

Personalized analysis of breast cancer using individual-level 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, individual-level gene-gene interaction networks of breast cancer samples were established by using an individual-level 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 individual-level network perspective.

Personalized analysis of breast cancer using individual-level networks

Ke Zhu¹, Cong Pian¹, Qiong Xiang¹, Xin Liu¹, Yuanyuan Chen^{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^{1,2}

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

Email address: chenyuanyuan@njau.edu.cn

Abstract

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, individual-level gene-gene interaction networks of breast cancer samples were established by using an individual-level 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 individual-level network perspective.

31

32 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. The
40 gold standard for the molecular typing of breast cancer is PAM50 molecular typing based on the
41 expression profile of 50 genes, which classifies breast cancer into the Normal-like, LuminalA,
42 LuminalB, Basal-like, and Her2 subtypes [**Error! Reference source not found.**].

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 (progesterone 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 individual gene expression
63 and mutation information. However, this information cannot fully reflect the personalized
64 interaction and regulation among genes. Because onset and progression of cancer are often caused
65 by the disruption of important biological networks such as cell cycle and apoptosis, but not a single
66 gene. Gene-gene interaction networks can reveal the interaction relations and regulatory
67 mechanisms among genes. And they have the irreplaceable function of the single-gene monitoring
68 of information (such as expression and mutation) in many aspects [7]. Therefore, the mechanism

of the occurrence and development of breast cancer can be explored through changes in the interactions between genes. In this paper, we constructed individual-level networks of breast cancer samples by calculating the correlation coefficient of protein-coding gene pairs to explore the gene-gene interaction networks related to breast cancer stages and subtypes (see Fig.1).

The survival time of different patients with breast cancer is significantly different. At present, the 5-year survival rate of breast cancer patients in China has reached 83.2%. However, the 5-year survival rate of advanced cancer patients and Basal-like cancer patients is significantly lower, so it is necessary to study the biomarkers that affect the prognosis of breast cancer. In 2009, Joel S. Parker et al. established a single-gene level survival analysis model to improve the prognosis of breast cancer and predict the efficacy of chemotherapy [8]. However, the robustness of the gene-based model is not very high. Thus, this paper aims to establish a more stable prognostic analysis model of breast cancer patients through gene-gene interactions. We used the differential correlation coefficients to model the prognosis of breast cancer. Lasso regression is suitable for data analysis and model construction with many independent variables but a limited sample size [Error! Reference source not found.]. In this study, we used a Lasso regression model to effectively reduce the dimensionality of large gene pairs and then identified the gene interactions related to the prognosis of breast cancer. Finally, a multivariate Cox proportional hazards regression analysis based on the gene interactions was carried out to predict the survival of patients with breast cancer (see Fig.1B). An prognosis model was established and it also performed very well on an independent validation data.

Materials & Methods

Datasets

In this paper, the RNA sequencing (RNA-seq) data of 290 normal breast tissues was downloaded from the GTEx database (<https://gtexportal.org/home/>), and the RNA-seq data of 1093 breast cancer samples was downloaded from the TCGA database (<https://portal.gdc.cancer.gov/>). A human protein-protein interaction network was from the STRING database version 11.0 (<https://string-db.org/>), and gene sets of 186 KEGG pathways were downloaded from the GSEA/MSigDB database (<http://software.broadinstitute.org/gsea/msigdb>). In addition, the clinical information of the breast cancer patients was downloaded from the TCGA database, including TNM stage, prognosis survival time and other information. The 290 normal breast tissues were used as reference samples. The gene expression data sets of normal and cancer samples were both converted to the TPM form and contain 18006 genes in total. The independent validation data of the prognosis model was from the GSE3494 set in GEO Datasets, which contains 251 expression profiles of breast tumors by array.

Construction of individual-level networks

In this study, gene-gene interactions with high confidence (comprehensive score >0.9) were selected from the STRING database, which includes regulatory, physical and co-expression

protein-protein interaction networks. Furthermore, the above gene-gene interactions with both genes in one of the 186 KEGG pathways were used as the background network (or template network). In short, the individual-level networks to be constructed are actually networks with significant perturbations of co-expression.

In the following analysis, individual-level networks for breast cancer samples were constructed based on gene expression profiles by using the method introduced in reference [7] (see Fig.1A). First, using the gene expression data of n reference samples, the reference network can be constructed by calculating the correlation coefficient PCC_n (the Pearson correlation coefficient (PCC)) of the gene pairs connected in the background network. The weights of the edges in the reference network are the PCC of the corresponding gene pairs. Then, the expression data of an individual breast cancer sample was added to the reference samples, and the perturbed network of the single sample was constructed by calculating the new correlation coefficient PCC_{n+1} of the gene pairs in the background network. For the single breast cancer sample, the differential correlation coefficients of each edge between the perturbed network and the reference network were calculated: $\Delta PCC_n = PCC_{n+1} - PCC_n$. In reference [7], Liu et al. have proved that ΔPCC_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 large enough. The significance level of each ΔPCC was determined by the Z-test. The statistical Z-value is calculated as follows with the null hypothesis that ΔPCC_n is equal to 0:

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

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

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

Only the gene pairs that are perturbed significantly in most breast cancer samples are considered to be related to breast cancer. Then, the edges that are perturbed significantly in more than 90% of the samples by the binomial right-sided test (P -value < 0.05) constitute a gene-gene interaction network related to breast cancer. A slight change in the expression of high-degree genes in the network may cause disturbances of the entire network. Thus, these genes with high degree are considered to be the key genes for the onset and development of breast cancer. We selected genes with degrees > 5 in the identified breast cancer-related network for the subsequent enrichment analysis. Furthermore, we identified the specific networks and key genes related to each TNM stage and PAM50 subtype with the same method. Here, because of the small number of stage V samples, stage V was combined with stage IV.

Pathway enrichment analysis

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

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

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

Survival analysis by the Cox regression model

According to the clinical data of patients with breast cancer, we utilized the “survival” package and “survminer” package in R/Bioconductor to establish a univariate Cox proportional hazards regression model by setting patient survival conditions (survival time and survival status) as the dependent variables and the gene pairs connected in the background network as the covariates. Gene pairs with P -values < 0.05 were considered to be related to the prognosis of breast cancer [9].

A large number of covariates may cause overfitting in establishing a multivariate Cox proportional hazards regression model; thus, using the least absolute shrinkage and selection operator (LASSO), we further selected the key gene pairs from these significant ones obtained by the univariate Cox proportional hazards analysis. LASSO is a common method used in high-dimensional data regression, which can select prognosis-related gene pairs of breast cancer by shrinking regression coefficients. The tuning parameter (λ) with the smallest mean-square error was selected by four-fold cross-validation to establish an optimal LASSO regression model. Then, the coefficients of most gene pairs reduced to zero, and smaller number of gene pairs with nonzero coefficients were considered to be closely correlated with the prognosis of breast cancer.

LASSO Cox analysis was performed by using the “glmnet” package in R. Then, the risk score for each sample was calculated by the LASSO Cox regression model. According to the median risk score, breast cancer patients were divided into two groups (a high-risk group and a low-risk group). In addition, 234 breast tumors with relapse free survival information in the validation data set were analyzed used the above individual-level network method, and they were also divided into two groups in the same way. Finally, the corresponding Kaplan-Meier survival curves were plotted by using the packages “survminer” and “survival” in R.

Results

Breast cancer-related gene interaction networks

The background network consisted of 46916 edges and 3237 genes. In addition, 2190 gene pairs were identified as significantly related to breast cancer, which constituted the gene-gene

interaction network related to breast cancer (including 915 genes in total). We use the Cytoscape software to visualize the breast cancer-related network (see Fig.2).

Genes with degrees > 5 in the breast cancer-related network (198 in total which are shown in Table S1). Among them, some genes with higher degrees (> 20) have been shown to be related to breast cancer. For example, CCNB1 has strong power to predict the survival of breast cancer patients with the phenotype of ER positive [11]. The overexpression of GRB2 has been demonstrated to be significantly associated with the occurrence and poor prognosis of breast cancer [12]. PCNA has been proven to be a marker of proliferation in the diagnosis of breast cancer [13], SF3B4 has been shown to be a tumour suppressor, and somatic inactivating mutations occasionally occur in breast cancer [14]. UBE2C may promote the development of breast cancer [15]. High Cdc20 and securin immune expression are associated with extremely poor outcomes in breast cancer patients [16], and overexpression of RPL17 affects breast cancer-associated brain metastases [17]. MAD2L1 may have great effect on breast cancer progression, and its expression might help to predicting breast cancer prognosis [18]. The high expression of TRA2B is closely related to the cancer cell survival and therapeutic sensitivity of breast cancer [19]. GTF2H4 has been identified to be related to the survival risk of breast cancer [20].

Stage-related gene interaction networks

The results of the four stage-related gene-gene interaction networks are shown in Fig.S1A-D. There are obvious similarities and differences among the four stage-related gene interaction networks. There are 81 key genes shared by all stages (see Table S3), among which RPL17, CCNB1, and SF3B4 are genes that are highly (with degrees > 25) related to breast cancer. Stage I has 5 specific genes: PSMC5, SDHB, RPL11, SDHA, and RPL13. Stage II has 3 specific genes: STX6, CCNA2, and CDC25C. Stage III has 4 specific genes: NDUFA6, EPN1, SF3A3, and LSM7. Stage IV has the largest number of specific genes, with a total of 38, among which CDC42, LSM2, NDUFS6, and CDC25A are strongly associated with it. And these stage-specific key genes are shown in Table S4.

Subtype-related gene interaction networks

The results of the four subtype-related gene interaction networks are shown in Fig.S1E-G. The four subtype-related networks share similar and different characteristics. There are 34 key genes shared by the four subtypes (see Table S7). Among them, RPL17 and CCNB1 have higher degrees. The Luminal A subtype has 11 specific genes, including RPL23A, RPL10 and PRPF6, which are greatly related to it with higher degrees. The Luminal B subtype has 17 specific genes, including COX6C, EGFR and CLTC, which are related to it with higher degrees. The Her2 subtype has 3 specific genes, NDUFA6, CCR8, and CASP3. The Basal-like subtype has the largest number of specific genes, 17 in total, including LSM2, DDX5, SF3A3 and MAGOH, with higher degrees. And these subtype-specific key genes are shown in Table S8.

Pathways enriched in breast cancer patients

There were 41 pathways (see Table S2) enriched in the breast cancer samples according to the pathway enrichment analysis, including some immune-related pathways, such as the Toll-like receptor signaling pathway, antigen processing and presentation, complement and coagulation cascades, the [RIG-I-like receptor signaling pathway](#), and the cytosolic DNA-sensing pathway. Some important signal transduction and signal molecular interaction pathways were also included, such as the MAPK signaling pathway, Wnt signaling pathway, cytokine-cytokine receptor interaction, and ECM-receptor interaction pathways. Breast cancer is closely related to endocrine disorders [21], and two endocrine-related pathways, adipocytokine signaling pathway and PPAR signaling pathway, have also been identified as being related to breast cancer. In addition, some metabolic pathways, especially lipid metabolism pathways, have also been identified as being associated with breast cancer [22], such as the steroid hormone biosynthesis, arachidonic acid metabolism, arginine and proline metabolism pathway, and glycerolipid metabolism. Additionally, pathways in cancer was also enriched. The enrichment results are shown in Fig.3A.

Most of these pathways have been documented to be related to breast cancer. For example, the dysregulation of the steroid hormone biosynthesis pathway may affect steroid hormone levels and may thus be related to the susceptibility to breast cancer [21]. The PPAR signaling pathway may play an important role in the neoadjuvant chemotherapy response of breast cancer [23]. Mounting preclinical evidence supports targeting the MAPK signaling pathway in the triple negative breast cancer (TNBC) [24]. AMPK activators inhibit breast cancer cell proliferation by inhibiting DVL3-promoted Wnt/ β -catenin signaling pathway activity [25]. Toll-like receptors may play dual roles in human cancers [26]. The co-activation of the Hedgehog and Wnt signaling pathways is a poor prognostic marker in TNBC [27]. Prl-3 is closely related to cell migration and invasion in TNBC [28]. The YHD inhibition of 4T1 breast tumour growth may be related to the negative regulation of the JAK/STAT3 pathway by repressing the expression of IL-6 and TGF- β [29].

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 pathway, regulation of actin cytoskeleton, calcium signaling pathway, pathways in cancer, and cell adhesion molecules, which have been shown to have a high correlation with breast cancer [24, 25, 30-32]. 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 pathway, 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 [23, 33-36].

As we expected, stage IV was specifically enriched the most pathways (18 in total, see Table S6) different from the 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.S2). A multivariate Cox proportional hazards regression model with these gene pairs as independent variables was constructed as follows.

$$\text{Score} = 206.3 * (\text{ENO}, \text{PGK2}) + 35.9 * (\text{ENO}, \text{PKLR}) + 4.1 * (\text{EBP}, \text{HSD17B7}) + 5.5 * (\text{CYP1B}, \text{HSD17B1}) - 3.4 * (\text{NDUFB2}, \text{NDUFB4}) - 0.6 * (\text{ATP6V1A}, \text{ATP6V1B1}) + \dots$$

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

Discussion

At present, research on cancer pathology is limited to gene expression and mutation information. However, the model of one gene to one disease is no longer suitable for the study of complex diseases. In fact, genes do not exist in isolation but participate in some complex biological networks, such as gene-gene interaction networks. Gene mutations or surroundings changes often affect the balance of gene interaction networks and the perturbation of the networks then affect the onset and development of complex diseases. Therefore, network analysis can provide a more

comprehensive and systematic point of view, to better understand the human disease onset and development mechanism.

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

In this paper, individual-level networks of breast cancer samples were established to explore the genes-gene interaction networks related to the stages and PAM50 subtypes of breast cancer. Then, the pathways related to breast cancer were identified by hypergeometric testing. Through the same method, we also obtained the stage-related pathways and subtype-related pathways. Finally, the edge biomarkers (gene pairs) that are closely related to the prognosis of breast cancer were determined by using the LASSO regression model, and then a more stable prognostic analysis model was established by using these biomarkers. Our results indicate that the prognosis model has the robust and strong generalization capability, and it can be used in different gene expression data sets.

Many studies have shown that network-based methods are more robust and effective than single-gene-based methods. The advantages of network-based methods have been well documented and accepted in the analysis of noisy high-throughput data. Different from the usual network-based method, we made better use of the gene interaction relations not only the gene sets in the background network to explore the individual-level perturbation networks. This work helps us to better understand the heterogeneity and mechanism of breast cancer from an individual-level network angle. Precision medicine advocates the development of individualized treatment according to the unique features of the patients. Therefore, identifying the unique pathogeny embedded in each patient is important to develop a treatment strategy for each patient. Our individual-level network analysis of breast cancer will promote the development of precision medicine.

Acknowledgements

We are very grateful to Prof. Zhaohui Qin for his useful advice.

Additional information and declarations

Funding: This work was supported by the China Postdoctoral Science Foundation (No. 2019M651658), and the National college students' innovation training program (No. 201910307068Z).

Author contributions: All authors discussed the results and commented on the manuscript.

Conflicts of Interest: The authors declare no conflicts of interest.

Supplementary Materials: Supplementary information for this article are available in the supplemental files. **Table S1:** Key genes of BRCA (the genes with degrees > 5 in the breast cancer-related gene interaction network). **Table S2:** Pathways enriched in BRCA (the pathways with P -value < 0.05 by the enrichment analysis for all breast cancer samples). **Table S3:** Key genes shared by four stages. **Table S4:** Stage-specific key genes. **Table S5:** Enriched pathways shared by four stages. **Table S6:** Stage-specific enriched pathways. **Table S7:** Key genes shared by four subtypes. **Table S8:** Subtype-specific key genes. **Table S9:** Enriched pathways shared by four subtypes. **Table S10:** Subtype-specific enriched pathways. **Figure S1:** Gene interaction networks related to breast cancer TNM stages and PAM50 subtypes. **Figure S2:** Establishment of the LASSO regression model.

References

1. Bray F, et al. Global cancer statistics 2018: GL OBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA Cancer J Clin*, 2018. 68(6):394- 424.
2. Perou CM, Sorlie T, Eisen MB, et al. Molecular portraits of human breast tumors. *Nature*, 2000. 406(6797): 747-752.
3. Hammond ME, Hayes DF, Dowsett M, et al. American Society of Clinical Oncology/College of American Pathologists guideline recommendations for immunohistochemical testing of estrogen and progesterone receptors in breast cancer. *J Clin Oncol*, 2010.28(16):2784.
4. Fitzgibbons PL, Murphy DA, Hammond MEH, et al. Recommendations for validating estrogen and progesterone receptor immunohistochemistry assays. *Arch Pathol Lab Med*, 2010. 134(6):930.
5. Cheang MCU, Chia SK, David V, et al. Ki-67 index, HER2 status, and prognosis of patients with luminal B breast cancer. *J Natl Cancer Inst*, 2009. 101(10):736-750.
6. Goldhirsch A, Wood WC, Coates AS, et al. Strategies for subtypes-dealing with the diversity of breast cancer: highlights of the St.Gallen International Expert Consensus on the Primary Therapy of Early Breast Cancer. *Ann Oncol*, 2011. 22(8):1736-1747.
7. Liu XP, Chen LN, et al. [Personalized characterization of diseases using individual-level networks](#). *Nucleic Acids Research*, 2016. 44(22):e164.
8. Joel SP, Michael M. Supervised Risk Predictor of Breast Cancer Based on Intrinsic Subtypes. *J Clin Oncol*, 2009. 27(8): 1160–1167.
9. Zou H, Wanggou SY, Ye NR, Li YW, Huang Q, Liu HW, Xiong ZJ, Li XJ. Survival risk prediction model for patients with glioma. *Journal of International Neurology and Neurosurgery*, 2019. 46(1): 1-6.
10. Cheng P. A prognostic 3-long noncoding RNA signature for patients with gastric cancer. *J Cell Biochem*. 2018. 119(11):9261-9269.
11. Ding K, Li W, Zou Z, Zou X, Wang C2. CCNB1 is a prognostic biomarker for ER+ breast cancer. *Med Hypotheses*. 2014. 83(3):359-64.

12. Zhang Y, Xu G, Liu G, Ye Y, Zhang C, Fan C, Wang H, Cai H, Xiao R, Huang Z, Luo Q. miR-411-5p inhibits proliferation and metastasis of breast cancer cell via targeting GRB2. *Biochem Biophys Res Commun*. 2016. 476(4):607-613.
13. Juríková M, et.al. Ki67, PCNA, and MCM proteins: Markers of proliferation in the diagnosis of breast cancer. *Acta Histochem*. 2016. 118(5):544-52.
14. Denu RA, Burkard ME. Synchronous Bilateral Breast Cancer in a Patient With Nager Syndrome. *Clin Breast Cancer*. 2017. 17(3):151-153.
15. Mo CH, Gao L, Zhu XF, Wei KL, Zeng JJ, Chen G, Feng ZB. The clinicopathological significance of UBE2C in breast cancer: a study based on immunohistochemistry, microarray and RNA-sequencing data. *Cancer Cell Int*. 2017. 17:83.
16. Karra H, Repo H, et.al. Cdc20 and securin overexpression predict short-term breast cancer survival. *Br J Cancer*. 2014. 110(12):2905-13.
17. Yuan F, Wang W, Cheng H. Co-expression network analysis of gene expression profiles of HER2⁺ breast cancer-associated brain metastasis. *Oncol Lett*. 2018. 16(6):7008-7019.
18. Wang Z, Katsaros D, Shen Y, Fu Y, Canuto EM, Benedetto C, Lu L, Chu WM, Risch HA, Yu H. Biological and Clinical Significance of MAD2L1 and BUB1, Genes Frequently Appearing in Expression Signatures for Breast Cancer Prognosis. *PLoS One*. 2015. 10(8):e0136246.
19. Best A, Dagliesh C, Ehrmann I, Kheirollahi-Kouhestani M, Tyson-Capper A, Elliott DJ. Expression of Tra2 β in Cancer Cells as a Potential Contributory Factor to Neoplasia and Metastasis. *Int J Cell Biol*. 2013. 2013:843781.
20. Ge J, Liu H, Qian D, Wang X, Moorman PG, Luo S, Hwang S, Wei Q. Genetic variants of genes in the NER pathway associated with risk of breast cancer: A large-scale analysis of 14 published GWAS datasets in the DRIVE study. *Int J Cancer*. 2019. 145(5):1270-1279.
21. Sakoda LC, Blackston C, Doherty JA, Ray RM, Lin MG, Stalsberg H, Gao DL, Feng Z, Thomas DB, Chen C. Polymorphisms in steroid hormone biosynthesis genes and risk of breast cancer and fibrocystic breast conditions in Chinese women. *Cancer Epidemiol Biomarkers Prev*. 2008. 17(5):1066-73.
22. Merdad A, Karim S, Schulten HJ, Jayapal M, Dallol A, Buhmeida A, Al-Thubaity F, et al. Transcriptomics profiling study of breast cancer from Kingdom of Saudi Arabia revealed altered expression of Adiponectin and Fatty Acid Binding Protein4: Is lipid metabolism associated with breast cancer? *BMC Genomics*. 2015. 16 Suppl 1:S11.
23. Chen YZ, Xue JY, Chen CM, Yang BL, Xu QH, Wu F, Liu F, Ye X, Meng X, Liu GY, Shen ZZ, Shao ZM, Wu J. PPAR signaling pathway may be an important predictor of breast cancer response to neoadjuvant chemotherapy. *Cancer Chemother Pharmacol*. 2012. 70(5):637-44.
24. Giltane JM, Balko JM. Rationale for targeting the Ras/MAPK pathway in triple-negative breast cancer. *Discov Med*. 2014. 17(95):275-83.
25. Zou YF, Xie CW, Yang SX, Xiong JP. AMPK activators suppress breast cancer cell growth by inhibiting DVL3-facilitated Wnt/ β -catenin signaling pathway activity. *Mol Med Rep*. 2017. 15(2):899-907.

26. [Khademalhosseini M, Arababadi MK](#). Toll-like receptor 4 and breast cancer: an updated systematic review. *Breast Cancer*. 2019. 26(3):265-271.
27. [Bhateja P, Cherian M, Majumder S, Ramaswamy B](#). The Hedgehog Signaling Pathway: A Viable Target in Breast Cancer? *Cancers (Basel)*. 2019. 11(8):1126.
28. [Gari HH, DeGala GD, Ray R, Lucia MS, Lambert JR](#). PRL-3 engages the focal adhesion pathway in triple-negative breast cancer cells to alter actin structure and substrate adhesion properties critical for cell migration and invasion. *Cancer Lett*. 2016. 380(2):505-12.
29. [Mao D, Feng L, Gong H](#). The Antitumor and Immunomodulatory Effect of Yanghe Decoction in Breast Cancer Is Related to the Modulation of the JAK/STAT Signaling Pathway. *Evid Based Complement Alternat Med*. 2018. 2018:8460526.
30. [Kazazian K, Go C, Wu H, Brashavitskaya O, Xu R, Dennis JW, Gingras AC, Swallow CJ](#). Plk4 Promotes Cancer Invasion and Metastasis through Arp2/3 Complex Regulation of the Actin Cytoskeleton. *Cancer Res*. 2017. 77(2):434-447.
31. [Woltmann A, Chen B, Lascorz J, Johansson R, Eyfjörd JE, Hamann U, Manjer J, Enquist-Olsson K, Henriksson R, Herms S, Hoffmann P, Hemminki K, Lenner P, Försti A](#). Systematic pathway enrichment analysis of a genome-wide association study on breast cancer survival reveals an influence of genes involved in cell adhesion and calcium signaling on the patients' clinical outcome. *PLoS One*. 2014. 9(6): e98229.
32. [Saadatmand S, de Kruijf EM, Sajet A, Dekker-Ensink NG, van Nes JG, Putter H, Smit VT, van de Velde CJ, Liefers GJ, Kuppen PJ](#). Expression of cell adhesion molecules and prognosis in breast cancer. *Br J Surg*. 2013. 100(2):252-60.
33. [Bao Y, Wang L, Shi L, Yun F, Liu X, Chen Y, Chen C, Ren Y, Jia Y](#). Transcriptome profiling revealed multiple genes and ECM-receptor interaction pathways that may be associated with breast cancer. *Cell Mol Biol Lett*. 2019. 24:38.
34. [Yang Y, Liu L, Fang M, Bai H, Xu Y](#). The chromatin remodeling protein BRM regulates the transcription of tight junction proteins: Implication in breast cancer metastasis. *Biochim Biophys Acta Gene Regul Mech*. 2019. 1862(5):547-556.
35. [Tang X, Shi L, Xie N, Liu Z, Qian M, Meng F, Xu Q, Zhou M, Cao X, Zhu WG, Liu B](#). SIRT7 antagonizes TGF- β signaling and inhibits breast cancer metastasis. *Nat Commun*. 2017. 8(1):318.
36. [Peng L, Hu Y, Chen D, Linghu R, Wang Y, Kou X, Yang J, Jiao S](#). Ubiquitin specific protease 21 upregulation in breast cancer promotes cell tumorigenic capability and is associated with the NOD-like receptor signaling pathway. *Oncol Lett*. 2016. 12(6):4531-4537.

Figure Legends

Figure 1

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

(A) Construction of individual-level 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 individual-level network for sample sd_x is composed by those edges with significant ΔPCC_n .

(B) The framework to identify the breast cancer-related gene interaction network based on gene expression. Using the individual-level network analysis method, m cancer individual-level networks were constructed. Then, these constructed individual-level 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 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 $-\log_{10}(p)$.
- (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 interaction networks related to breast cancer stages and 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 ≥ 15 , and the blue ones are the genes with degrees < 15 .

- (A-D) Gene interaction networks associated with Stage I, II, III, and IV respectively.
- (E-H) Gene interaction networks related to the LumA, LumB, Her2, and Basal-like subtypes respectively.

Figure S2. Establishment of the LASSO regression model.

- (A) Four-fold cross-validation for tuning parameter (λ) selection in the LASSO model.
- (B) LASSO coefficient profiles of 272 gene interactions.

Figure 1

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

(A) Construction of individual-level 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 individual-level network for sample sd_x is composed by those edges with significant ΔPCC_n . **(B)** The framework to identify the breast cancer-related gene interaction network based on gene expression. Using the individual-level network analysis method, m cancer individual-level networks were constructed. Then, these constructed individual-level 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.

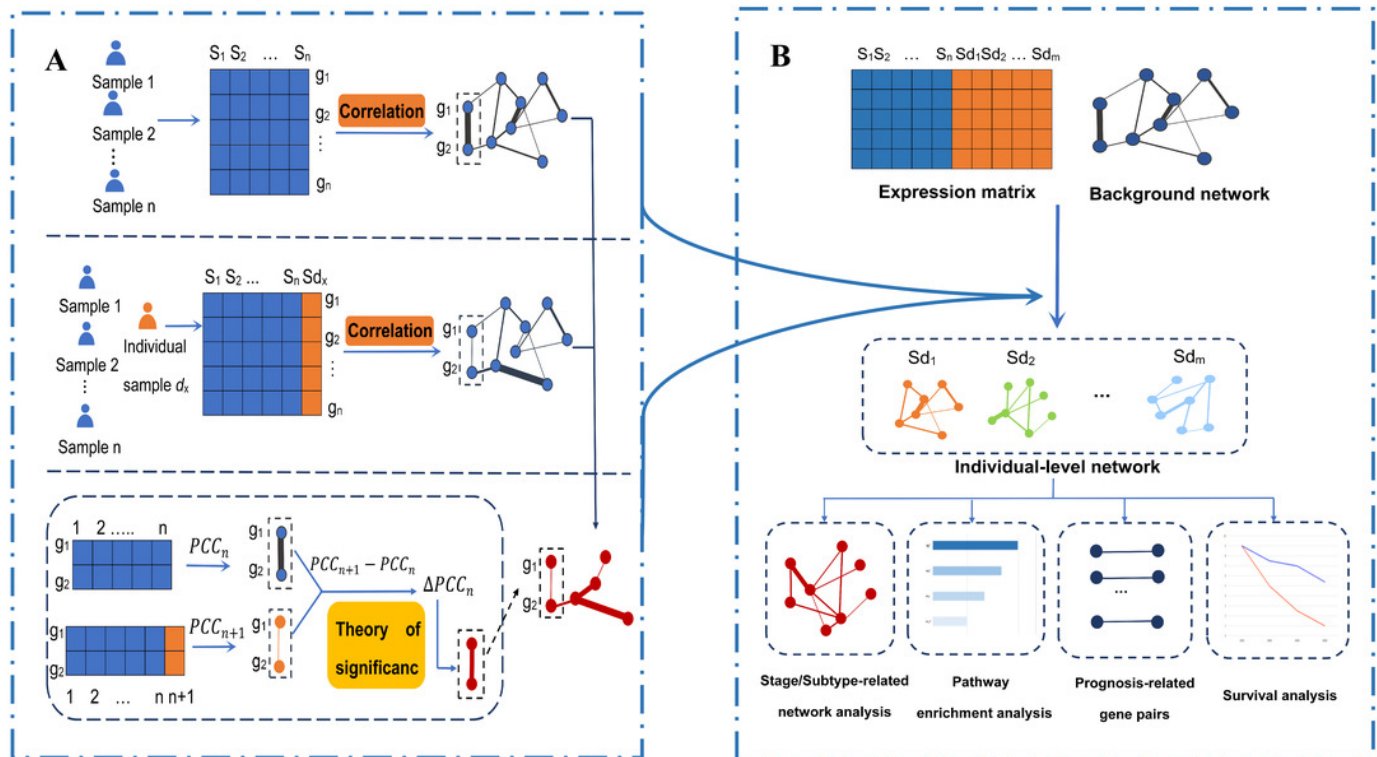
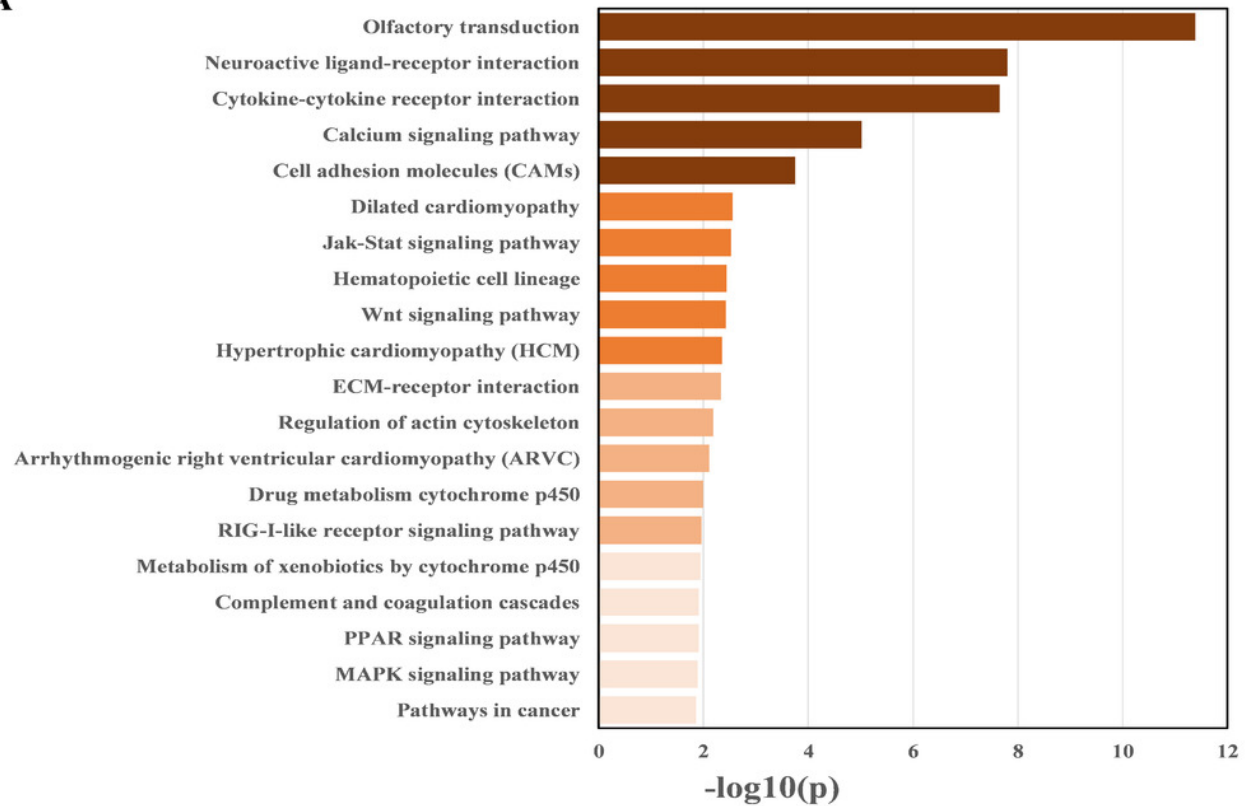


Figure 3

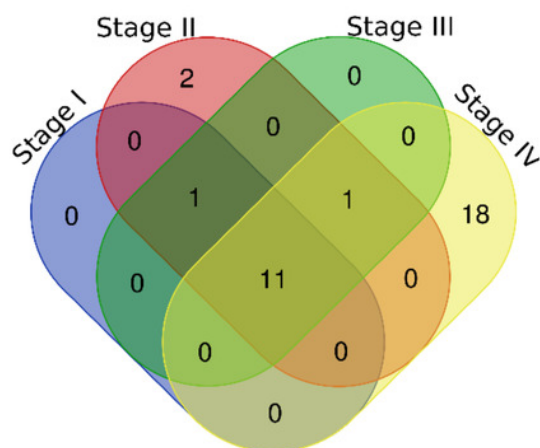
Pathways enriched in breast cancer, as well as different stages and subtypes of it.

(A) KEGG pathways enriched in breast cancer samples, ranked by $-\log_{10}(p)$. **(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.

A



B



C

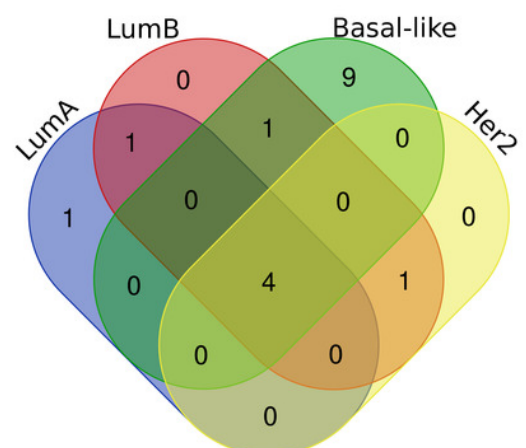


Figure 4

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.

