# Personalized analysis of breast cancer using sample-specific networks

**Ke Zhu** [Equal first author, 1] , **Cong Pian** [Equal first author, 1] , **Qiong Xiang** [1] , **Xin Liu** [1] , **Yuanyuan Chen** [Corresp. 1, 2]

[1] College of Science, Nanjing Agricultural University, Nanjing, Jiangsu, China

[2] State Key Laboratory of Bioelectronics, School of Biological Science and Medical Engineering, Southeast University, Nanjing, Jiangsu, China

Corresponding Author: Yuanyuan Chen
Email address: chenyuanyuan@njau.edu.cn

Breast cancer is a disease with high heterogeneity. Cancer is not usually caused by a single gene, but by multiple genes and their interactions with others and surroundings. Estimating breast cancer-specific gene-gene interaction networks is critical to elucidate the mechanisms of breast cancer from a biological network perspective. In this study, sample-specific gene-gene interaction networks of breast cancer samples were established by using a sample-specific network analysis method based on gene expression profiles. Then, gene-gene interaction networks and pathways related to breast cancer and its subtypes and stages were further identified. The similarity and difference among these subtype-related (and stage-related) networks and pathways were studied, which showed highly specific for subtype Basal-like and Stage IV and V. Finally, gene pairwise interactions associated with breast cancer prognosis were identified by a Cox proportional hazards regression model, and a risk prediction model based on the gene pairs was established, which also performed very well on an independent validation data set. This work will help us to better understand the mechanism underlying the occurrence of breast cancer from the sample-specific network perspective.

1 # Personalized analysis of breast cancer using sample-
2 # specific networks

3

4

5 Ke Zhu[1], Cong Pian[1], Qiong Xiang[1], Xin Liu[1], Yuanyuan Chen[1,2]

6

7 [1] College of Science, Nanjing Agricultural University, Nanjing, Jiangsu, China

8 [2] State Key Laboratory of Bioelectronics, School of Biological Science and Medical
9 Engineering, Southeast University, Nanjing, Jiangsu, China

10

11 Corresponding Author:

12 Yuanyuan Chen[1,2]

13 Xuanwu District Xiaolingwei Street Weigang No.1, Nanjing, Jiangsu, 210095, China

14 Email address: chenyuanyuan@njau.edu.cn

15

## 16 Abstract

17    Breast cancer is a disease with high heterogeneity. Cancer is not usually caused by a single
18 gene, but by multiple genes and their interactions with others and surroundings. Estimating breast
19 cancer-specific gene-gene interaction networks is critical to elucidate the mechanisms of breast
20 cancer from a biological network perspective. In this study, sample-specific gene-gene interaction
21 networks of breast cancer samples were established by using a sample-specific network analysis
22 method based on gene expression profiles. Then, gene-gene interaction networks and pathways
23 related to breast cancer and its subtypes and stages were further identified. The similarity and
24 difference among these subtype-related (and stage-related) networks and pathways were studied,
25 which showed highly specific for subtype Basal-like and Stage IV and V. Finally, gene pairwise
26 interactions associated with breast cancer prognosis were identified by a Cox proportional hazards
27 regression model, and a risk prediction model based on the gene pairs was established, which also
28 performed very well on an independent validation data set. This work will help us to better
29 understand the mechanism underlying the occurrence of breast cancer from the sample-specific
30 network perspective.

31

## Introduction

33      According to the latest data from the survey of the International Agency for Research on
34 Cancer (IARC) in 2018, the incidence of breast cancer is 24.2% among women worldwide, ranking
35 first in female cancers [1]. At present, the incidence of breast cancer is the highest, and its mortality
36 ranks fourth in China. Breast cancer has strong heterogeneity. Based on the TNM staging system,
37 breast cancer can be divided into Stages I, II, III, IV, and V. There are many clinical types of breast
38 cancer according to pathological classification and molecular classification. The pathological
39 classification generally divides breast cancer into invasive and non-invasive breast cancer. And
40 the gold standard for the molecular typing of breast cancer is PAM50 molecular typing based on
41 the expression profile of 50 genes, which classifies breast cancer into the Normal-like, LuminalA,
42 LuminalB, Basal-like, and Her2 subtypes [2].

43      The molecular typing of breast cancer has important reference value for clinical treatment of
44 breast cancer. However, molecular typing requires transcription sequencing which is difficult to
45 promote clinically. Currently, the diagnosis of breast cancer classification is mainly through
46 immunohistochemistry (IHC), namely, diagnosis by the expression of four markers, ER (oestrogen
47 receptor), PR (progestin receptor), HER2 gene (human epidermal growth factor receptor 2) and
48 Ki-67 protein (proliferating cell nuclear antigen). ER and PR are important indicators for endocrine
49 therapy and prognosis evaluation in breast cancer. Studies have shown that their expression are
50 positively correlated with total survival, treatment failure time, endocrine therapy response time,
51 and recurrence time [3, 4]. In 2009, Cheang used GEP (gene expression analysis) to determine
52 14% as the threshold of Ki-67, which could be used to divide patients into two groups with good
53 and bad prognoses [5]. In 2011, the St. Gallen International Expert Consensus agreed to include
54 Ki-67 as an important standard for molecular typing, which is the key to distinguishing the Luminal
55 A and Luminal B subtypes [6]. In the growth and metastasis of breast cancer, HER2 is one of the
56 most important factors, and its status can be used to predict the effect of drug treatment for breast
57 cancer. Early detection and diagnosis and timely treatment are of great significance to improve the
58 survival rate of breast cancer patients.

59      The aetiology of breast cancer is still not clear, and there are many related factors, such as
60 individual differences and a lack of effective treatments. With the development of biomedicine,
61 personalized medicine is becoming the direction of breast cancer treatment in the future. At
62 present, the medical plan can only be formulated through the study of single gene expression and
63 mutation information. However, this information cannot fully reflect the personalized interaction
64 and regulation among genes. Because onset and progression of cancer are often caused by the
65 disruption of important biological networks such as cell cycle and apoptosis, but not a single gene.
66 Indeed, there is a new and cutting-edge field of medical research, called network medicine, whose
67 basic idea is that human diseases are rarely caused by single molecular determinant, but more
68 likely influenced by a network of interacting molecular determinants with the propensity to cluster

69 together in the human interactome [7-9]. Gene-gene interaction networks can reveal the interaction
70 relations and regulatory mechanisms among genes. And they have the irreplaceable function of
71 the single-gene monitoring of information (such as expression and mutation) in many aspects [10].
72 Therefore, the mechanism of the occurrence and development of breast cancer can be explored
73 through changes in the interactions between genes. In this paper, we constructed sample-specific
74 networks of breast cancer samples by calculating the correlation coefficient of protein-coding gene
75 pairs to explore the gene-gene interaction networks related to breast cancer stages and subtypes
76 (see Fig.1).

77　　　 The survival time of different patients with breast cancer is significantly different. At present,
78 the 5-year survival rate of breast cancer patients in China has reached 83.2%. However, the 5-year
79 survival rate of advanced cancer patients and Basal-like cancer patients are significantly lower, so
80 it is necessary to study the biomarkers that affect the prognosis of breast cancer. In 2009, Joel S.
81 Parker et al. established a single-gene level survival analysis model to improve the prognosis of
82 breast cancer and predict the efficacy of chemotherapy [11]. However, the robustness of the gene-
83 based model is not very high. Thus, this paper aims to establish a more stable prognostic analysis
84 model of breast cancer patients through gene-gene interactions. We used the differential
85 correlation coefficients to model the prognosis of breast cancer. Lasso regression is suitable for
86 data analysis and model construction with many independent variables but a limited sample size
87 [12]. In this study, we used a Lasso regression model to effectively reduce the dimensionality of
88 large gene pairs and then identified the gene interactions related to the prognosis of breast cancer.
89 Finally, a multivariate Cox proportional hazards regression analysis based on the gene interactions
90 was carried out to predict the survival of patients with breast cancer (see Fig.1B). A prognosis
91 model was established and it also performed very well on an independent validation data.
92

93 ## Materials & Methods

94 ### Datasets

95　　　 In this paper, the RNA sequencing (RNA-seq) data of 290 normal breast tissues was
96 downloaded from the GTEx database (https://gtexportal.org/home/), and the RNA-seq data of
97 1093 breast cancer samples was downloaded from the TCGA database
98 (https://portal.gdc.cancer.gov/). A human protein-protein interaction network was from the
99 STRING database version 11.0 (https://string-db.org/), and gene sets of all available186 KEGG
100 pathways were downloaded from the GSEA/MSigDB database (http:
101 //software.broadinstitute.org/gsea/msigdb). In addition, the clinical information of the breast
102 cancer patients was downloaded from the TCGA database, including TNM stage, prognosis
103 survival time and other information. The 290 normal breast tissues were used as reference samples.
104 The gene expression data sets of normal and cancer samples were both converted to the TPM form
105 and contain 18006 genes in total. The independent validation data of the prognosis model was
106 from the GSE3494 set in GEO Datasets, which contains 251 expression profiles of breast tumors
107 by array.

108 **Construction of sample-specific networks**

109    In this study, gene-gene interactions with high confidence (comprehensive score >0.9) were
110 selected from the STRING database, which includes regulatory, physical and co-expression
111 protein-protein interaction networks. Furthermore, the above gene-gene interactions with both
112 genes in one of the 186 KEGG pathways were used as the background network (or template
113 network), which contained 3257 genes in total. The sample-specific network method aimed to
114 calculate the difference of the gene co-expression when the single cancer sample was added to a
115 bunch of normal samples. In short, the sample-specific networks to be constructed are actually
116 networks with significant perturbation edges of gene co-expression.

117    In the following analysis, sample-specific networks for breast cancer samples were
118 constructed based on gene expression profiles by using the method introduced in reference [10]
119 (see Fig.1A). First, using the gene expression data of $n$ reference samples, namely all the normal
120 breast tissues data, the reference network can be constructed by calculating the correlation
121 coefficient $PCC_n$ (the Pearson correlation coefficient (PCC)) of the gene pairs connected in the
122 background network. The weights of the edges in the reference network are the PCC of the
123 corresponding gene pairs. Then, the expression data of a single breast cancer sample was added to
124 the reference samples, and the perturbed network of the single sample was constructed by
125 calculating the new correlation coefficient $PCC_{n+1}$ of the gene pairs in the background network.
126 For the single breast cancer sample, the differential correlation coefficients of each edge between
127 the perturbed network and the reference network were calculated as: $\Delta PCC_n = PCC_{n+1} - PCC_n$,
128 which called differential network for the sample. In reference [10], Liu et al. have proved that
129 $\Delta PCC_n$ follows a normal distribution with a mean value of 0 and a variance of $\frac{1 - PCC_n^2}{n-1}$ when $n$ is
130 large enough. The significance level of each $\Delta PCC$ was determined by the $Z$-test. The statistical $Z$-
131 value is calculated as follows with the null hypothesis that $\Delta PCC_n$ is equal to 0:

132
$$Z = \frac{\Delta PCC_n}{(1 - PCC_n^2)/(n-1)}.$$

133    Then, we can obtain the $P$-value for each gene pair from the $Z$-value. Gene pairs (or edges)
134 were considered statistically significant if their $P$-values < 0.01. All significant edges constitute
135 the sample-specific network. Thus, adding the expression data of 1093 breast cancer samples to
136 the reference samples one at a time, we finally constructed 1093 sample-specific networks.

137 **Identification of stage/subtype-related gene-gene interaction networks**

138    Only the gene pairs that are perturbed significantly in the most breast cancer samples are
139 considered to be related to breast cancer. Then, the edges that are perturbed significantly in more
140 than 90% of the samples by the binomial right-sided test ($P$-value <0.05) constitute a gene-gene
141 interaction network related to breast cancer. Specifically, we firstly divided the above breast cancer
142 samples into different stages or subtypes based on TNM staging and PAM50 subtype system,
143 secondly selected the perturbed significantly edges in more than 90% samples of different stages
144 or subtypes, and then the stage/subtype gene-gene interaction networks were constructed.

145    A slight change in the expression of high-degree genes in the network may cause disturbances
146    of the entire network. Thus, these genes with high degree are considered to be the key genes for
147    the onset and development of breast cancer. We selected genes with degrees >5 in the identified
148    breast cancer-related network for the subsequent enrichment analysis. Furthermore, we also
149    identified the key genes related to each TNM stage and PAM50 subtype with the same method.
150    Here, because of the small number of stage V samples, stage V was combined with stage IV.

151    **Pathway enrichment analysis**

152    For the pathway enrichment analysis, we used the hypergeometric test as follows:

153
$$p(m,M,N,n) = 1 - \sum_{i=0}^{m-1} \frac{\binom{M}{i}\binom{N-M}{n-i}}{\binom{N}{n}},$$

154    where $N$ is the total number of genes in the background network, $M$ represents the number of key
155    genes related to breast cancer (or a stage or subtype of breast cancer), $n$ accounts for the number
156    of genes in a pathway, and $m$ represents the number of genes that both in the pathway and in key
157    genes related to breast cancer (or a stage or subtype of breast cancer). Then, the pathway with $P$-
158    value <0.05 was considered as significantly enriched in the breast cancer (or a stage or subtype of
159    breast cancer) samples. Otherwise, we regarded that the pathway is not enriched in the
160    corresponding group.

161    **Survival analysis by the Cox regression model**

162    Different from the usual survival analysis based on gene expression, the perturbation of gene
163    co-expression $\Delta PCC$ (i.e. gene pairs or edges) was used to survival analysis. According to the
164    clinical data of patients with breast cancer, we utilized the "survival" package and "survminer"
165    package in R/Bioconductor to establish a univariate Cox proportional hazards regression model by
166    setting patients' survival conditions (survival time and survival status) as the dependent variables
167    and the $\Delta PCC$ of gene pairs in the differential network for each breast cancer samples as the
168    covariates. Gene pairs with $P$-values < 0.05 were considered to be related to the prognosis of breast
169    cancer [13].

170    A large number of covariates may cause overfitting in establishing a multivariate Cox
171    proportional hazards regression model; thus, using the least absolute shrinkage and selection
172    operator (LASSO), we further selected the key gene pairs from these significant ones obtained by
173    the univariate Cox proportional hazards analysis. LASSO is a common method used in high-
174    dimensional data regression, which can select prognosis-related gene pairs of breast cancer by
175    shrinking regression coefficients. The tuning parameter ($\lambda$) with the smallest mean-square error
176    was selected by four-fold cross-validation to establish an optimal LASSO regression model. Then,
177    the coefficients of most gene pairs reduced to zero, and a smaller number of gene pairs with
178    nonzero coefficients were considered to be closely correlated with the prognosis of breast cancer.
179    LASSO Cox analysis was performed by using the "glmnet" package in R. Then, the risk score
180    for each sample was calculated by the LASSO Cox regression model. According to the median
181    risk score, breast cancer patients were divided into two groups (a high-risk group and a low-risk

182  group). In addition, 234 breast tumors with relapse free survival information in the validation data
183  set were analyzed by using the above sample-specific network method, and risk scores were
184  calculated by the Cox regression model based on 1093 samples in TCGA. Then the validation
185  samples were also divided into two groups in the same way. Finally, the corresponding Kaplan-
186  Meier survival curves were plotted by using the packages "survminer" and "survival" in R.
187

## Results

189  **Breast cancer-related gene-gene interaction networks**

190      The background network consisted of 46916 edges and 3237 genes. In addition, 2190 gene
191  pairs were identified as significantly related to breast cancer, which constituted the gene-gene
192  interaction network related to breast cancer (including 915 genes in total). We use the Cytoscape
193  software to visualize the breast cancer-related network (see Fig.2).

194      Genes with degrees > 5 in the breast cancer-related gene-gene interaction network (198 in
195  total which are shown in Table S1). Among them, some genes with higher degrees (> 20) have
196  been shown to be related to breast cancer. For example, CCNB1 has strong power to predict the
197  survival of breast cancer patients with the phenotype of ER positive [1414]. The overexpression
198  of GRB2 has been demonstrated to be significantly associated with the occurrence and poor
199  prognosis of breast cancer [15]. PCNA has been proven to be a marker of proliferation in the
200  diagnosis of breast cancer [16], SF3B4 has been shown to be a tumour suppressor, and somatic
201  inactivating mutations occasionally occur in breast cancer [1717]. UBE2C may promote the
202  development of breast cancer [18]. High Cdc20 and securin immune expression are associated
203  with extremely poor outcomes in breast cancer patients [19], and overexpression of RPL17 affects
204  breast cancer-associated brain metastases [20]. MAD2L1 may have great effect on breast cancer
205  progression, and its expression might help to predicting breast cancer prognosis [2121]. The high
206  expression of TRA2B is closely related to the cancer cell survival and therapeutic sensitivity of
207  breast cancer [22]. GTF2H4 has been identified to be related to the survival risk of breast cancer
208  [23].

209  **Stage-related gene-gene interaction networks**

210      The results of the four stage-related gene-gene interaction networks are shown in Fig.S1A-D.
211  And the top 10 genes with the highest degrees in these four networks are displayed in Fig.S1E-H.
212  There are obvious similarities and differences among the four stage-related gene interaction
213  networks. There are 81 key genes shared by all stages (see Table S3), among which RPL17,
214  CCNB1, and SF3B4 are genes that are highly (with degrees > 25) related to breast cancer. Stage I
215  has 5 specific genes: PSMC5, SDHB, RPL11, SDHA, and RPL13. Stage II has 3 specific genes:
216  STX6, CCNA2, and CDC25C. Stage III has 4 specific genes: NDUFA6, EPN1, SF3A3, and
217  LSM7. Stage IV has the largest number of specific genes, with a total of 38, among which CDC42,

218 LSM2, NDUFS6, and CDC25A are strongly associated with it. And these stage-specific key genes
219 are shown in Table S4.

**Subtype-related gene-gene interaction networks**

221 The results of the four subtype-related gene interaction networks are shown in Fig.S2A-D.
222 And the top 10 genes with the highest degrees in these four networks are displayed in Fig.S2E-H.
223 The four subtype-related networks share similar and different characteristics. There are 34 key
224 genes shared by the four subtypes (see Table S7). Among them, RPL17 and CCNB1 have higher
225 degrees. The Luminal A subtype has 11 specific genes, including RPL23A, RPL10, and PRPF6,
226 which are greatly related to it with higher degrees. The Luminal B subtype has 17 specific genes,
227 including COX6C, EGFR, and CLTC, which are related to it with higher degrees. The Her2
228 subtype has 3 specific genes, NDUFA6, CCR8, and CASP3. The Basal-like subtype has the largest
229 number of specific genes, 17 in total, including LSM2, DDX5, SF3A3, and MAGOH, with higher
230 degrees. And these subtype-specific key genes are shown in Table S8.

**Pathways enriched in breast cancer patients**

232 There were 41 pathways (see Table S2) enriched in the breast cancer samples according to
233 the pathway enrichment analysis, including some immune-related pathways, such as the Toll-like
234 receptor signaling pathway, antigen processing and presentation, complement and coagulation
235 cascades, the RIG-I-like receptor signaling pathway, and the cytosolic DNA-sensing pathway.
236 Some important signal transduction and signal molecular interaction pathways were also included,
237 such as the MAPK signaling pathway, Wnt signaling pathway, cytokine-cytokine receptor
238 interaction, and ECM-receptor interaction pathways. Breast cancer is closely related to endocrine
239 disorders [24], two endocrine-related pathways, adipocytokine signaling pathway, and PPAR
240 signaling pathway, have also been identified as being related to breast cancer. In addition, some
241 metabolic pathways, especially lipid metabolism pathways, have also been identified as being
242 associated with breast cancer [25], such as the steroid hormone biosynthesis, arachidonic acid
243 metabolism, arginine and proline metabolism pathway, and glycerolipid metabolism. Additionally,
244 pathways in cancer was also enriched. The enrichment results are shown in Fig.3A.

245 Most of these pathways have been documented to be related to breast cancer. For example,
246 the dysregulation of the steroid hormone biosynthesis pathway may affect steroid hormone levels
247 and may thus be related to the susceptibility to breast cancer [24]. The PPAR signaling pathway
248 may play an important role in the neoadjuvant chemotherapy response of breast cancer [26].
249 Mounting preclinical evidence supports targeting the MAPK signaling pathway in the triple
250 negative breast cancer (TNBC) [27]. AMPK activators inhibit breast cancer cell proliferation by
251 inhibiting DVL3-promoted Wnt/β-catenin signaling pathway activity [28]. Toll-like receptors may
252 play dual roles in human cancers [29]. The co-activation of the Hedgehog and Wnt signaling
253 pathways is a poor prognostic marker in TNBC [30]. Prl-3 is closely related to cell migration and
254 invasion in TNBC [31]. The YHD inhibition of 4T1 breast tumour growth may be related to the
255 negative regulation of the JAK/STAT3 pathway by repressing the expression of IL-6 and TGF-$β$
256 [32].

**Stage-related pathways**

The overlapping of pathways enriched in the four TNM stages are shown in Fig.3B. The proportion of enriched pathways shared by the four stages (see Table S5) is relatively high, including the Wnt signaling pathway, MAPK signaling pathway27, regulation of actin cytoskeleton, calcium signaling pathway34, pathways in cancer, and cell adhesion molecules, which have been shown to have a high correlation with breast cancer [27, 28, 33-35]. The pathways enriched in different stages are slightly different, especially Stage IV of breast cancer, which has 18 specific enriched pathways, among which the PPAR signaling pathway26, ECM-receptor interaction, tight junction, TGF-beta signaling pathway, NOD-like receptor signaling pathway, and other signaling pathways are mostly related to the metastases of breast cancer [26, 36-39].

As we expected, stage IV was specifically enriched the most pathways (18 in total, see Table S6) different from other stages. This result is probably because Stage IV breast cancer patients are the most serious, and their cancer cells are likely to have deteriorated and metastasized. Therefore, the disruption of the biological system balance of breast cancer patients at this stage is larger than that of other stages. Thus, the specific enriched pathways of Stage IV are correspondingly more.

**Subtype-related pathways**

The overlap and difference of the enriched pathways in the four PAM50 subtypes are shown in Fig.3C. There are slight differences in the subtype-related pathways. There are 4 enriched pathways shared by the four subtypes (see Table S9) including the cytokine-cytokine receptor interaction. As a special subtype of breast cancer, the Basal-like subtype (or TNBC) is characterized by high histological differentiation, a high risk of metastasis, a high recurrence rate, and a low survival rate. Probably due to the higher risk of Basal-like subtype, there are 9 specific pathways enriched in it, including the leukocyte transendothelial migration and chemokine signaling pathway. The subtype-specific enriched pathways are shown in Table S10.

**Prognosis-related gene pairs**

A total of 5652 gene pairs significantly related to the survival and prognosis of breast cancer were found by the univariate Cox proportional hazards model. In addition, 272 gene pairs were further identified by Lasso regression (see Fig.S3). A multivariate Cox proportional hazards regression model with these gene pairs as independent variables was constructed as follows.

$$Score = 206.3 * (ENO, PGK2) + 35.9 * (EN0, PKLR) + 4.1 * (EBP, HSD17B7) + 5.5 * (CYP1B, HSD17B1) - 3.4 * (NDUFB2, NDUFB4) - 0.6 * (ATP6V1A, ATP6V1B1) + ...$$

The risk scores of the 1093 breast cancer patients in TCGA were calculated by this model. The median of the risk scores divided all patients into two groups. The corresponding Kaplan-Meier survival curve is shown in Fig.4A. Of note, survival analysis indicates that overall survival probability of patients with high risk scores is significantly lower than that with low risk scores ($P$-value < 0.0001).

In addition, the risk scores of the 234 breast tumors in the validation data set were also calculated by the above model with 264 gene pairs (8 gene pairs were omitted since these genes were not included in the expression profile of the validation data). In the same way, there are two groups with different scores. The relapse free survival probabilities of the two groups are significantly different ($P$-value < 0.0001), and the relapse free survival status of tumors in the low

298    score group are all "alive" (see Fig. 4B). This result indicates that the prognosis model based on
299    gene pairs can well predict the survival time of breast tumors in the independent validation data
300    set.
301

## Discussion

303        At present, research on cancer pathology is limited to gene expression and mutation
304    information. However, the model of one gene to one disease is no longer suitable for the study of
305    complex diseases. In fact, genes do not exist in isolation but participate in some complex biological
306    networks, such as gene-gene interaction networks. Gene mutations or surroundings changes often
307    affect the balance of gene interaction networks and the perturbation of the networks then affect the
308    onset and development of complex diseases. Studies have shown that some genetic elements of
309    breast cancer are related to nearby gene expression, such as some repetitive DNA in ER+/HER2-
310    breast cancer and transposable elements [40, 41]. Therefore, network analysis can provide a more
311    comprehensive and systematic point of view, to better understand the human disease onset and
312    development mechanism.

313        Based on personalized medicine, Precision Medicine is a new medical concept and medical
314    model, which needs to grasp the specific characteristic of different cancer samples accurately. The
315    analysis of the biological network disturbance for each cancer patient conforms to the concept of
316    precision medicine. In addition, the personalized medical treatment of breast cancer is in a
317    relatively slow development stage.

318        In this paper, sample-specific networks of breast cancer samples were established to explore
319    the gene-gene interaction networks related to the TNM stages and PAM50 subtypes of breast
320    cancer. Then, the pathways related to breast cancer were identified by hypergeometric test.
321    Through the same method, we also obtained the stage-related pathways and subtype-related
322    pathways. Finally, the edge biomarkers (gene pairs) that are closely related to the prognosis of
323    breast cancer were determined by using the LASSO regression model, and then a more stable
324    prognostic analysis model was established by using these biomarkers. Our results indicate that the
325    prognosis model has the robust and strong generalization capability, and it can be used in different
326    gene expression data sets.
327        Many studies have shown that network-based methods are more robust and effective than
328    single-gene-based methods, such as SWIM and WGCNA [42, 43]. SWIM is a tool able to extract
329    from complex correlation networks the so-called "switch genes" that could be associated to the
330    transition from physiological to a pathological condition. The WGCNA method plans to exploit
331    the correlation patterns among genes. The advantages of network-based methods have been well
332    documented and accepted in the analysis of noisy high-throughput data. Different from the usual
333    network-based method, we made better use of a prior background network to explore the sample-
334    specific networks. And the sample-specific networks are actually networks with significant
335    perturbations edges of gene co-expression in our study, which is really very different from

336    WGCNA. This study helps us to better understand the heterogeneity and mechanism of breast
337    cancer from an individual-level perspective. Precision medicine advocates the development of
338    individualized treatment according to the unique features of patients. Therefore, identifying the
339    unique pathogeny embedded in each patient is important to develop a treatment strategy for each
340    patient. Our sample-specific network analysis of breast cancer will promote the development of
341    precision medicine.

342

## Conclusions

344    In this paper, the sample-specific network of each breast cancer sample was constructed based
345    on network analysis, and further breast cancer (subtype/stage)-related gene-gene interaction
346    networks were identified. The edge biomarkers (gene pairs) related to the prognosis of breast
347    cancer were also identified and a risk prediction model was established based on these edge
348    biomarkers finally.

349    This study develops an individualized network analysis for each patient which would
350    promote a new train of thought and method for the precision medicine. This whole process of
351    sample-specific network analysis using co-expression can also be used to analyze other
352    cancers. However, the co-expression perturbation which used to construct sample-specific
353    network, does not roundly measure the changes of gene interactions. So, we will consider further
354    designing a method which can characterize the perturbation of gene interactions comprehensively.
355    In addition, how to obtain subtype-specific networks (or stage-specific networks) from sample-
356    specific networks based on network structure is still a problem worth considering.

357

## Acknowledgements

360

## Additional information and declarations

365    **Author contributions:** All authors discussed the results and commented on the manuscript.
366    **Conflicts of Interest:** The authors declare no conflicts of interest.
367    **Supplementary Materials:** Supplementary information for this article are available in the
368    supplemental files. **Table S1**: Key genes of BRCA (the genes with degrees > 5 in the breast cancer-

369 related gene-gene interaction network). **Table S2**: Pathways enriched in BRCA (the pathways with
370 *P*-value <0.05 by the enrichment analysis in all breast cancer samples). **Table S3**: Key genes
371 shared by four stages. **Table S4**: Stage-specific key genes. **Table S5**: Enriched pathways shared
372 by four stages. **Table S6**: Stage-specific enriched pathways. **Table S7**: Key genes shared by four
373 subtypes. **Table S8**: Subtype-specific key genes. **Table S9**: Enriched pathways shared by four
374 subtypes. **Table S10**: Subtype-specific enriched pathways. **Figure S1**: Gene-gene interaction
375 networks and bar charts of gene degrees related to breast cancer TNM stages. **Figure S2**: Gene-
376 gene interaction networks and bar charts of gene degrees related to breast cancer PAM50 subtypes.
377 **Figure S3**: Establishment of the LASSO regression model.
378

# References

380 1.  Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A. Global cancer statistics
381     2018: GL OBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185
382     countries. *CA Cancer J Clin*, 2018. 68(6):394- 424.
383 2.  Perou CM, Sorlie T, Eisen MB, van de Rijn M, Jeffrey SS, Rees CA, Pollack JR, Ross DT,
384     Johnsen H, Akslen LA, Fluge O, Pergamenschikov A, Williams C, Zhu SX, Lønning PE,
385     Børresen-Dale AL, Brown PO, Botstein D. Molecular portraits of human breast tumors.
386     *Nature*, 2000. 406(6797): 747-752.
387 3.  Hammond ME, Hayes DF, Dowsett M, Allred DC, Hagerty KL, Badve S, Fitzgibbons PL,
388     Francis G, Goldstein NS, Hayes M, Hicks DG, Lester S, Love R, Mangu PB, McShane L,
389     Miller K, Osborne CK, Paik S, Perlmutter J, Rhodes A, Sasano H, Schwartz JN, Sweep FC,
390     Taube S, Torlakovic EE, Valenstein P, Viale G, Visscher D, Wheeler T, Williams RB, Wittliff
391     JL, Wolff AC. American Society of Clinical Oncology/College of American Pathologists
392     guideline recommendations for immunohistochemical testing of estrogen and progesterone
393     receptors in breast cancer. *J Clin Oncol*, 2010.28(16):2784-95.
394 4.  Fitzgibbons PL, Murphy DA, Hammond ME, Allred DC, Valenstein PN. Recommendations
395     for validating estrogen and progesterone receptor immunohistochemistry assays. *Arch Pathol
396     Lab Med,* 2010. 134(6):930-5.
397 5.  Cheang MC, Chia SK, Voduc D, Gao D, Leung S, Snider J, Watson M, Davies S, Bernard PS,
398     Parker JS, Perou CM, Ellis MJ, Nielsen TO. Ki-67 index, HER2 status, and prognosis of
399     patients with luminal B breast cancer. *J Natl Cancer Inst*, 2009. 101(10):736-750.
400 6.  Goldhirsch A, Wood WC, Coates AS, Gelber RD, Thürlimann B, Senn HJ; Panel members.
401     Strategies for subtypes-dealing with the diversity of breast cancer: highlights of the St.Gallen
402     International Expert Consensus on the Primary Therapy of Early Breast Cancer. *Ann Oncol,*
403     2011. 22(8):1736-1747.
404 7.  Albert-László Barabási, Natali Gulbahce, Joseph Loscalzo. Network medicine: a network-
405     based approach to human disease. *Nature Reviews Genetics*. 2011. 12:56-68.

406   8.   Conte F, Fiscon G, Licursi V, Bizzarri D, D'Antò T, Farina L, Paci P. A paradigm shift in
407        medicine: A comprehensive review of network-based approaches. *Biochim Biophys Acta*
408        *Gene Regul Mech.* 2019. 194416.
409   9.   Giulia Fiscon, Federica Conte, Lorenzo Farina, Paola Paci. Network-Based Approaches to
410        Explore Complex Biological Systems towards Network Medicine. *Genes.* 2018. 9(9):437.
411   10.  Liu X, Wang Y, Ji H, Aihara K, Chen L. Personalized characterization of diseases using
412        sample-specific networks. *Nucleic Acids Research*, 2016. 44(22): e164.
413   11.  Parker JS, Mullins M, Cheang MC, Leung S, Voduc D, Vickery T, Davies S, Fauron C, He
414        X, Hu Z, Quackenbush JF, Stijleman IJ, Palazzo J, Marron JS, Nobel AB, Mardis E, Nielsen
415        TO, Ellis MJ, Perou CM, Bernard PS. Supervised Risk Predictor of Breast Cancer Based on
416        Intrinsic Subtypes. *J Clin Oncol,* 2009. 27(8): 1160–1167.

417   12.  Zhang, Shangli & Zhang, Lili & Qiu, Kuanmin & Lu, Ying & Cai, Baigen. Variable Selection
418        in Logistic Regression Model. Chinese Journal of Electronics. 2015. 24(4):813-817.

419   13.  Cheng P. A prognostic 3-long noncoding RNA signature for patients with gastric cancer. *J*
420        *Cell Biochem.* 2018. 119(11):9261-9269.
421   14.  Ding K, Li W, Zou Z, Zou X, Wang C. CCNB1 is a prognostic biomarker for ER+ breast
422        cancer. *Med Hypotheses.* 2014. 83(3):359-64.
423   15.  Zhang Y, Xu G, Liu G, Ye Y, Zhang C, Fan C, Wang H, Cai H, Xiao R, Huang Z, Luo Q.
424        miR-411-5p inhibits proliferation and metastasis of breast cancer cell via targeting GRB2.
425        *Biochem Biophys Res Commun.* 2016. 476(4):607-613.
426   16.  Juríková M, Danihel Ľ, Polák Š, Varga I. Ki67, PCNA, and MCM proteins: Markers of
427        proliferation in the diagnosis of breast cancer. *Acta Histochem.* 2016. 118(5):544-52.
428   17.  Denu RA, Burkard ME. Synchronous Bilateral Breast Cancer in a Patient With Nager
429        Syndrome. *Clin Breast Cancer.* 2017. 17(3):151-153.
430   18.  Mo CH, Gao L, Zhu XF,Wei KL, Zeng JJ, Chen G, Feng ZB. The clinicopathological
431        significance of UBE2C in breast cancer: a study based on immunohistochemistry, microarray
432        and RNA-sequencing data. *Cancer Cell Int.* 2017. 17:83.
433   19.  Karra H, Repo H, Ahonen I, Löyttyniemi E, Pitkänen R, Lintunen M, Kuopio T, Söderström
434        M, Kronqvist P. Cdc20 and securin overexpression predict short-term breast cancer survival.
435        *Br J Cancer.* 2014. 110(12):2905-13.
436   20.  Yuan F, Wang W, Cheng H. Co-expression network analysis of gene expression profiles of
437        HER2[+] breast cancer-associated brain metastasis. *Oncol Lett.* 2018. 16(6):7008-7019.
438   21.  Wang Z, Katsaros D, Shen Y, Fu Y, Canuto EM, Benedetto C, Lu L, Chu WM, Risch HA, Yu
439        H. Biological and Clinical Significance of MAD2L1 and BUB1, Genes Frequently Appearing
440        in Expression Signatures for Breast Cancer Prognosis. *PLoS One.* 2015. 10(8):e0136246.
441   22.  Best A, Dagliesh C, Ehrmann I, Kheirollahi-Kouhestani M, Tyson-Capper A, Elliott DJ.
442        Expression of Tra2 β in Cancer Cells as a Potential Contributory Factor to Neoplasia and
443        Metastasis. *Int J Cell Biol.* 2013. 2013:843781.

444  23. Ge J, Liu H, Qian D, Wang X, Moorman PG, Luo S, Hwang S, Wei Q. Genetic variants of
445      genes in the NER pathway associated with risk of breast cancer: A large-scale analysis of 14
446      published GWAS datasets in the DRIVE study. *Int J Cancer.* 2019. 145(5):1270-1279.
447  24. Sakoda LC, Blackston C, Doherty JA, Ray RM, Lin MG, Stalsberg H, Gao DL, Feng Z,
448      Thomas DB, Chen C. Polymorphisms in steroid hormone biosynthesis genes and risk of breast
449      cancer and fibrocystic breast conditions in Chinese women. *Cancer Epidemiol Biomarkers*
450      *Prev.* 2008. 17(5):1066-73.
451  25. Merdad A, Karim S, Schulten HJ, Jayapal M, Dallol A, Buhmeida A, Al-Thubaity F, GariI
452      MA, Chaudhary AG, Abuzenadah AM, Al-Qahtani MH. Transcriptomics profiling study of
453      breast cancer from Kingdom of Saudi Arabia revealed altered expression of Adiponectin and
454      Fatty Acid Binding Protein4: Is lipid metabolism associated with breast cancer? *BMC*
455      *Genomics.* 2015.  16 Suppl 1: S11.
456  26. Chen YZ, Xue JY, Chen CM, Yang BL, Xu QH, Wu F, Liu F, Ye X, Meng X, Liu GY, Shen
457      ZZ, Shao ZM, Wu J. PPAR signaling pathway may be an important predictor of breast cancer
458      response to neoadjuvant chemotherapy. *Cancer Chemother Pharmacol.* 2012. 70(5):637-44.
459  27. Giltnane JM, Balko JM. Rationale for targeting the Ras/MAPK pathway in triple-negative
460      breast cancer. *Discov Med.* 2014. 17(95):275-83.
461  28. Zou YF, Xie CW, Yang SX, Xiong JP. AMPK activators suppress breast cancer cell growth
462      by inhibiting DVL3-facilitated Wnt/β-catenin signaling pathway activity. *Mol Med Rep.* 2017.
463      15(2):899-907.
464  29. Khademalhosseini M, Arababadi MK. Toll-like receptor 4 and breast cancer: an updated
465      systematic review. *Breast Cancer.* 2019. 26(3):265-271.
466  30. Bhateja P, Cherian M, Majumder S, Ramaswamy B. The Hedgehog Signaling Pathway: A
467      Viable Target in Breast Cancer? *Cancers (Basel)*. 2019. 11(8):1126.
468  31. Gari HH, DeGala GD, Ray R, Lucia MS, Lambert JR. PRL-3 engages the focal adhesion
469      pathway in triple-negative breast cancer cells to alter actin structure and substrate adhesion
470      properties critical for cell migration and invasion. *Cancer Lett.* 2016. 380(2):505-12.
471  32. Mao D, Feng L, Gong H. The Antitumor and Immunomodulatory Effect of Yanghe Decoction
472      in Breast Cancer Is Related to the Modulation of the JAK/STAT Signaling Pathway. *Evid*
473      *Based Complement Alternat Med.* 2018. 2018:8460526.
474  33. Kazazian K, Go C, Wu H, Brashavitskaya O, Xu R, Dennis JW, Gingras AC, Swallow CJ.
475      Plk4 Promotes Cancer Invasion and Metastasis through Arp2/3 Complex Regulation of the
476      Actin Cytoskeleton. *Cancer Res.* 2017. 77(2):434-447.
477  34. Woltmann A, Chen B, Lascorz J, Johansson R, Eyfjörd JE, Hamann U, Manjer J, Enquist-
478      Olsson K, Henriksson R, Herms S, Hoffmann P, Hemminki K, Lenner P, Försti A. Systematic
479      pathway enrichment analysis of a genome-wide association study on breast cancer survival
480      reveals an influence of genes involved in cell adhesion and calcium signaling on the patients'
481      clinical outcome. *PLoS One.* 2014. 9(6): e98229.

482    35. Saadatmand S, de Kruijf EM, Sajet A, Dekker-Ensink NG, van Nes JG, Putter H, Smit VT,
483        van de Velde CJ, Liefers GJ, Kuppen PJ. Expression of cell adhesion molecules and prognosis
484        in breast cancer. *Br J Surg.* 2013. 100(2):252-60.
485    36. Bao Y, Wang L, Shi L, Yun F, Liu X, Chen Y, Chen C, Ren Y, Jia Y. Transcriptome profiling
486        revealed multiple genes and ECM-receptor interaction pathways that may be associated with
487        breast cancer. *Cell Mol Biol Lett.* 2019. 24:38.
488    37. Yang Y, Liu L, Fang M, Bai H, Xu Y. The chromatin remodeling protein BRM regulates the
489        transcription of tight junction proteins: Implication in breast cancer metastasis. *Biochim
490        Biophys Acta Gene Regul Mech.* 2019. 1862(5):547-556.
491    38. Tang X, Shi L, Xie N, Liu Z, Qian M, Meng F, Xu Q, Zhou M, Cao X, Zhu WG, Liu B. SIRT7
492        antagonizes TGF-β signaling and inhibits breast cancer metastasis. *Nat Commun.* 2017.
493        8(1):318.
494    39. Peng L, Hu Y, Chen D, Linghu R, Wang Y, Kou X, Yang J, Jiao S. Ubiquitin specific protease
495        21 upregulation in breast cancer promotes cell tumorigenic capability and is associated with
496        the NOD-like receptor signaling pathway. *Oncol Lett.* 2016. 12(6):4531-4537.

497    40. Yandım C, Karakülah G. Dysregulated expression of repetitive DNA in ER+/HER2- breast
498        cancer. *Cancer Genet.* 2019. 239:36-45.
499    41. Karakülah G, Arslan N, Yandım C, Suner A. TEffectR: an R package for studying the potential
500        effects of transposable elements on gene expression with linear regression model. *PeerJ.* 2019.
501        7: e8192.

502    42. Paci P, Colombo T, Fiscon G, Gurtner A, Pavesi G, Farina L. SWIM: a computational tool to
503        unveiling crucial nodes in complex biological networks. *Sci Rep.* 2017. 7(1):44797.
504    43. Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network
505        analysis. *BMC Bioinformatics.* 2008. 9(1):559.
506

# Figure Legends

## Figure 1

An integrative framework identifying breast cancer-related gene-gene interaction networks.

**(A)** Construction of sample-specific networks based on gene expression data. A reference network can be established based on the expression profiles of $n$ reference samples by calculating the correlation coefficients $PCC_n$ of gene pairs. Then, adding a new sample $sd_x$ into the reference samples, a perturbed network is established by calculating the new correlation $PCC_{n+1}$ of the $n + 1$ samples. Because of sample $sd_x$, the perturbed network is different from the reference network, and the difference $\Delta PCC_n($ $PCC_{n+1} - PCC_n)$ of each edge in the background network constitutes the differential network. Then, the significance of each edge can be quantified by a statistical Z-test. The sample-specific network for sample $sd_x$ is composed by those edges with significant $\Delta PCC_n$.

**(B)** The framework to identify the breast cancer-related gene-gene interaction network based on gene expression. Using the sample-specific network analysis method, $m$ cancer sample-specific networks were constructed. Then, these constructed sample-specific networks were analysed to identify breast cancer-related networks, stage-related networks and subtype-related networks, as well as gene-interaction biomarkers associated with the prognosis of breast cancer. Moreover, pathway enrichment analysis based on KEGG pathways and survival analysis based on the LASSO regression model were performed.

## Figure 2

Gene-gene interaction networks related to breast cancer. Nodes in these networks stand for genes, and the size of the nodes corresponds to the degree of the genes in the network. The purple nodes represent the genes with degrees ≥ 15, and the blue ones are the genes with degrees < 15.

# Figure 3

Pathways enriched in breast cancer, as well as different stages and subtypes of it.

**(A)** KEGG pathways enriched in breast cancer samples, ranked by *-log10(p)*.

**(B)** Overlap and difference of the enriched pathways in the four breast cancer stages. There are 11 commonly enriched pathways in the four stages. The number of Stage IV-specific pathways was 18.

**(C)** Overlap and difference of the enriched pathways in the four PAM50 subtypes. There are 4 commonly enriched pathways in the four PAM50 subtypes. The number of Basal-like specific pathways is 9.

# Figure 4

Kaplan-Meier survival analysis.

**(A)** Kaplan-Meier survival plots for two different groups of breast cancer patients in TCGA.The X axis is survival days. The Y axis is overall survival rate.

**(B)** Kaplan-Meier survival plots for two different groups of breast tumors in the independent validation data set. The X axis is relapse free survival time (days). The Y axis is relapse free survival rate.

**Figure S1.** Gene-gene interaction networks related to breast cancer TNM stages. Nodes in these networks stand for genes, and the size of the nodes corresponds to the degree of the genes in the network. The purple nodes represent the genes with degrees ⩾ 15, and the blue ones are the genes with degrees < 15.

**(A-D)** Gene-gene interaction networks associated with Stage I, II, III, and IV respectively.

**(E-H)** The bar charts of top 10 genes with the highest degrees in gene-gene interaction networks related to Stage I, II, III, and IV.  The Y axis is gene, and the X axis is the gene degree.

**Figure  S2.** Gene-gene  interaction  networks  related  to  breast  cancer  PAM50 subtypes.  Nodes  in  these  networks  stand  for  genes,  and  the  size  of  the  nodes

corresponds to the degree of the genes in the network. The purple nodes represent the genes with degrees $\geqslant$ 15, and the blue ones are the genes with degrees < 15.

**(B-D)** Gene-gene interaction networks associated with LumA, LumB, Her2, and Basal-like subtypes respectively.

**(E-H)** The bar charts of top 10 genes with the highest degrees in gene-gene interaction networks related to LumA, LumB, Her2, and Basal-like subtypes. The Y axis is gene, and the X axis is the gene degree.

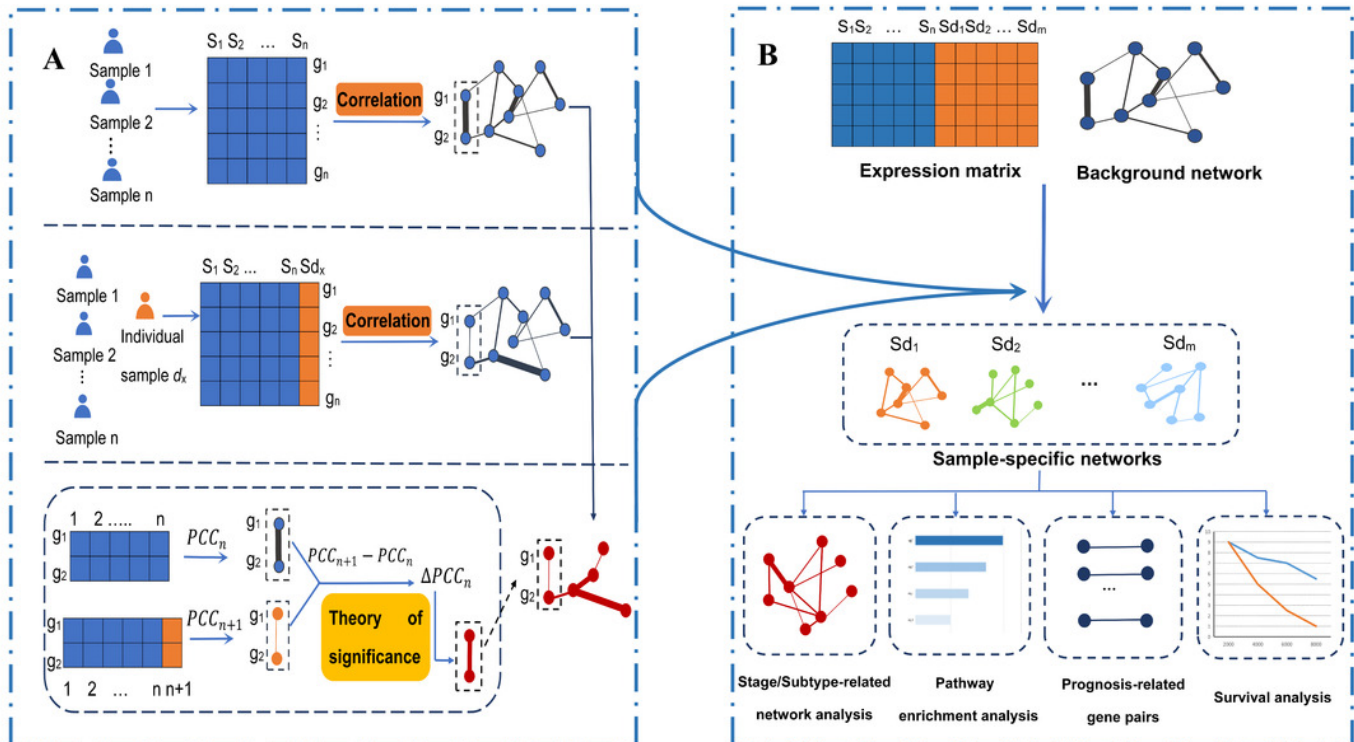## Figure S3. Establishment of the LASSO regression model.

**(A)** Four-fold cross-validation for tuning parameter ($\lambda$) selection in the LASSO model.

**(B)** LASSO coefficient profiles of 272 gene interactions.

# Figure 1

An integrative framework identifying breast cancer-related gene-gene interaction networks.

(A) Construction of sample-specific networks based on gene expression data. A reference network can be established based on the expression profiles of $n$ reference samples by calculating the correlation coefficients $PCC_n$ of gene pairs. Then, adding a new sample $sd_x$ into the reference samples, a perturbed network is established by calculating the new correlation $PCC_{n+1}$ of the $n+1$ samples. Because of sample $sd_x$, the perturbed network is different from the reference network, and the difference $\Delta PCC_n$ ( $PCC_{n+1}$ - $PCC_n$) of each edge in the background network constitutes the differential network. Then, the significance of each edge can be quantified by a statistical Z-test. The sample-specific network for sample sdx is composed by those edges with significant $\Delta PCC_n$. (B) The framework to identify the breast cancer-related gene-gene interaction network based on gene expression. Using the sample-specific network analysis method, $m$ cancer sample-specific networks were constructed. Then, these constructed sample-specific networks were analysed to identify breast cancer-related networks, stage-related networks and subtype-related networks, as well as gene-interaction biomarkers associated with the prognosis of breast cancer. Moreover, pathway enrichment analysis based on KEGG pathways and survival analysis based on the LASSO regression model were performed.

# Figure 2

Gene-gene interaction networks related to breast cancer.

Nodes in these networks stand for genes, and the size of the nodes corresponds to the degree of the genes in the network. The purple nodes represent the genes with degrees ≥ 15, and the blue ones are the genes with degrees < 15.
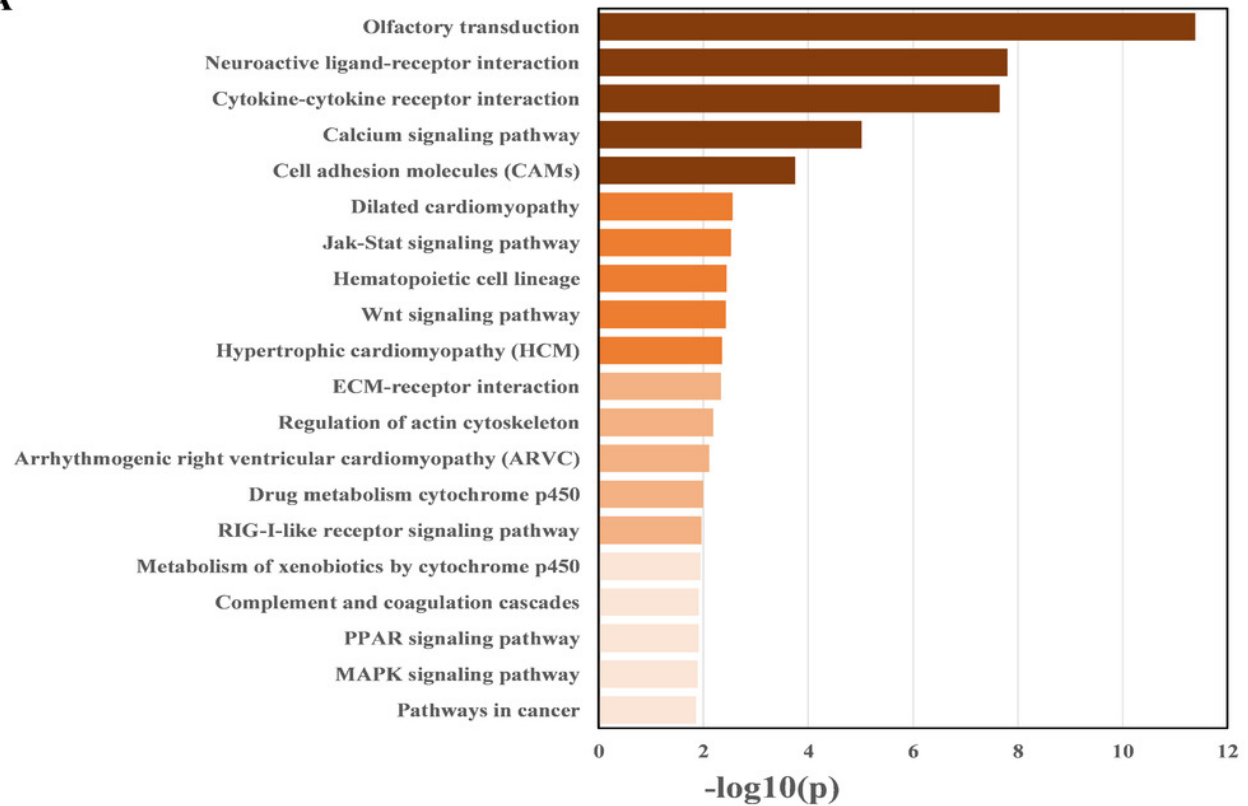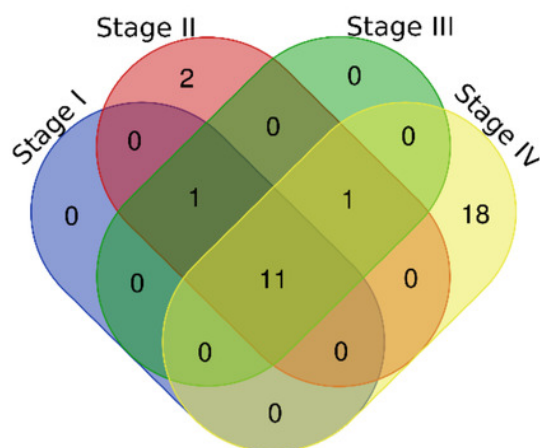
# Figure 3

Pathways enriched in breast cancer, as well as different stages and subtypes of it.

**(A)** KEGG pathways enriched in breast cancer samples, ranked by *-log10(p)*.**(B)** Overlap and difference of the enriched pathways in the four breast cancer stages. There are 11 commonly enriched pathways in the four stages. The number of Stage IV-specific pathways was 18.**(C)** Overlap and difference of the enriched pathways in the four PAM50 subtypes. There are 4 commonly enriched pathways in the four PAM50 subtypes. The number of Basal-like specific pathways is 9.
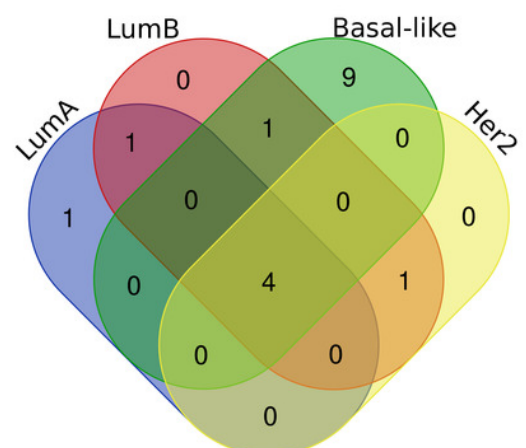
# Figure 4

Kaplan-Meier survival analysis.

**(A)** Kaplan-Meier survival plots for two different groups of breast cancer patients in TCGA.The X axis is survival days. The Y axis is overall survival rate. **(B)** Kaplan-Meier survival plots for two different groups of breast tumors in the independent validation data set. The X axis is relapse free survival time (days). The Y axis is relapse free survival rate.