

Integrating Protein Interaction Networks and Somatic Mutation Data to Detect Driver Modules in Pan-Cancer

Hao Wu (✉ haowu@sdu.edu.cn)

Shandong University <https://orcid.org/0000-0003-2340-9258>

Zhong-Li Chen

Northwest A&F University: Northwest Agriculture and Forestry University

Ying-Fu Wu

Northwest A&F University: Northwest Agriculture and Forestry University

Hong-Ming Zhang

Northwest A&F University: Northwest Agriculture and Forestry University

Quan-Zhong Liu

Northwest A&F University: Northwest Agriculture and Forestry University

Research

Keywords: driver modules, node similarity, random walk with restart, complex networks

Posted Date: May 5th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-482282/v1>

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

[Read Full License](#)

Version of Record: A version of this preprint was published at Interdisciplinary Sciences: Computational Life Sciences on September 7th, 2021. See the published version at <https://doi.org/10.1007/s12539-021-00475-y>.

RESEARCH

Integrating Protein Interaction Networks and Somatic Mutation Data to Detect Driver Modules in Pan-cancer

Hao Wu^{2*}, Zhong-Li Chen^{1,3†}, Ying-Fu Wu¹, Hong-Ming Zhang^{1*} and Quan-Zhong Liu¹

Sample of title note

*Correspondence:

haowu@sdu.edu.cn;

zhm@nwsuaf.edu.cn

¹College of Information

Engineering, Northwest A&F University, Yangling, 712100, China

Full list of author information is available at the end of the article

[†]Zhong-Li Chen and Hao Wu contributed equally to this work.

Abstract

Background: With the constant update of large-scale sequencing data and the continuous improvement of cancer genomics data such as the cancer genome atlas ICGC and TCGA, it gains increasing importance how to detect the functional high-frequency mutation gene set in cells that causes cancer within the field of medicine.

Methods: In this study, to solve the issue of mutated gene heterogeneity and improve the accuracy of driver modules, we propose a new recognition method of driver modules, named ECSWalk, based on the human protein interaction networks and pan-cancer somatic mutation data. This study firstly utilizes high mutual exclusivity and high coverage between mutation genes and topological structure similarity of the nodes in complex networks to calculate interaction weights between genes. Secondly, the method of random walk with restart is utilized to construct a weighted directed network, and the strong connectivity principle of the directed graph is utilized to create the initial candidate modules with a certain number of genes. Finally, the large modules in the candidate modules are reasonably split using the way of the induced subgraph, and the small modules are expanded using a greedy strategy to obtain the optimal driver modules.

Results: This method is applied to the analysis of TCGA pan-cancer data, and the experimental results show that ECSWalk can detect driver modules more effectively and accurately, and can identify new candidate gene sets with higher biological relevance and statistical significance than MEXCOWalk and HotNet2.

Conclusions: ECSWalk is of theoretical guidance and practical value for cancer diagnosis, treatment and drug targets.

Keywords: driver modules; node similarity; random walk with restart; complex networks

1 Introduction

- 2 In recent years, with the continuous advancement of cancer research at the bio-
- 3 molecular level, targeted therapy for cancer-causing genes has become a hot research

4 area in cancer [1]. For the driver genes that have functional mutations in cancer,
5 specific drugs can be used to control their expression and transcription levels, and
6 effectively control the development and deterioration of cancer. However, we cannot
7 effectively identify all driver mutations in cancer by statistical analysis of a single
8 mutation gene[2]. Therefore, compared with the screening of a single driver muta-
9 tion gene, detection of the mutated driver gene sets in cancer is of high biological
10 relevance and statistical significance [3, 4]. In particular, they can not only explore
11 the pathogenesis of cancer in depth, but also provide drug targets for the clinical
12 treatment of cancer by inferring the interaction of upstream and downstream genes
13 in the driver modules, and they provide reliable theoretical basis and data support
14 for precision medicine and personalized medicine [5].

15 Previous studies have found that screening for high-frequency mutations is con-
16 ducive to identifying the group of mutated driver genes in cancer [6, 7, 8], such
17 as EGFR, TP53, PIK3CA and other high-frequency mutation genes, which are
18 screened as driver factors during tumorigenesis. However, it normally involves a
19 large amount of time and effort to identify high-frequency mutation gene sites in
20 a large number of tumor samples for biological experimental verification, which is
21 difficult to achieve under the current technical level, experimental conditions, and
22 research capabilities. Also, the current methods tend to ignore the problem of mu-
23 tation heterogeneity in complex diseases [9, 10]. To effectively solve this problem,
24 most of the previous studies utilized the high mutual exclusivity and high coverage
25 widely existed in the genome map to explore the driver pathways that lead to the
26 occurrence of cancer[6, 7, 11, 12].

27 Based on the characteristics of high mutual exclusivity and high coverage, Vandin
28 et al.[6] proposed the Dendrix algorithm to identify driver pathways from somatic
29 mutation data. This method introduces a penalty overlap and a reward coverage
30 mechanism to solve the maximum weight sub-matrix problem of driver pathway
31 identification based on gene mutation data. The Dendrix algorithm not only im-
32 proves coverage but also guarantees mutual exclusivity among the genes in a driver
33 pathway. However, this iterative search method is prone to produce local optimal
34 solutions, and it requires specifying the number of genes in a driver pathway in
35 advance. Leiserson et al.[7]proposed a Multi-Dendrix algorithm based on the Den-
36 drix algorithm, which detects multiple driver pathways at the same time. The al-
37 gorithm utilizes linear regression to simultaneously detect multiple gene sets that
38 meet high coverage and high mutual exclusivity, and can generate a globally opti-
39 mal solution. However, the Multi-Dendrix algorithm needs to pre-specify both the
40 maximum number of genes in a driver pathway and the number of driver pathways.
41 Therefore, the two algorithms do not have good universality and robustness, and
42 have certain limitations when applied to different data sets.

43 Pan-cancer data analysis provides new ideas and methods for the study of clinical
44 diagnosis and treatment across cancer types [13, 14, 15, 16]. As one application,
45 Leiserson et al.[14] proposed the HotNet2 algorithm by integrating differential ex-
46 pression genes, significant mutation genes and protein interaction networks to detect
47 combinations of rare somatic mutations. Using the principle of thermal diffusion and
48 the random walk with restart model, the edge weight that becomes stable after the
49 random walk is restarted as the weight of the directed edge. By removing the di-

50 rected edge with a small weight, the method detects strongly connected components
51 as driver modules in the directed graph. Although the algorithm reduces the output
52 of false positive results and improves the accuracy of the prediction results, the con-
53 version probability just considers the degree of the vertex during the random walk
54 process, but ignores the mutual exclusivity among genes. Therefore, the problem
55 of gene mutation heterogeneity cannot be effectively solved in the pan-cancer data
56 with large differences.

57 Based on the HotNet2 algorithm, Rafsan et al. [16] proposed a MEXCOWalk al-
58 gorithm to mine driver modules by using split and expansion techniques. The algo-
59 rithm utilizes mutual exclusivity between genes and coverage scores of gene set to
60 reflect the edge weight of the network. Then the random walk with restart strat-
61 egy is used to construct a weighted directed network, and the split and expansion
62 techniques are utilized to mine driver modules. Although the set of driver genes
63 identified by this algorithm has high mutual exclusivity and high coverage, the al-
64 gorithm ignores the relationships of mutual exclusivity between expanded leaf nodes
65 and seed modules in the expansion stage of small modules. Therefore, this algorithm
66 reduces the accuracy of driver module identification to a certain extent.

67 As mentioned above, although the previous methods can detect the gene sets with
68 high mutual exclusivity and high coverage, they just focus on the mutual exclusivity
69 and coverage between genes, instead of the topological structure of complex net-
70 works. To effectively solve the problem of mutated gene heterogeneity and improve
71 the accuracy of driver modules, this paper proposes a driver module detection algo-
72 rithm (ECSWalk) based on gene mutation and human protein interaction network.

73 The algorithm optimizes the definition in terms of the following three characteris-
74 tics of high mutual exclusivity, high coverage among genes and high similarity of
75 topological structure. Firstly, the complex network topology analysis method is used
76 in the human protein interaction network data to calculate the topology similarity
77 between network nodes, and then it is combined with the two characteristics of high
78 coverage and high mutual exclusivity of the mutated genes to obtain the weight of
79 the vertices and edges in the human protein network. Among them, the weights of
80 vertices in the human protein network are obtained according to the coverage of
81 mutated gene; the random walk with restart strategy is utilized to calculate the
82 weights of edges in the network by the three characteristics, namely the coverage,
83 the mutual exclusivity, and the similarity of the topological structure between the
84 nodes. Secondly, based on the weighted network constructed in the previous step,
85 the large modules are split into several candidate gene sets using the method of the
86 induced subgraph. And the greedy strategy is utilized to add the nodes in the leaf
87 module to the seed module to achieve the optimal gene sets. These mutated gene
88 sets with high mutual exclusivity, high coverage and high similarity of the topo-
89 logical structure are likely to work as driver modules in cancer[17, 18]. This study
90 not only applies the analysis method of complex network topology to the biological
91 network, but also improves the method of determining module size based on split
92 and expansion. The study clarifies the interactions between genes in the detected
93 driver modules, which promote the study of cancer pathogenesis and drug targets.

94 **Methods**

95 **Mutual exclusivity and coverage**

96 To accurately detect and classify a large number of genes in the cancer genome map,
 97 and reduce the impact from error in the actual biological sequencing experiment
 98 process, this study introduces the definition of mutual exclusivity and coverage[16].

99 Let $G(V, E)$ denote the protein interaction network (PPI), each vertex $u_i \in V$
 100 corresponds to a protein in PPI network, and each protein u_i corresponds to a
 101 mutated gene g_i . The undirected edge $(u_i, u_j) \in E$ in PPI network corresponds to
 102 the interaction between gene pair (g_i, g_j) . Therefore, the node g_i represents both
 103 a gene and the corresponding protein in G . The sample with gene g_i mutated is
 104 represented by S_i , and $M \subseteq V$ is a subset of genes. For any pair of genes $g_i, g_j \in M$,
 105 $g_i \neq g_j$, if $S_i \cap S_j = \emptyset$, the genes in M are mutually exclusive.

The mutual exclusivity of gene subset M is represented as:

$$ED(M) = \frac{|\cup_{g_i \in M} S_i|}{\sum_{g_i \in M} |S_i|} \quad (1)$$

106 if $ED(M) = 1$, then the genes in the subset M are mutually exclusive, that is, at
 107 most one gene within the subset M is mutated in each sample.

The coverage of gene subset M is represented as:

$$CD(M) = \frac{|\cup_{g_i \in M} S_i|}{|\cup_{g_i \in V} S_i|} \quad (2)$$

108 If $CD(M) = 1$, then the gene subset M completely covers all patients, that is, at
 109 least one gene within the subset M is mutated in each sample.

110 Node similarity

111 The abnormal local area network composed of nodes and their neighbor nodes
 112 may affect the performance of the entire network in complex networks. Therefore,
 113 to measure the topological relationship of each mutated gene in the same driver
 114 module, we select the local area network with node g_i as the center and its direct
 115 neighbor node as the radius. Combined with the characteristics of the local area
 116 network structure, we propose node similarity based on Jensen-Shannon (JS) di-
 117 vergence and analyze the topological structure of the protein interaction network
 118 [19, 20, 21]. The similarity is mainly defined by the discrete probability set, and the
 119 index construction process mainly includes two steps, the first step is to construct
 120 the probability set[22], and the second step is to define the node similarity according
 121 to the JS divergence.

Construction of the probability set. Let d_i be the degree of the i th gene node, and
 d_{max} be the maximum node degree in the local area network. Suppose there are N
 nodes in the probability set of one node, $N = d_{max} + 1$. The sum of node degrees
 D_{g_i} of gene node g_i in its local area network is expressed as follows:

$$D_{g_i} = \sum_{j=1}^n d(j) \quad (3)$$

122 Where N is the number of genes in the local area network, and $d(j)$ is the degree
 123 of the j th gene in the local area network of gene g_i .

In the local area network, the discrete probability of gene g_i is expressed as follows:

$$p(i) = \frac{d_i}{D_{g_i}} \quad (4)$$

124 Standardize the discrete probabilities of gene g_i , and sort the discrete probabilities
 125 of genes in the local area network of gene g_i from large to small, and obtain the set
 126 of discrete probabilities as $P(i)$.

$$P(i) = (p_i(1), p_i(2), \dots, p_i(n), \dots, p_i(N)) \quad (5)$$

127 Where $p_i(n)$ represents the discrete probability value of the n th gene in the local
 128 area network with N gene nodes ($n \leq N$). Therefore, in the discrete probability set
 129 $P(i)$, the N elements in set $P(i)$ are the discrete probability values of each node in
 130 the local area network.

Suppose that two adjacent genes g_i and g_j in the constructed gene network corre-
 spond to two different probability sets $P(i)$ and $P(j)$ with the same number of genes
 in the local area network. According to the set of discrete probability constructed
 in formula (5), the Kullback-Leibler (KL) divergence value between genes g_i and g_j
 is expressed as:

$$D_{KL}(P(i)||P(j)) = \sum_{k=1}^N p_i(k) \log \frac{p_i(k)}{p_j(k)} \quad (6)$$

Due to the asymmetry of KL divergence, different results can be obtained by exchanging the positions of $P(i)$ and $P(j)$. To solve this asymmetry problem, the same result can be obtained by exchanging the positions of $P(i)$ and $P(j)$, The JS divergence value is obtained based on the KL divergence value, which is expressed as:

$$D_{JS}(P(i)||P(j)) = \frac{1}{2}KL(P(i)||\frac{P(i)+P(j)}{2}) + \frac{1}{2}KL(P(j)||\frac{P(i)+P(j)}{2}) \quad (7)$$

Therefore, we use an analysis of gene pair similarity based on the network topology, which is defined as follows:

$$SIM(g_i, g_j) = 1 - D_{JS}(P(i)||P(j)) \quad (8)$$

131 Obviously, the larger the $SIM(g_i, g_j)$ value, the more similar the two gene nodes
 132 in the network. It can be seen from formula (8) that the value range of $SIM(g_i, g_j)$
 133 is $[0, 1]$, and $SIM(g_i, g_j) = 1$ indicates that the two gene nodes in the network have
 134 the same topological structure.

135 Construction of edge-weighted networks

136 Given a PPI network $G(V, E)$, where node set $V = (u_1, u_2, u_3, \dots, u_n)$ represents the
 137 set of mutated genes corresponding to the PPI network, edge set $E = \{e = (u_i, u_j)\}$
 138 represents the set of protein interaction relationships.

139 Construction of a weighted undirected graph G_ω . For each vertex $g_i \in V$, the
 140 coverage of vertex g_i represents the weight of vertex g_i , that is, $\omega(g_i) = CD(g_i)$.

141 Obviously, the more the mutated samples in gene g_i , the greater the weight of the
 142 vertex g_i .

Taking into account the chance of increasing the coexistence of a gene and its surrounding genes, this study defines the set of node g_i and its direct neighbor nodes as the local area network $Ne(g_i)$ as follows:

$$Ne(g_i) = \{g_i\} \cup (\cup_{\forall (g_i, g_j) \in E} g_j) \quad (9)$$

To balance the mutual exclusivity between genes and the opportunities for the coexistence between a gene and its surrounding genes, this study utilizes the average value of $ED(Ne(g_i))$ and $ED(Ne(g_j))$ as the mutual exclusivity $ED(g_i, g_j)$ of gene pairs in the network, as shown below.

$$ED(g_i, g_j) = \frac{ED(Ne(g_i)) + ED(Ne(g_j))}{2} \quad (10)$$

To reduce the chance of a single gene with large coverage added to the edge weight, the product of the coverage of two genes is used to represent the coverage $CD(g_i, g_j)$ between gene pairs, as shown below.

$$CD(g_i, g_j) = CD(\{g_i\}) \times CD(\{g_j\}) \quad (11)$$

The study integrates the three characteristics of mutual exclusivity, coverage, and similarity among gene pairs to calculate the edge weight of the weighted undirected

graph as follows:

$$\omega(g_i, g_j) = \begin{cases} SIM(g_i, g_j) \neq 0 & \\ \frac{2 \times SIM(g_i, g_j)}{\frac{1}{ED(g_i, g_j)} + \frac{1}{CD(g_i, g_j)}} & ED(g_i, g_j) \neq 0 \\ CD(g_i, g_j) \neq 0 & \\ 0 & otherwise \end{cases} \quad (12)$$

The principle of thermal diffusion is utilized to construct a weighted directed graph by performing a random walk with restart on G_ω [14]. The random walk with restart means that the source node gene g_i transfers to its neighboring nodes with a certain probability, and they utilize the restart probability to transfer to the source node again. This process is repeated until it reaches a stable state. It is expressed as follows:

$$F_{t+1} = (1 - \beta)PF_t + \beta F_0 \quad (13)$$

Where F_0 is the initial state of the source node gene g_i , and $F_0 = CD(g_i) \cdot F_t$ is the probability distribution at time t ; β is the restart probability of its neighbor nodes transferring to the source node, and $0 \leq \beta \leq 1$, which is used to control the heat of the source node diffusion to the rest nodes of the network. It is necessary to choose a suitable β that all source nodes retain most of the heat in their direct neighbor nodes [14]. According to [14, 16], the value of β in the study is set to be 0.4; E represents the transition probability matrix of the restart random walk

process, which is positively correlated with the edge weight, as shown below:

$$P(g_i, g_j) = \begin{cases} \frac{\omega(g_i, g_j)}{\sum_k \omega(g_i, g_j)} & (g_i, g_j) \in E \\ 0 & otherwise \end{cases} \quad (14)$$

143 Where $\sum_k \omega(g_i, g_j)$ represents the sum of the edge weights between the source node
144 g_i and its direct neighbor nodes.

As the value of t increases, F_{t+1} gradually converges, and then the random walk restarts until it reaches a stable state [23]. The edge weight value F is calculated according to the following formula [14]:

$$F = \beta(I - (1 - \beta)(P(g_i, g_j)))^{-1} F_0 \quad (15)$$

145 Where I is the identity matrix. Restart the random walk to create a directed edge
146 with weight F for each pair of gene pair g_i and g_j ($i \neq j$)[14], and finally realize
147 the construction of a weighted directed graph G_d . The algorithm is described as follows.

Algorithm 1 Construction of the weighted directed network

Input: a PPI network $G(V, E)$, gene mutation sample set S

Output: $G(V, E, F)$

- 1: *Initialization* : $j = |E|, i = 1, \beta = 0.4$
 - 2: **for** i to j **do**
 - 3: **compute** $ED(g_i, g_j), CD(g_i, g_j), SIM(g_i, g_j), \omega(g_i, g_j)$
 - 4: **if** $ED(g_i, g_j) \neq 0$ and $CD(g_i, g_j) \neq 0$ and $SIM(g_i, g_j) \neq 0$ **then**
 - 5: $genenetwork[g_i][g_j] = \omega(g_i, g_j)$
 - 6: $i = i + 1$
 - 7: **if** $(g_i, g_j) \in E$ **then**
 - 8: $F = \beta(I - (1 - \beta)(\frac{\omega(g_i, g_j)}{\sum_k \omega(g_i, g_j)}))^{-1} F_0$
-

149 Driver modules detection

150 In this study, the strongly connected component (SCC) division method of the
 151 directed graph is utilized to generate the driver modules [14]. The process is divided
 152 into three steps.

153 The first step is to create a set of initial candidate modules. The SCC is firstly
 154 employed as an initial set of candidate modules. The minimum weight edge in G_d is
 155 iteratively deleted until strongly connected subgraphs are generated from G_d , and
 156 then the strongly connected subgraph is added to the initial module set P . Finally,
 157 all the modules in P whose gene number is less than min_module_size are removed.
 158 The above process is carried out iteratively until the number of genes in P decreases
 159 to $total_genes$. We finally obtain the initial module set $P = (M_1, M_2, \dots, M_r)$.

The second step is to split the large and medium-sized modules into module set
 P [16]. For the weighted directed graph G_d and a module M_q , $G_d(M_q)$ denotes the
 gene set of derived subgraphs in the directed graph G_d (the derived subgraph is
 different from the strongly connected subgraph), which corresponds to the genes in
 M_q . L denotes the set of derived subgraphs, as shown below.

$$L = \{G_d(M_q)\} \quad (16)$$

160 Let $split_size$ be the degree of the node with the largest value in the subgraph
 161 derived by module M_q , Modules with more nodes than $split_size$ will be split as
 162 large modules. In the splitting process, $G_c \in L$ is a subgraph derived from directed
 163 graph M_q , v' is the node with the largest out-degree value in G_c , and $IN(v')$

164 represents the local area network of v' in G_c . If the number of nodes in $IN(v')$ is
 165 not less than min_module_size , then they will be divided into seed modules, or else
 166 they will be divided into leaf modules. All the strongly connected subgraphs that
 167 meet the conditions in the directed graph G_c are split in the same way.

The third step is to add leaf modules to the seed module. The leaf node g_m
 connected to any node in the seed module is selected to extend the seed module by
 utilizing the greedy strategy, and the extension function is defined as follows:

$$G(g_m) = \overline{G^{in}(g_m)} - \overline{G^{out}(g_m)} \quad (17)$$

168 Where $\overline{G^{in}(g_m)}$ represents the average weight of the edge between node g_m in the
 169 leaf module and the node in the seed module, $\overline{G^{out}(g_m)}$ represents the average
 170 weight of the edge between node g_m in the leaf module and the rest of the nodes
 171 in the leaf module. If $\overline{G^{in}(g_m)}$ is not less than $\overline{G^{out}(g_m)}$, then the node is added to
 172 the seed module.

173 **Experimental results and analysis**

174 **Data preprocessing**

175 This study utilizes the somatic mutation data and the combined human PPI net-
 176 work data from HINT+HI2012 [14]. Somatic mutation data comes from the TCGA
 177 pan-cancer dataset containing 12 cancer types, which are composed of 3281 sam-
 178 ples with 20472 SNVs and 4334 samples with 720 CNAs. According to the data
 179 preprocessing method of [14], the samples with hyper-mutation and the genes that
 180 are low-expression in all tumor types are firstly screened out. And then MutSigCV

Algorithm 2 Driver modules detection

Input: $G_d(V, E, F)$, $total_genes$, min_module_size **Output:** Driver module set P

```

1: repeat
2:    $P = SCC(G_d)$  and  $M_q \in P$ 
3:    $P = P - M_q$  with  $|M_q| < min\_module\_size$ ,  $E = E - e$  with  $e = E_{min}$ 
4: until ( $|\cup_{M_q \in P} M_q| == total\_genes$ )
5:  $split\_size = outdeg(G_d(M_q))_{max}$ 
6: for  $M_q \in P$  and  $|M_q| > split\_size$  do
7:    $L = \{G_d(M_q)\}$  and  $G_c \in L$ 
8:   while  $L = \emptyset$  do
9:      $IN(v') \subseteq G_c$  with  $v' = outdeg(G_c)_{max}$ 
10:     $L = L - G_c$  and  $G_c = G_c - IN(v')$ 
     $seed_q = seed_q \cup IN(v')$  or  $leaf_q = leaf_q \cup IN(v')$ 
11:    for  $M_j \in SCC(G_c)$  do
12:       $seed_q = seed_q \cup M_j$  or  $leaf_q = leaf_q \cup M_j$  or  $L = L \cup M_j$ 
13:    for  $M_i \in leaf_q$  do
14:       $seed_q = seed_q \cup M_i$  with  $G(g_m) \geq 0$ 
15:     $P = seed_q$ 

```

181 tool is applied to filter the genes without obvious mutations in SNVs, after that 218
182 genes are deleted, 5 samples without SNV and 1973 samples with only CNA are
183 deleted, and 7894 genes with <3 RNA-seq reads in $> 30\%$ of tumors of each cancer
184 type are deleted. Finally, we obtain the data set containing 3110 samples and a
185 total of 11565 somatic mutation genes. (2) HINT+HI2012 as the PPI network data
186 includes high-quality interaction database (HINT) and human interaction database
187 (HI2012). The data preprocessing is as follows, a merge operation is firstly performed
188 based on the interaction relationship in the HINT and HI2012 database. Then, we
189 delete closed loops and duplicate edges in the new PPI network. Finally, we obtain
190 a protein interaction network composed of 9858 proteins and 40704 interactions.

191 Parameter setting

192 If the number of genes in a driver module is less than 3, the gene set is not usually
193 considered as a driver module [10, 14]. Therefore, *min_module_size* as the minimum
194 module size is set to be 3.

195 Enrichment analysis

196 To verify the enrichment effect of the modules identified by ECSWalk, we utilize
197 functional annotation tools (DAVID) to analyze the enrichment of the driver mod-
198 ules detected by ECSWalk, HotNet2, and MEXCOWalk algorithm. The results in
199 nine types of cancer are shown in Figure 1. The value is set to be 100 to obtain
200 the driver modules in the ECSWalk and MEXCOWalk, and HotNet2 obtains 14
201 consensus driver modules.

202 As shown in Figure 1, the enrichment effect of the driver modules detected by
203 ECSWalk is better than that by HotNet2 and MEXCOWalk algorithms among the
204 nine types of cancer, especially in GBM, Melanoma, UCEC, NSCLC, PAAD and
205 CML. Therefore, the driver modules identified by ECSWalk show extremely high
206 statistical significance and biological relevance.

207 Comparison of module accuracy

208 To evaluate the accuracy of the driver modules detected by ECSWalk, this study
209 uses the Accuracy and F-measure evaluation indices to measure the accuracy of
210 the driver modules based on the known pathways[24]. The higher the Accuracy
211 value, the better the classification effect. The higher the F-measure value, the more
212 driver modules can be enriched in the known biological pathways, indicating that

213 the method is more accurate in mining driver modules. The calculation formulas of
 214 Accuracy and F-measure are as follows:

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (18)$$

$$Precision = \frac{TP}{TP + FP} \quad (19)$$

$$Recall = \frac{TP}{TP + FN} \quad (20)$$

$$F - measure = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (21)$$

215 Where TP indicates the number of modules in which positive classes are predicted
 216 to be positive classes; FP indicates the number of modules in which negative classes
 217 are predicted to be positive classes; TN indicates the number of modules in which
 218 negative classes are predicted to be negative classes; FN indicates the number of
 219 modules in which positive classes are predicted to be negative classes.

The following formula is utilized to calculate the enrichment of the driver modules in one known biological pathway.

$$p - value = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}} \quad (22)$$

220 Where N represents the total number of genes, K represents the number of genes
221 in a known biological pathway, n represents the number of genes in a driver mod-
222 ule, and k represents the number of overlapping genes between a known biological
223 pathway overlaps and a driver module. The driver modules with p-value <0.01 are
224 set to be the positive type, and the driver modules with p-value ≥ 0.01 are set to
225 be the negative type. All the p-values are corrected using the Benjamin-Hochberg
226 method.

227 Figure 2 shows the enrichment performance of the driver modules obtained by
228 ECSWalk, MEXCOWalk and HotNet2 based on the known biological pathways
229 obtained by using the DAVID tool for enrichment analysis. It can be seen from
230 Figure 2A, the ACC value of the ECSWalk is 1 in the nine types of cancer, which
231 shows that ECSWalk has extremely high accuracy in detecting driver modules.
232 Specifically, the ACC value of the ECSWalk is 100%, 200% and 50% higher than
233 that of the HotNet2 respectively in the three cancers of GBM, BLCA and PAAD,
234 and the ACC value of the ECSWalk is 50%, 200%, 50%, 33.3%, 100% and 100%
235 higher than MEXCOWalk respectively in GBM, BLCA, UCEC, NSCLC, PAAD
236 and LAML. Also of note, the ACC value of the ECSWalk is 1, but the ACC values

237 of the MEXCOWalk and HotNet2 are 0 in CDAD. Therefore, ECSWalk has better
238 performance than MEXCOWalk and HotNet2 in detecting driver modules in cancer.

239 As can be seen in Figure 2B, the F-measure values of the ECSWalk algorithm
240 are 1 in the nine types of cancer, which means that the accuracy and recall of the
241 ECSWalk algorithm are both 1 in these nine cancers. This shows all the predicted
242 driver modules are positive. It is interesting to note that the F-measure values of
243 gene sets detected by MEXCOWalk and HotNet2 algorithms are both 0 in CDAD.
244 Therefore, ECSWalk has a better capability in detecting driver modules in nine
245 types of cancers.

246 Comparison of the optimal modules

247 *EGFR module*

248 The weighted diagram of the EGFR module detected by ECSWalk is shown in
249 Figure 3A. The EGFR module detected by ECSWalk, MEXCOWalk and HotNet2
250 is shown in Table 1.

251 It can be seen from Table 1, the p-value of the EGFR module detected by EC-
252 SWalk is 1.09E-13, and the p-values of the EGFR module detected by MEXCOWalk
253 and HotNet2 are 8.1E-07 and 6.9E-06, respectively. Although the number of genes
254 in the EGFR module detected by ECSWalk is less than that by MEXCOWalk, the
255 EGFR module detected by ECSWalk has higher coverage than those of the other
256 two algorithms. Therefore, the gene set detected by ECSWalk has higher biological
257 relevance and statistical significance than those by the other two algorithms.

258 The EGFR module detected by ECSWalk mutates in 62.77% (1952/3110) sam-
259 ples, in which PIK3CA, EGFR, APC, ERBB2 and PIK3R1 have high mutation

260 frequency, and mutate in 19.36%, 8.39%, 7.52%, 6.37% and 4.98% samples respec-
261 tively. This module is enriched in a variety of cancers, especially the p-value of this
262 module in UCEC is 1.09E-13 according to the DAVID enrichment analysis. It can be
263 also seen from Figure 3 that PIK3CA and NRAS, PIK3R1 have large edge weights,
264 and PIK3R1 and HRAS have a large edge weight, which indicates that the four
265 mutated genes have strong interaction relationships. In the ErbB signaling pathway
266 (Figure 4a), EGFR and IGF1R promote the expression of PIK3CA and PIK3R1,
267 and phosphorylation promotes the expression of HRAS and NRAS. At the same
268 time, HRAS and NRAS promote the expression of PIK3CA and PIK3R1. Besides,
269 CTNNB1 and APC have large weight values, which indicates that the two mutated
270 genes have a strong interaction relationship. In the Wnt signaling pathway (Fig-
271 ure 4b), CTNNB1 promotes the expression of itself and APC. This shows that the
272 mutations of CTNNB1 and APC constitute two key genes in the Wnt/ β -Catenin
273 pathway, and thus can be used as predictors of cervical cancer susceptibility [25].

274 *PPFIA1 module*

275 The weighted diagram of the PPFIA1 module detected by ECSWalk is shown in
276 Figure 3B. The PPFIA1 module detected by ECSWalk, MEXCOWalk and HotNet2
277 is shown in Table 2. HotNet2 does not mine any genes contained in the module,
278 and the gene set detected by ECSWalk contains one more gene PPFIA1 than that
279 by MEXCOWalk. The coverage of this module is 7.72%, and the mutual exclusiv-
280 ity is 96.77%. Although the PPFIA1 modules detected by the three algorithms do
281 not show enrichment information in the DAVID enrichment analysis, studies have
282 shown that there is an interaction between genes within the module [26, 27]. The

283 accuracy of the module can be verified in the previous research [27] in which the
284 surface tyrosine phosphatase receptor PTPRF and its adaptor PPFIA1 drive active
285 $\alpha 5 \beta 1$ integrin recycling and controls fibronectin fibrillogenesis and vascular mor-
286 phogenesis. Besides, the proteins encoded by PTPRF and PTPRS are important
287 members of the PTP family of protein tyrosine phosphatases. PTPS is a signal
288 molecule that regulates cell growth cycle, differentiation process, mitosis, gene mu-
289 tations and other cell life processes. The lack of PTPRS and PTPRF affects the
290 proliferation of mandibular cells, and results in craniofacial deformities. The WNT
291 and BMP signaling pathways are dysregulated in cells deficient in PTPRS and PT-
292 PRF [26]. Therefore, PPFIA1, PTPRS and PTPRF may be mutated in the same
293 pathway and cause cancer. In summary, the PPFIA1 module detected by ECSWalk
294 is of more biological relevance than those by the other two algorithms.

295 *TP53 module*

296 The weighted diagram of the TP53 module detected by ECSWalk is shown in Figure
297 3C. The TP53 module detected by ECSWalk, MEXCOWalk and HotNet2 is shown
298 in Table 3. The TP53 module detected by ECSWalk has higher coverage than
299 those by MEXCOWalk and HotNet2. The p-value of the TP53 module detected
300 by ECSWalk is $3.25E-20$ according to the DAVID enrichment analysis, and the
301 p-values detected by MEXCOWalk and HotNet2 are $5.6E-08$ and $3.84E-06$ respec-
302 tively. Therefore, the gene set detected by ECSWalk has higher biological relevance
303 and statistical significance than those by the other two algorithms.

304 The TP53 module detected by ECSWalk mutates in 70.09% (2180/3110) samples.
305 This module has significant enrichment in CML cancer, and the p-value of this mod-

306 ule in CML cancer is 1.1E-11 according to the DAVID enrichment analysis, which
307 indicates that the module has high biological relevance. It can be seen from Figure
308 3C, TP53, ATM, CHEK2, MDM2, MDM4, PTEN, CDKN2A, CDK4 and CDK6
309 have large edge weights, which indicates that the nine mutated genes have strong
310 relationships. In the P53 signaling pathway (Figure 4d), phosphorylation of ATM
311 promotes the expression of CHEK2, and phosphorylation of CHEK2 promotes the
312 expression of TP53; TP53 inhibits the expression of CDK4 and CDK6 by promot-
313 ing the expression of CDKN1A, and TP53 promotes the expression of upstream
314 MDM2 and the expression of downstream MDM2 and PTEN; MDM4 promotes
315 the expression of MDM2, but inhibits the expression of TP53; CDKN2A inhibits
316 the expression of MDM2, and MDM2 inhibits the expression of MDM4. This indi-
317 cates that regulating the p53 signaling pathway can interfere with CML, thereby
318 regulating the occurrence and development of CML [28].

319 Analysis of the remaining modules

320 *BAP1 module*

321 The weighted diagram of the BAP1 module detected by ECSWalk is shown in Figure
322 3D, and the mutual exclusivity of the BAP1 module is 96.95%. The BAP1 mod-
323 ule is a PR-DUB protein complex composed of BAP1 and ASXL1, which activates
324 its downstream tumor suppressor genes such as SOCS1/2, VHL and TXNIP by
325 combining with transcription factors[29]; In blood tumor cells, the BAP1 mutation
326 makes the C-terminal truncated ASXL1 mutation protein completely lose its ability
327 to bind to transcription factors, causes the ASXL1 mutation protein to significantly
328 weaken the transcription and regulation functions of the BAP1-ASXL1-FOXK1/K2

329 complex through a dominant negative mutation effect, reduces the expression of tu-
330 mor suppressor genes, thereby regulating glucose metabolism, hypoxia perception
331 and JAK-STAT and other tumor-related signaling pathways. This thus contributes
332 to promoting the proliferation and self-renewal of leukemia cells, and further in-
333 hibiting cell apoptosis under hypoxia [29]. It can be seen that ASXL1 and BAP1
334 have high biological relevance, and may exist in the same driver pathway to work
335 together and promote the occurrence of cancer.

336 *NOTCH3 module*

337 The weighted diagram of the NOTCH3 module is shown in Figure 3E. The NOTCH3
338 gene has high coverage and is altered in 129 samples. The mutual exclusivity of this
339 module is 98.82%, and it has been reported that the NOTCH signaling pathway
340 plays a vital role in the tumor microenvironment, and the ligand genes DLL1 and
341 JAG1 of the NOTCH signaling pathway have mutations in a variety of cancers, such
342 as DLL1, JAG1 and NOTCH3 existing high probability of mutations in patients
343 with colon cancer [30]. Besides, DLL1 and JAG1 can control the rate of neural
344 development in the NOTCH3 gene expression domain, and JAG1 can activate the
345 NOTCH signal transduction in the V1 and DL6 domains. DLL1 can also send
346 signals to nerve cells outside the NOTCH gene expression domain [31]. Therefore,
347 DLL1, JAG1 and NOTCH3 genes may exist in the same driver pathway.

348 *PAF1 module*

349 The weighted diagram of the PAF1 module detected by ECSWalk is shown in Figure
350 3F. The PAF1 module detected by ECSWalk mutates in 4.34% (135/3110) sam-

351 ples, and the PAF1 gene mutates in 104 samples. The mutual exclusivity between
352 genes in the module is 100%. This module has a significant enrichment relationship
353 in Cdc73/Paf1 complex, and the p-value of this module in Cdc73/Paf1 complex is
354 2.9E-9 according to GO enrichment analysis, which indicates that the genes in this
355 module have high biological relevance. Besides, CTR9, as the main component of
356 the RTF1 complex, participates in the assembly of the PAF1 complex by the TRP
357 domain, and Cdc73/Paf1 is a polyprotein complex associated with RNA polymerase
358 II and general RNA polymerase II transcription factor complexes[32], and it may
359 be involved in transcription initiation and extension. Furthermore, the mutated
360 genes PAF1 and CTR9 in the Cdc73/Paf1 complex can cause a variety of diseases
361 including malignant tumors. The PAF1 module also has a significant enrichment
362 relationship in Transcription elongation from RNA polymeraseII promoter path-
363 ways. The p-value of this module is 2.7E-8 according to GO enrichment analysis.
364 It has been reported in [33] that CTR9, RTF1 and LEO1 are the main members of
365 the Paf1/RNA polymerase II complex, PAF1, Rtf1 and LEO1 have frequent inter-
366 actions with PAF1, Cdc73 and PolIII, and the deletion of PAF1 or CTR9 will lead
367 to similar severe pleiotropic phenotypes. Therefore, the genes in this module have
368 high biological relevance and may exist in the same driver pathway to cause cancer.

369 *MAP3K1 module*

370 The weighted diagram of the MAP3K1 module detected by ECSWalk is shown in
371 Figure 3G. This module has significant enrichment in the MAPK signaling pathway,
372 and its p-value is 1.3E-3 according to the KEGG enrichment analysis. The phospho-
373 rylation of MAP3K1 in this module activates MAP2K4, MAP2K1 and MAP2K2,

374 and the phosphorylation of BRAF activates MAP2K1 and MAP2K2. It has been
375 reported in [34] that inhibiting the MAPK pathway can lead to lung cancer cell
376 apoptosis. Therefore, the genes in this module have high biological relevance and
377 may exist in the same driver pathway to cause cancer.

378 *MCL1 module*

379 The weighted diagram of the MCL1 module detected by ECSWalk is shown in
380 Figure 3H. It has been reported in [35] that USP9X and MCL1 in this module have
381 high biological relevance. For example, USP9X stabilizes the expression of MCL1
382 and promotes cell survival, and high expression of USP9X leads to an increase in
383 the amount of MCL1 protein in human follicular lymphoma and diffuse large B-cell
384 lymphoma, and vice versa. Therefore, USP9X is usually used as a target for clinical
385 treatment by maintaining the stability of the amount of MCL1 and other proteins
386 in human malignant tumors [35]. It can be seen that MCL1 and USP9X have high
387 biological relevance and may exist in the same driver pathway.

388 *MYC module*

389 The weighted diagram of the MYC module detected by ECSWalk is shown in Figure
390 3I. The module has more significant enrichment in the nucleoplasm, and the p-value
391 of this module in nucleoplasm is $8.3E-5$ according to GO enrichment analysis, which
392 indicates that this module has high biological relevance. It can be seen from Figure
393 3I that MYC has the largest coverage in the module, and MYC is altered in 9.10%
394 (283/3110) samples. It has been reported in [36] that MYC-encoded protein is a
395 multifunctional nucleolar phosphate protein, which acts as a transcription factor

396 regulating target genes during the cell life cycle and cell transformation process.
397 Overexpression, mutation, translocation and rearrangement of MYC are closely
398 related to the occurrence and development of a variety of cancers. It can be seen
399 from Figure 3I that FBXW7 and MYC have a large edge weight value, so these
400 two genes may likely have biological relevance. It has been reported in [37] that
401 FBXW7 loss of function contributes to worse overall survival and is associated with
402 accumulation of MYC in muscle invasive bladder cancer.

403 Discussion

404 This method identifies more new candidate gene sets for biological verification, and
405 the findings indicate that the identified candidate gene sets have high biological
406 relevance and statistical significance. The enrichment effect of the driver modules
407 detected by ECSWalk is better than that by HotNet2 and MEXCOWalk algorithms
408 among the nine types of cancer, especially in GBM, Melanoma, UCEC, NSCLC, and
409 CML. The EGFR module, PPFIA1 module and TP53 module detected by ECSWalk
410 have higher coverage than those by the other two algorithms. In particular, the
411 PPFIA1 module detected by ECSWalk contains more PTPRS genes compared with
412 the other two algorithms, which makes the PPFIA1 module more biologically and
413 statistically significant. Besides, we found that the EGFR module has extremely
414 high enrichment in UCEC(p-value: 1.09E-13), the TP53 module has extremely high
415 enrichment in CML(p-value: 1.1E-11). The results are of theoretical guidance and
416 practical value for cancer diagnosis, treatment and drug targets. There are certain
417 limitations in this study, although the greedy strategy can accurately identify the

418 driver genes belonging to the same module. This also leads to a decrease in the
419 operating efficiency of the algorithm to a certain extent.

420 **Conclusion**

421 This study proposes a carcinogenic driver module detection algorithm (ECSWalk)
422 by integrating somatic mutation data and protein interaction network. This study
423 firstly calculates the similarity between connected nodes in the protein interaction
424 network. Then, a weighted network is created by calculating mutual exclusivity,
425 coverage and topological structure similarity between mutation genes, and a restart
426 random walk clustering method is utilized to detect the carcinogenic driver mod-
427 ules. Finally, an induced subgraph method is utilized to split the large modules,
428 and a greedy strategy is utilized to expand the small modules to generate a set
429 of driver modules. The experimental results show that the ECSWalk can detect
430 more accurate driver modules than the other two algorithms. The driver modules
431 detected by ECSWalk have a lower p-value by the DAVID and Go enrichment analy-
432 sis, which indicates that the results of the algorithm have higher biological relevance
433 and statistical significance. It also shows that ECSWalk can achieve good results
434 by combining biological characteristics and complex network characteristics in de-
435 tecting carcinogenic driver modules. Therefore, the application of complex network
436 topology to biological networks is conducive to the study of the pathogenesis of
437 cancer based on the inherent properties of the data itself, and it is also conducive to
438 researchers and medical practitioners when carrying out research from the aspect
439 of complex network topology. Besides, ECSWalk can identify the target gene set in
440 some common cancers accurately, and analyze the biological relevance and statisti-

441 cal significance of the driver modules, which thus helps to enrich our understanding
442 of the pathogenesis of cancer.

443 **Acknowledgements**

444 We thank Jihua Dong for her careful proofreading, and also thank Bing Zhou, Zhaoheng A, Mengdi Liu, Pengyu
445 Zhang and Haoru Zhou for their helpful advice and discussions.

446 **Funding**

447 The work was supported by the National Natural Science Foundation of China (Grant No.61972322), the Natural
448 Science Foundation of Shaanxi Province (Grant No. 2021JM-110) and the Humanities and Social Science Fund of
449 Ministry of Education of China (Grant No.18YJCZH190). The funders had no role in study design, data collection
450 and analysis, decision to publish, or preparation of the manuscript.

451 **Abbreviations**

452 **ECSWalk**: A carcinogenic driver module detection method based on a network model.

453 **ICGC**: International Cancer Genome Consortium.

454 **TCGA**: The Cancer Genome Atlas.

455 **HotNet2**: An algorithm for finding significantly altered subnetworks in a large gene interaction network.

456 **MEXCOWalk**: Mutual exclusion and coverage based random walk to identify cancer modules.

457 **PPI**: Protein-protein interaction networks.

458 **HINT+HI2012**: A combination of high-quality protein-protein interactions from HINT and the recent HI-2012
459 set of protein-protein interactions.

460 **SCC**: Strongly connected component of the directed graph.

461 **KL**: Kullback-Leibler, a method of describing the difference between two probability distributions.

462 **JS**: Jensen-Shannon, an improved method based on KL divergence.

463 **DAVID**: The Database for Annotation, Visualization and Integrated Discovery.

464 **TP**: True Positive

465 **FP**: False Positive

466 **TN**: True Negative

467 **FN**: False Negative

468 **Availability of data and materials**

469 The data and materials we used can be download from <https://github.com/raphael-group/hotnet2>.(M.D.M.
470 Leiserson*, F. Vandin*, H.T. Wu, et al. (2014) Pan-Cancer Network Analysis Identifies Combinations of Rare
471 Somatic Mutations across Pathways and Protein Complexes. Nature Genetics 47, 106–114 (2015))

472 **Ethics approval and consent to participate**

473 Not applicable

474 Competing interests

475 The authors declare that they have no competing interests.

476 Consent for publication

477 Not applicable

478 Authors' contributions

479 Conceive and design the experiments: HW ZC. Perform the experiments: HW ZC. Analyze the data: ZC YW.

480 Contribute reagents/materials/analysis tools: ZC YW QL. Write the paper: HW ZC. Consult on the final version of

481 the paper and edit the paper: HW ZC HZ. The authors read and approve the final version of the manuscript.

482 Authors' information

483 **Hao Wu**, born in 1979, Ph.D., associate professor. His main research interests are within data mining, deep

484 learning, computational bioinformatics, complex networks and complex diseases, particularly in cancer genomics and

485 network biology. Email address: haowu@sdu.edu.cn.

486 **Zhong-Li Chen**, born in 1994, M.S. candidate. His main research interests are within computational

487 bioinformatics and biological big data mining. Email address: czl@nwfufu.edu.cn.

488 **Ying-Fu Wu**, born in 1998, M.S. candidate. His main research interests are within computational bioinformatics

489 and biological big data mining. Email address: wuyingfunwsuaf@163.com.

490 **Hong-Ming Zhang**, born in 1979, Ph.D., professor. His main research interests are within spatial big data

491 analysis and precision agriculture. Email address: zhm@nwsuaf.edu.cn.

492 **Quan-Zhong Liu**, born in 1978, Ph.D., associate professor. His main research interests are within bioinformatics

493 and data mining. Email address: liuqzhong@nwsuaf.edu.cn.

494 Author details

495 ¹College of Information Engineering, Northwest A&F University, Yangling, 712100, China. ²School of Software,

496 Shandong University, Jinan, 250100, China. ³Tibet Center for Disease Control and Prevention, Lhasa, 850000,

497 China.

498 References

499 1. Mckeage, M., Shepherd, P., Yozu, M., D, R.L.: Tumour mutation profiling with high-throughput multiplexed

500 genotyping: A review of its use for guiding targeted cancer therapy. *Current Cancer Therapy Reviews* **9**(4),

501 236–244 (2013)

502 2. Yu, X., Zeng, T., Li, G.: Integrative enrichment analysis: a new computational method to detect dysregulated

503 pathways in heterogeneous samples. *BMC Genomics* **16**(1), 918 (2015)

504 3. Zhang, J., Wu, L., Zhang, S., Zhang, S.: Discovery of co-occurring driver pathways in cancer. *BMC*

505 *Bioinformatics* **15**(1), 1–14 (2014)

506 4. Zhao, J., Zhang, S., Wu, L., Zhang, X.: Efficient methods for identifying mutated driver pathways in cancer.

507 *Bioinformatics* **28**(22), 2940 (2012)

- 508 5. Greenman, C., Stephens, P., Smith, R., Dalgliesh, G., Hunter, C., Bignell, H. Gand Davies, Teague, J., Butler,
509 A., Stevens, C., Edkins, S., O'Meara, S., Vastrik, I., Schmidt, E., Avis, T., Barthorpe, S., Bhamra, G., Buck,
510 G., Choudhury, B., Clements, J., Cole, J., Dicks, E., Forbes, S., Gray, K., Halliday, K., Harrison, R., Hills, K.,
511 Hinton, J., Jenkinson, A., Jones, D., Menzies, A., Mironenko, T., Perry, J., Raine, K.: Patterns of somatic
512 mutation in human cancer genomes. *Nature* **446**(7132), 153–158 (2007)
- 513 6. Vandin, F., Upfal, E., Raphael, B.: De novo discovery of mutated driver pathways in cancer. *Genome Research*
514 **22**(2), 175–181 (2012)
- 515 7. Leiserson, M., Blokh, D., Sharan, R., J. Raphael, B.: Simultaneous identification of multiple driver pathways in
516 cancer. *PLOS Computational Biology* **9**(5), 1003054 (2013)
- 517 8. Hou, J., Ma, J.: Dawnrank: discovering personalized driver genes in cancer. *Genome Medicine* **6**(7), 5 (2014)
- 518 9. Srihari, S., Ragan, M.: Systematic tracking of dysregulated modules identifies novel genes in cancer.
519 *Bioinformatics* **29**(12), 1553–1561 (2013)
- 520 10. Wu, H., Gao, L., Dong, J., Yang, X.: Detecting overlapping protein complexes by rough-fuzzy clustering in
521 protein-protein interaction networks. *Plos One* **9**(3), 91856 (2014)
- 522 11. Miller, C., Settle, S., Sulman, E., Aldape, K., Milosavljevic, A.: Discovering functional modules by identifying
523 recurrent and mutually exclusive mutational patterns in tumors. *BMC Medical Genomics* **4**(1), 34 (2011)
- 524 12. Wu, H.: Algorithm for detecting driver pathways in cancer based on mutated gene networks. *Chinese Journal of*
525 *Computers* **41**(6), 1400–1414 (2018)
- 526 13. Kim, Y., Cho, D., Dao, P., Przytycka, T.: Memcover: integrated analysis of mutual exclusivity and functional
527 network reveals dysregulated pathways across multiple cancer types. *Bioinformatics* **31**(12), 284–292 (2015)
- 528 14. Leiserson, M., Vandin, F., Wu, H., Dobson, J., Eldridge, J., Thomas, J., Papoutsaki, A., Kim, Y., Niu, B.,
529 McLellan, M., Lawrence, M., Gonzalez-Perez, A., Tamborero, D., Cheng, Y., Ryslik, G., Lopez-Bigas, N., Getz,
530 G., Ding, L., Raphael, B.: Pan-cancer network analysis identifies combinations of rare somatic mutations across
531 pathways and protein complexes. *Nature Genetics* **47**(2), 106–114 (2015)
- 532 15. Reyna, M., Leiserson, M., Raphael, B.: Hierarchical hotnet: identifying hierarchies of altered subnetworks.
533 *Bioinformatics* **34**(17), 972–980 (2018)
- 534 16. Rafsan, A., Ilyes, B., Cesim, E., Evis, H., Hilal, K.: Mexcwalk: Mutual exclusion and coverage based random
535 walk to identify cancer modules. *Bioinformatics* **36**(3), 872–879 (2019)
- 536 17. Wu, H., Gao, L., Li, F., Yang, X., Kasabov, N.: Identifying overlapping mutated driver pathways by
537 constructing gene networks in cancer. *BMC Bioinformatics* **16**(5), 3 (2015)
- 538 18. Nepusz, T., Yu, H., Paccanaro, A.: Detecting overlapping protein complexes in protein-protein interaction
539 networks. *Nature Methods* **9**(5), 471–472 (2012)
- 540 19. Guo, M., Wang, S., Liu, X., Tian, Z.: Algorithm for predicting the associations between mirnas and diseases.
541 *Journal of Software* **28**(11), 3094–3102 (2017)
- 542 20. Tang, D., Zhu, Q., Yang, F., Chen, K.: Efficient cluster analysis method for protein sequences. *Journal of*
543 *Software* **22**(8), 1827–1837 (2011)

- 544 21. Hou, Y., Duan, L., Li, L., Lu, L., Tang, C.: Search of genes with similar phenotype based on disease
545 information network. *Journal of Software* **29**(3), 721–733 (2018)
- 546 22. Zhang, Q., Li, M., Deng, Y.: A new structure entropy of complex networks based on tsallis nonextensive
547 statistical mechanics. *International Journal of Modern Physics C* **27**(10), 440–450 (2016)
- 548 23. Hofree, M., Shen, J., Carte, H., Gross, A., Ideker, T.: Network-based stratification of tumor mutations. *Nature*
549 *Methods* **10**, 1108–1115 (2013)
- 550 24. Li, F., Gao, L., Wang, B.: Detection of driver modules with rarely mutated genes in cancers. *IEEE/ACM*
551 *Transactions on Computational Biology and Bioinformatics* **17**(2), 390–401 (2020)
- 552 25. Wang, B., Wang, M., Li, X., M, Y., Liu, L.: Variations in the wnt/ β -catenin pathway key genes as predictors of
553 cervical cancer susceptibility. *Pharmacogenomics and personalized medicine* **13**, 157–165 (2020)
- 554 26. Katherine, S., Noriko, U., Wiljan, H., Michel, L., Bouchard, M.: Inactivation of lar family phosphatase genes
555 ptpns and ptpnf causes craniofacial malformations resembling pierre-robin sequence. *Development* **140**(16),
556 3413–3422 (2013)
- 557 27. Mana, G., Clapero, F., Panieri, E., Panero, V., Böttcher, R., Tseng, H., Saltarin, F., Astanina, E., Wolanska,
558 K., Morgan, M., Humphries, M., Santoro, M., Serini, G., Valdembrì, D.: Ppfia1 drives active $\alpha 5 \beta 1$ integrin
559 recycling and controls fibronectin fibrillogenesis and vascular morphogenesis. *Nature Communications* **7**(1),
560 13546 (2016)
- 561 28. Li, H., Liu, L., Liu, C., Zhuang, J., Zhou, C., Yang, J., Gao, C., Liu, G., Lv, Q., Sun, C.: Deciphering key
562 pharmacological pathways of qingdai acting on chronic myeloid leukemia using a network pharmacology-based
563 strategy. *Med Sci Monit* **24**, 5668–5688 (2018)
- 564 29. Xia, Y., Zeng, Y., Zhang, M., Liu, P., Liu, F., Zhang, H., He, C., Sun, Y., Zhang, J., Zhang, C., Song, L.,
565 Ding, C., Tang, Y., Yang, Z., Yang, C., Wang, P., Guan, K., Xiong, Y., Ye, D.: Tumor-derived neomorphic
566 mutations in asxl1 impairs the bap1-asxl1-foxk1/k2 transcription network. *Protein & Cell* (2020)
- 567 30. Wang, X., Xi, X., Wu, J., Wan, Y., Hui, H., Cao, X.: MicroRNA-206 attenuates tumor proliferation and migration
568 involving the downregulation of notch3 in colorectal cancer. *Oncology Reports* **33**(3), 1402–1410 (2015)
- 569 31. Catarina, R., Susana, R., Claudia, G., Domingos, H.: Two notch ligands, dll1 and jag1, are differently restricted
570 in their range of action to control neurogenesis in the mammalian spinal cord. *Plos One* **5**(11), 15515 (2010)
- 571 32. Amrich, C., Davis, C., Rogal, W., Shirra, M., Heroux, A., Gardner, R., Arndt, K., VanDemark, A.: Cdc73
572 subunit of paf1 complex contains c-terminal ras-like domain that promotes association of paf1 complex with
573 chromatin. *The Journal of biological chemistry* **287**(14), 10863–75 (2012)
- 574 33. Mueller, C., Jaehning, J.: Ctr9, rtf1, and leo1 are components of the paf1/rna polymerase ii complex. *Molecular*
575 *and Cellular Biology* **22**(7), 1971–1980 (2002)
- 576 34. Tsujino, I., Nakanishi, Y., Shimizu, T., Obana, Y., Ohni, S., Takahashi, N., Nemoto, N., Hashimoto, S.: 999
577 correlation between differences in the increase in mapk (erk1/2) activity due to driver mutations and prognosis
578 in non-small-cell lung cancer. *European Journal of Cancer* **48**, 241–241 (20012)
- 579 35. Schwickart, M., Huang, X., Lill, J., Liu, J., Ferrando, R., French, D., Maecker, H., O'Rourke, K., Bazan, F.,

- 580 Eastham-Anderson, J., Yue, P., Dornan, D., Huang, D., Dixit, V.: Deubiquitinase usp9x stabilizes mcl1 and
581 promotes tumour cell survival. *Nature* **463**(7277), 103–107 (2010)
- 582 36. Sabò, A., Kress, T., Pelizzola, M., De, P., Gorski, M., Tesi, A., Morelli, M., Bora, P., Doni, M., Verrecchia, A.,
583 Tonelli, C., Fagà, G., Bianchi, V., Ronchi, A., Low, D., Müller, H., Guccione, E., Campaner, S., Amati, B.:
584 Selective transcriptional regulation by myc in cellular growth control and lymphomagenesis. *Nature* **511**,
585 488–492 (2014)
- 586 37. Matumoto, T., Chen, Y., Contreras-Sanz, A., Ikeda, K., Schulz, G., Gao, J., Oo, H., Roberts, M., Costa, J.,
587 Nykopp, T.: Fbxw7 loss of function contributes to worse overall survival and is associated with accumulation of
588 myc in muscle invasive bladder cancer. *Urologic Oncology: Seminars and Original Investigations* **38**(12),
589 904–905 (2010)

590 Figures

Figure 1: Comparison of enrichment effect

Figure 2: Comparison of modules accuracy A) ACC values, B) F-measure values

Figure 3: Driver modules detected by ECSWalk.
(The genes in the rough edge represent the names of the dysregulated modules.
The size of the circle is proportional to the mutation frequency, and the
thickness of the line segment is proportional to the weight between nodes.)

Figure 4: Gene interaction diagram

591 Tables

Table 1: EGFR module

Algorithm	Gene set	Number of modules	Coverage	P-Value
HotNet2	EGFR ERBB2 AREG ELF3 ERBB4 LRIG1 OSMR	7	17.94%	6.9E-06
MEXCOWalk	ERBB2 CDKN2A FLNA ERBB4 ATM MCM2 IGF1R MDM4 PHF17 VHL NPM1 CDH1 STK11 MDM2 EGFR TLN1	16	45.98%	8.1E-07
ECSWalk	PIK3CA ERBB2 EGFR TLN1 PIK3R1 CTNNB1 IGF1R NEDD9 APC VARS2 HRAS NRAS	12	47.85%	1.09E-13

Table 2: PPFIA1 module

Algorithm	Gene set	Number of modules	Coverage	P-Value
HotNet2	N/A	N/A	N/A	N/A
MEXCOWalk	PPFIA1 PPP2R1A PTPRF	3	6.82%	N/A
ECSWalk	PTPRS PPFIA1 PPP2R1A PTPRF	4	7.72%	N/A

Table 3: TP53 module

Algorithm	Gene set	Number of modules	Coverage	P-Value
HotNet2	CCND1 CDKN2A CUL9 NPM1 PTEN TP53 AGL2 ALS2CR8 AMOTL1 ANKRD12 BRPF1 CACHD1 CARNS1 CDKN2AIP CELSR3 CHD8 EPHA3 HECW2 IFT140 IWS1 MAGI2 MAST3 MDM4 MLL5 PLEKHA8 PRDM2 PRKRIR SCAPER SETD2 SMG1 SMG5 SMG7 SNRK SPTBN2 STK11IP STRADA SWAP70 TEP1 TRIP12 TTLL5 UACA ZFP91 ZMIZ1 ZNF227 ZNF668	45	68.39%	3.84E-6
MEXCOWalk	E4F1 PTGS2 BMP1 ELL PPP1R13L WDR33 SMYD2 HINFP HIPK2 EGR1 KAT8 STK4 NOC2L EHMT1 CUL9 SNRPN CABLES1 WT1 WWOX CCT5 ARID3A HSPA9 ZNF384 RFWD2 TOP1 PLK3 RNF20 ERCC6 TOPORS TP53 RB1CC1 CDKN1B CUL1 KDM5A BRCA1 CTNNB1 CDK4 CCND1 RB1 CDK6	40	49.58%	5.6E-8
ECSWalk	BLM CSNK2A1 PTK2 WDR33 ING1 USP7 ING5 BRCA1 STAT1 UBE3A EP300 MDM2 MDM4 MED1 ATM CHD3 CREBBP WRN CUL9 IKBKB WT1 HSP90AA1 CHEK2 BCL2L1 CDK4 CDK6 ZNF384 CDKN2A XPO1 SIN3A HDAC2 HTT DDX5 TP53 AKT1 NPM1 RB1 TP53BP1 PTEN SMARCA4	40	70.09%	3.25E-20

Figures

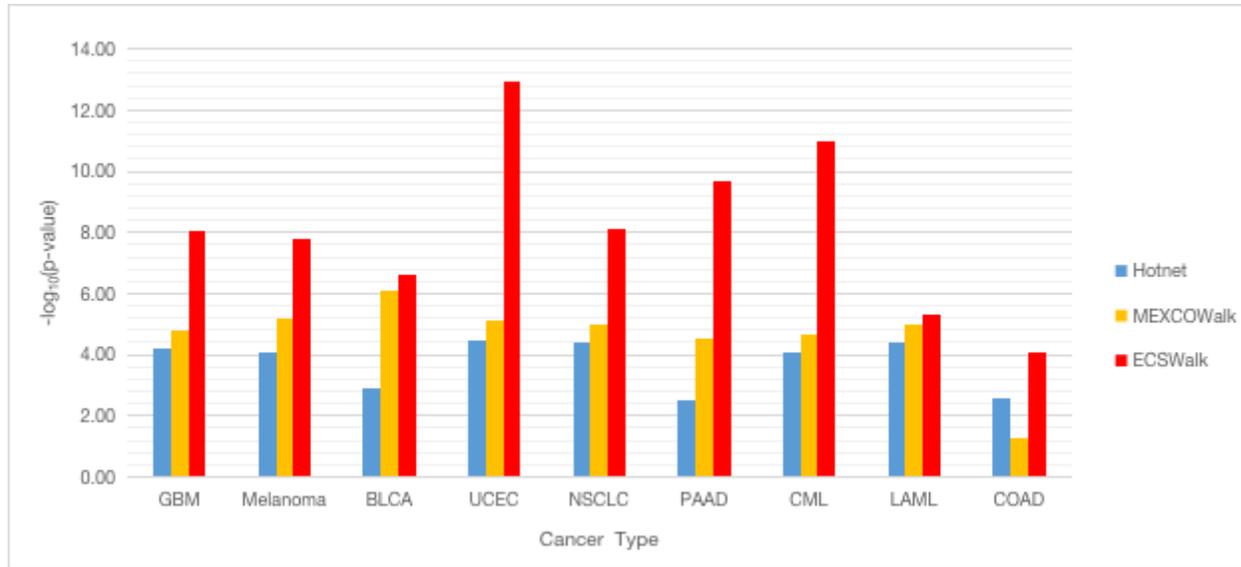


Figure 1

Comparison of enrichment effect

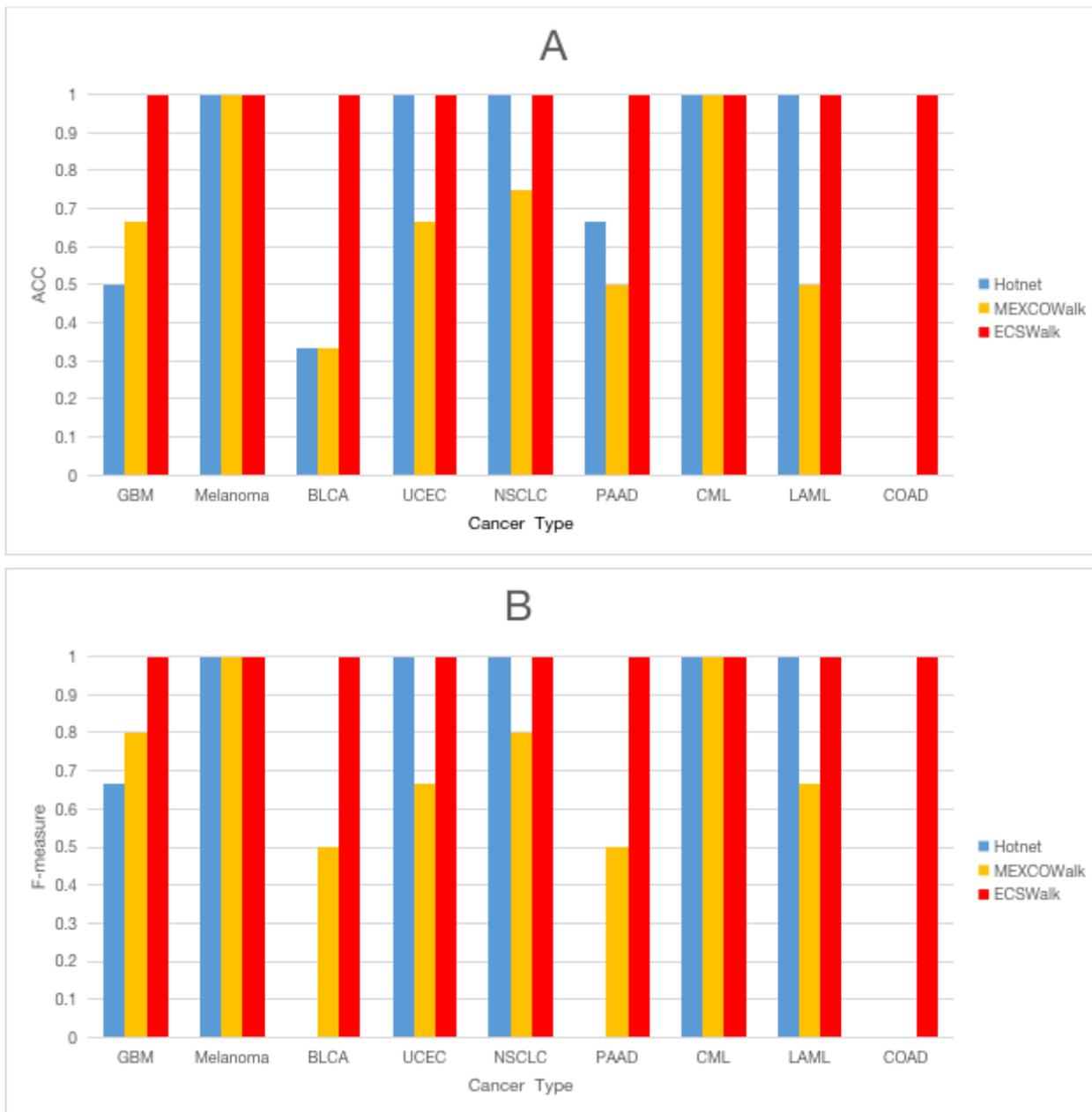


Figure 2

Comparison of modules accuracy A) ACC values, B) F-measure values

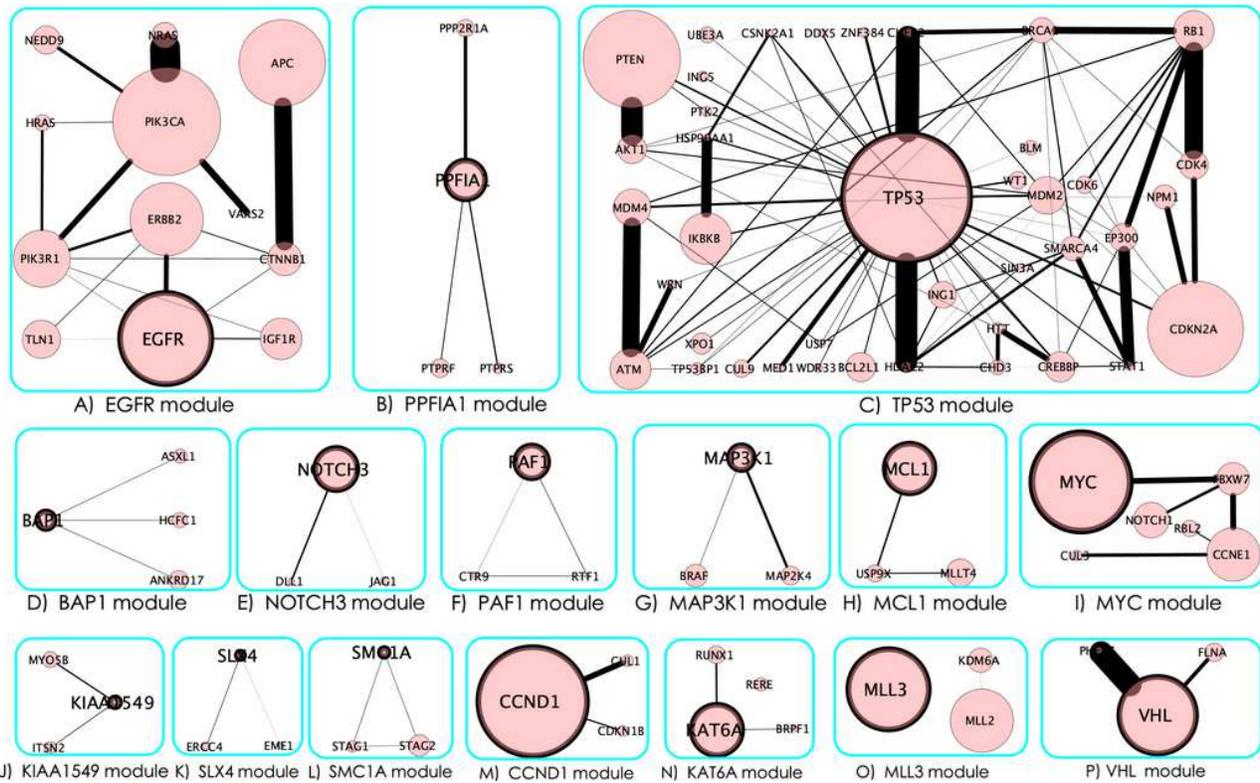


Figure 3

Driver modules detected by ECSWalk. (The genes in the rough edge represent the names of the dysregulated modules. The size of the circle is proportional to the mutation frequency, and the thickness of the line segment is proportional to the weight between nodes.)

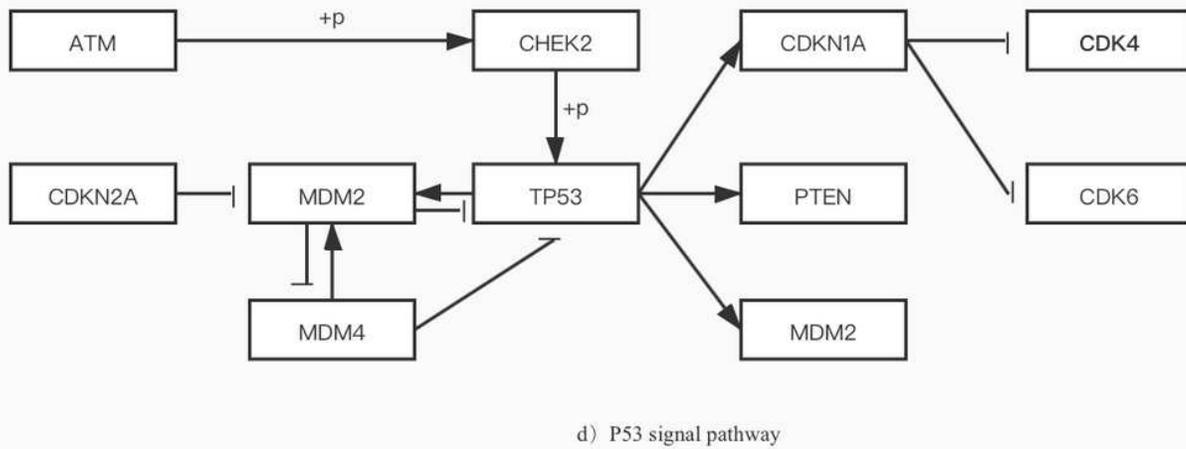
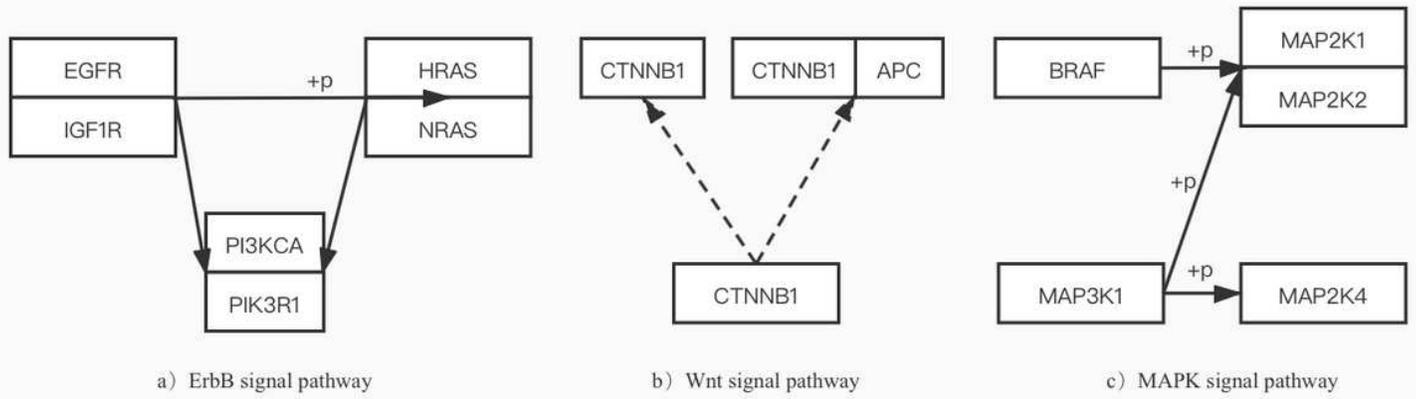


Figure 4

Gene interaction diagram