

- [BRCAcancergene.txt](#)
- [BRCAnormalgene.txt](#)
- [LIHCcancergene.txt](#)
- [LIHCnormalgene.txt](#)
- [LUSCnormalgene.txt](#)
- [SupplementFigureS1.pdf](#)
- [SupplementFigureS2.pdf](#)
- [SupplementFigureS3.pdf](#)
- [SupplementFigureS4.pdf](#)
- [SupplementFigureS5.pdf](#)
- [SupplementFigureS6.pdf](#)
- [SupplementNote.docx](#)
- [SupplementTableS1.docx](#)
- [SupplementTableS2.docx](#)
- [SupplementTableS3.docx](#)
- [SupplementTableS4.docx](#)
- [SupplementTableS5.docx](#)
- [LUSCcancergene.txt](#)
- [LUADnormalgene.txt](#)
- [LUADcancergene.txt](#)