# Network-based transfer of pan-cancer immunotherapy responses to guide breast cancer prognosis

Check for updates

Xiaobao Ding[1,2,3], Lin Zhang[1], Ming Fan[1] ✉ & Lihua Li [1,3] ✉

Breast cancer prognosis is complicated by tumor heterogeneity. Traditional methods focus on cancer-specific gene signatures, but cross-cancer strategies that provide deeper insights into tumor homogeneity are rarely used. Immunotherapy, particularly immune checkpoint inhibitors, results from variable responses across cancers, offering valuable prognostic insights. We introduced a network-based transfer (NBT) of pan-cancer immunotherapy responses to enhance breast cancer prognosis using node embedding and heat diffusion algorithms, identifying gene signatures netNE and netHD. Our results showed that netHD and netNE outperformed seven established breast cancer signatures in prognostic metrics, with netHD excelling. All nine gene signatures were grouped into three clusters, with netHD and netNE enriching the immune-related interferon-gamma pathway. Stratifying TCGA patients into two groups based on netHD revealed significant immunological differences and variations in 20 of 50 cancer hallmarks, emphasizing immune-related markers. This approach leverages pan-cancer insights to enhance breast cancer prognosis, facilitating insight transfer and improving tumor homogeneity understanding.

Breast cancer is the most common malignancy among women and remains the leading cause of tumor-related mortality worldwide[1], despite significant advances in treatment strategies. This disease is notably heterogeneous, with each patient potentially having a different prognosis[2]. Although many great efforts have been made to develop prognostic models for clinical decision-making, accurately predicting the prognosis remains a substantial challenge for patients with breast cancer.

To address this challenge, clinicians and researchers often consider several clinical and pathological features, including tumor size, lymph node status, and histological grade, when assessing prognosis[3,4]. However, these clinicopathological factors are insufficient for accurately predicting the outcomes of patients with breast cancer. Recognizing that tumor heterogeneity originates at the molecular level[5], molecular examinations of key biomarkers, including estrogen receptor (ER)[6], progesterone receptor (PR)[7], HER2[8], and Ki-67[9], are routinely performed to enhance the accuracy of prognostic assessments[10]. This approach underscores the critical role these molecular markers play in enabling more precise evaluations of breast cancer prognosis[11].

In recent years, the advancement of high-throughput molecular profiling technologies has expanded molecular examinations and significantly facilitated the construction of cancer prognosis models[11,12]. Therefore, gene signatures, which consist of multiple genes, are now widely used to predict breast cancer prognosis[13,14]. Researchers can identify these gene signatures through biological experiments or computational methods, both of which are applied within the framework of transcriptomics. For instance, the LM gene signature[15], which includes 54 genes associated with lung metastasis in breast cancer, was identified from a biological experiment using transcriptomic analysis of cell lines. Compared to biological experiments, computational methods have become prevalent due to their convenience in identifying gene signatures[16]. Numerous gene signatures have been developed, including the PAM50[17], which builds upon an expanded set of intrinsic genes identified in earlier studies to select those crucial for differentiating among five intrinsic breast cancer subtypes. Endo[18] conducted univariate Cox regression and ultimately selected eight genes of interest. Similarly, the RS[19] filters 16 cancer markers from a pool of 250 candidate genes. Additionally, Mamma[20] and GGI97[21] leveraged statistical analysis to pinpoint genes with differential expression between two distinct tumor types. Building on these advancements, single-cell RNA sequencing (scRNA-seq) has recently become a key tool for studying tumor heterogeneity at the cellular level[22]. Using scRNA-seq data, researchers have identified the scP.W[23] gene signature by exploring phenotypic variations associated with epithelial-mesenchymal transition (EMT) in tumor cells.

[1]Institute of Biomedical Engineering and Instrumentation, Hangzhou Dianzi University, Hangzhou, China. [2]Institute of Big Data and Artificial Intelligence in Medicine, School of Electronics and Information Engineering, Taizhou University, Taizhou, China. [3]School of Computer Science and Technology, Hangzhou Dianzi University, Hangzhou, China. ✉e-mail: ming.fan@hdu.edu.cn; lilh@hdu.edu.cn

These methods offer a variety of prognostic options for breast cancer, with several having received approval from the Food and Drug Administration for commercial use[24].

The aforementioned methods identify gene signatures specific to breast cancer prognosis and aim to bridge molecular and clinical phenotypes. These approaches focus primarily on breast cancer and associated genes. However, researchers are increasingly interested in molecular markers suitable for pan-cancer, as these may offer better insights into tumor homogeneity. These markers could also improve breast cancer prognosis. Immunotherapy, especially the use of immune checkpoint inhibitors (ICIs), is effective for pan-cancer treatment and yields distinct outcomes[25]. Patients treated with ICIs can exhibit four distinct response categories—complete response (CR), partial response (PR), stable disease (SD), and progressive disease (PD)—according to the modified RECIST (mRECIST) criteria[26,27]. Immunotherapy responses also reveal distinct prognoses that may advance cancer prognosis research.

Inspired by the success of immunotherapy across various cancers, which has led to distinct responses and prognoses in patients, we hypothesized that certain universal genes and their interactions can consistently predict cancer outcomes. These genes have the potential to guide breast cancer prognosis. However, transferring immunotherapeutic insights to breast cancer prognosis poses significant challenges. To address these challenges, we propose a theoretical model in which the biological network can bridge genes and clinical phenotypes. According to this model, genes and their interactions are pivotal in determining the ultimate clinical outcomes, influencing various clinical phenotypes across different cancer types. Our initial objective within this framework was to identify genes whose expressions anchor immunotherapy responses, which we term anchor genes. After identifying these anchor genes, we can utilize established protein–protein interaction (PPI) databases to gain insights into gene interactions (e.g., STRING[28]) and employ network analysis algorithms to discover additional genes functionally related to these anchor genes. Based on these efforts, we aimed to identify a clinical phenotype-specific gene set and improve the prognostic accuracy for breast cancer patients. The transfer of insights across clinical phenotypes, facilitated by biological networks, offers a new perspective for cancer research, extending beyond breast cancer prognosis.

## Results

### Overview of network-based breast cancer prognosis
We hypothesized that pan-cancer responses to immunotherapy are determined by specific genes and their interactions, which are applicable to many other clinical phenotypes across different cancer types. Therefore, we propose a network-based transfer (NBT) method utilizing immunotherapy responses to guide breast cancer prognosis. The initial step involved identifying clinical phenotype-specific genes, which we termed anchor genes (Fig. 1a and Supplementary Fig. S1). Specifically, the expression of these genes is directly anchored to immunotherapy responses. Subsequently, a PPI network was abstracted from STRING, with anchor genes mapped to this network (Fig. 1b). Here, we utilized two distinct network analysis algorithms—node embedding and heat diffusion—to generate a ranked gene list for gene signature identification. We then integrated multiple breast cancer cohorts and applied a user-defined greedy strategy to identify gene signatures (Fig. 1c). This process helped us generate distinct gene signatures: netNE, derived from node embedding, and netHD, derived from heat diffusion (Supplementary Table S1). Finally, we benchmarked all gene signatures derived from different methods by comparing them across various breast cancer datasets, including those from the TCGA, METBRIC, and GEO databases (Fig. 1d).

### Prognostic evaluation of pan-cancer immunotherapy cohorts
Immunotherapy represents a groundbreaking treatment approach that is effective for a wide range of cancer types. Despite its broad applicability, it leads to diverse prognostic outcomes among patients. Specifically, patients categorized under CR, PR, SD, and PD showed clear differences in prognosis

(Fig. 2a, b). In terms of overall survival or progression-free survival, patients who achieved a CR generally experienced the longest survival, indicating the most favorable prognosis, while those who achieved a PD exhibited the shortest survival, reflecting the least favorable outcomes. Patients with PR and SD fell into the intermediate prognostic category, with PR associated with slightly better outcomes than SD. These findings were consistent across the pan-cancer cohort. To precisely measure the concordance between survival time and immunotherapy responses, the concordance index (C-index) was calculated for each study. All datasets, except for Prins_2019, demonstrated a C-index above 0.7, indicating high concordance, with some studies even surpassing 0.8 (Fig. 2c, d). Immunotherapy elicits responses with distinct prognoses, offering valuable insights into breast cancer prognosis.

### The superior performance of network-based transfer method in risk prediction
To benchmark various methods for breast cancer prognosis, we collected signatures from seven different methods. Our network-based method employed two distinct strategies: netHD, which utilizes heat diffusion, and netNE, which is based on node embedding. The METABRIC dataset, which includes 1420 patients and is the largest cohort, was used to refine both the netHD and netNE. Ultimately, nine signatures were evaluated as risk prediction benchmarks. Both the METABRIC and TCGA datasets provided overall survival (OS) and relapse-free survival (RFS) data, whereas the GPL6098 dataset included only RFS data. We conducted 10-fold cross-validation on each dataset and aggregated the average concordance indices for each gene signature. Additionally, we ranked the gene signatures by their concordance indices and calculated the mean rank for each signature within each dataset (Table 1).

Notably, netHD and netNE demonstrated superior performance, with higher C-indices and lower mean rank values. Specifically, netHD and netNE achieved mean ranks of 1.2 and 1.6, respectively, with netHD performing the best. Following the network-based methods, the Endo method also demonstrated strong performance, particularly with the METABRIC and GEO datasets, where it performed nearly, as well as the network-based methods. However, its performance significantly decreased in the TCGA dataset, with the C-index decreasing to less than 0.6. Additionally, despite containing only eight genes, Endo outperformed larger signatures such as CGI97 (97 genes), Mamma (66 genes), and Pam50 (50 genes).
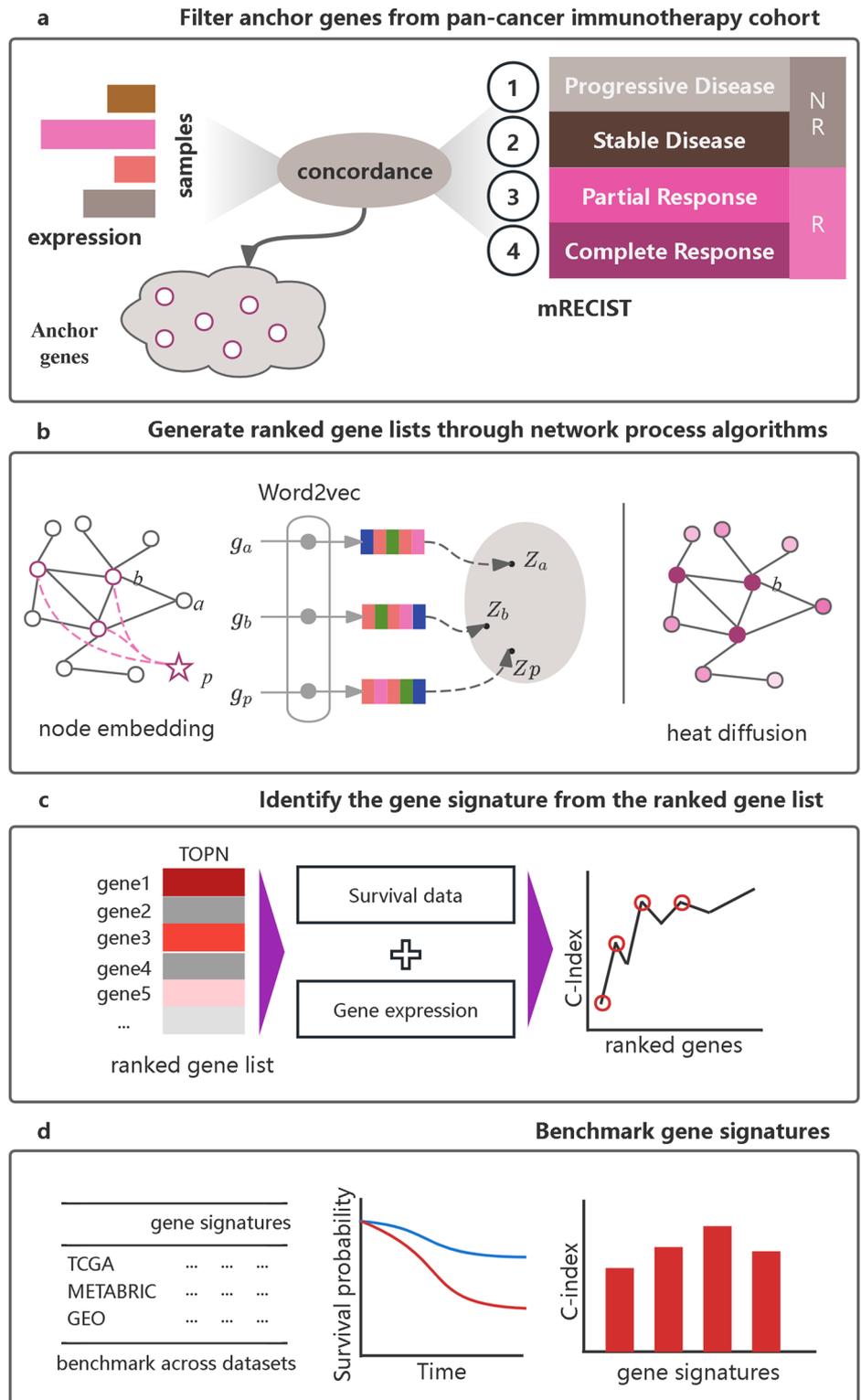
To further assess the robustness of our approach, we validated the above approach across ten cancer types in TCGA, including SKCM, STAD, and BLCA. Our method was applied in parallel across different cohorts, revealing variability in the C-indices of the prognostic models. Most C-indices exceeded 0.6, with some reaching as high as 0.7 for KIRC and even 0.8 for LGG (Supplementary Table S2). Relative to prognostic models for breast cancer, these results are robust and acceptable, suggesting that transferring immunotherapy responses from a pan-cancer context to inform specific cancer prognoses can effectively broaden our research perspective.

### The network-based transfer method performs well in the risk group
To further evaluate the risk stratification capability of the network-based method, we chose the hazard ratio (HR) as the criterion to benchmark it against seven other methods. For each method, we stratified patients into two groups based on the risk predictions calculated from the respective gene signatures. If a patient's risk score was greater than the median value, the patient was assigned to the high-risk group; otherwise, the patient was assigned to the low-risk group. HRs were calculated for all predictions generated by each method (Table 2). We observed that the two versions of the network-based method (netHD and netNE) perform well but not outstandingly. The Endo method performed the best, followed by the scP.W. Although the network-based methods did not achieve the best performance, they ranked reasonably well; specifically, they ranked 4th to 5th out of the nine methods according to the mean rank, which is considered acceptable.

**Fig. 1 | Overview of the NBT workflow.**
**a** Identification of anchor genes within a pan-cancer immunotherapy cohort. Genes with expression variations that align with immunotherapy responses, as defined by the modified response evaluation criteria in solid tumors (mRECIST), are selected as anchor genes. **b** Generation of ranked gene lists utilizes two distinct strategies: node embedding and heat diffusion. Anchor genes are mapped to the PPIs. In the node embedding strategy, a virtual phenotype node is introduced, with all anchor genes connected to it. The gene list is then ranked based on the similarity between the phenotype node and each gene node. Conversely, the heat diffusion strategy ranks the gene list according to the heat value of each node. **c** Gene signature identification: An improved greedy algorithm is used to identify a gene signature from the ranked gene list, utilizing gene expression and survival data. **d** Benchmark gene signatures: all gene signatures are benchmarked using the concordance index and survival metrics derived from Kaplan–Meier plots across various datasets.



**a** Filter anchor genes from pan-cancer immunotherapy cohort

**b** Generate ranked gene lists through network process algorithms

**c** Identify the gene signature from the ranked gene list

**d** Benchmark gene signatures

## Evaluation using an independent test

In the aforementioned benchmark analysis, we considered both the C-index and HR to evaluate the comprehensive performance of various methods. Among these, netHD, which utilizes a heat diffusion strategy, performed better. We further evaluated its performance with other methods using an independent test. The entire METBRIC dataset was used as the training dataset. All nine gene signatures were used to train prognostic models, which were then evaluated on the full TCGA dataset and across various cancer subtypes using the C-index. Subsequently, breast cancer patients from the TCGA dataset were divided into two groups according to their risk predictions. We then assessed the survival outcomes of these two groups.

The netHD achieved the best performance on both the TCGA (OS) and TCGA (RFS) datasets (Fig. 3a, b). Notably, for the TCGA (OS), the netHD method significantly outperformed the other methods, achieving a C-index of 0.72. Following closely, RS and scP.W also performed well, with concordance indices of 0.67 and 0.66, respectively, on the TCGA (OS), and
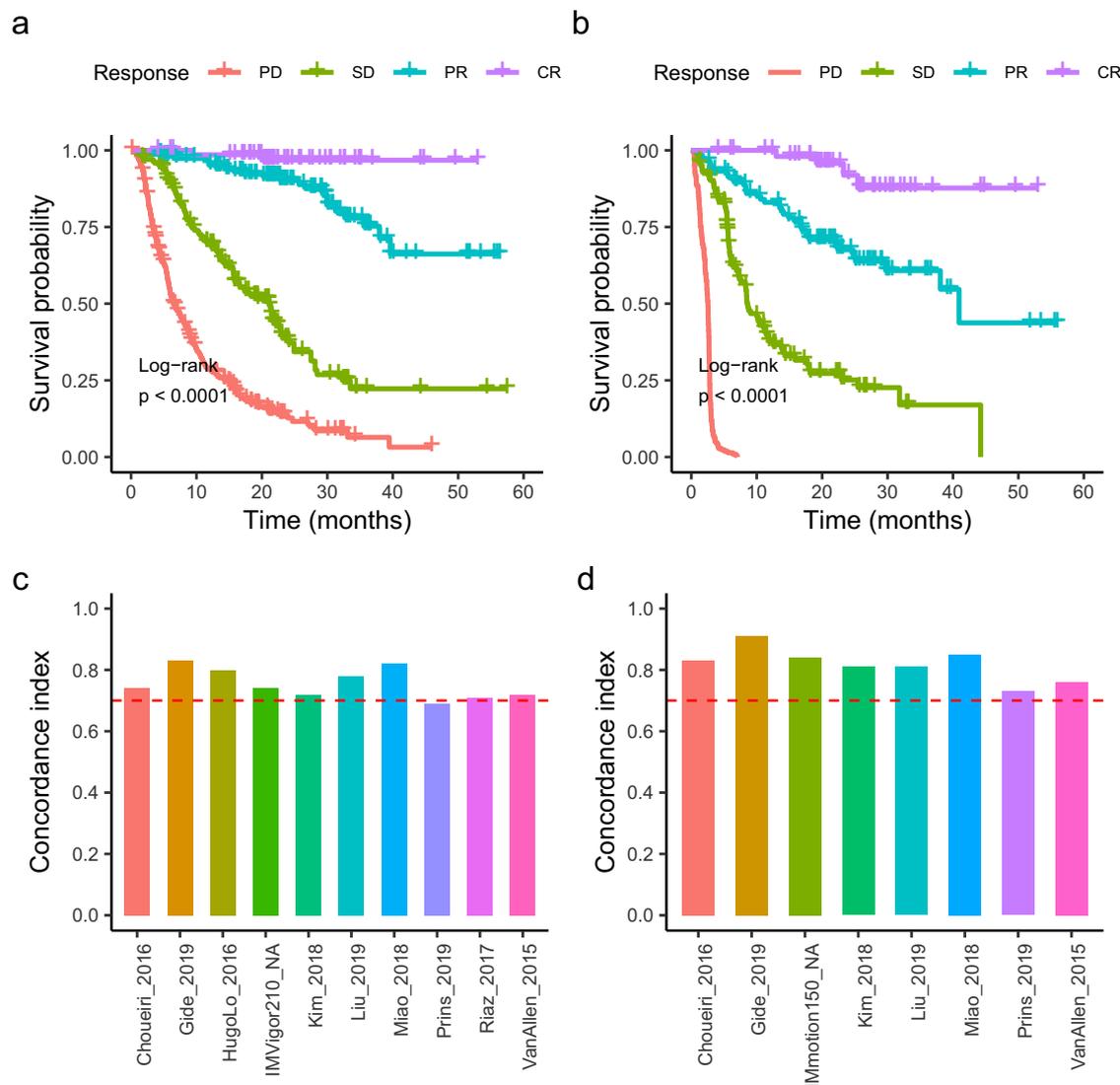
**Fig. 2 | Prognostic evaluation of the pan-cancer immunotherapy cohort.**
**a** Kaplan–Meier overall survival curves for the cohort, stratified by mRECIST response groups: complete response (CR), partial response (PR), stable disease (SD), and progressive disease (PD). **b** Kaplan–Meier progression-free survival curve for the cohort. **c** The concordance index for each dataset in the cohort, was calculated based on overall survival time. **d** Concordance index for each dataset in the cohort, was calculated based on progression-free survival time.

0.6 and 0.61 on the TCGA(RFS). Based on the risk scores, patients were categorized into high and low-risk groups, each showing distinct survival outcomes with a *p*-value of less than 0.0001(Fig. 3c, d).

Considering the heterogeneity of breast cancer, the test dataset was further categorized by ER, PR, and HER2 status (positive or negative) into four subtypes: ER+, PR+, HER2+, and triple-negative breast cancer (TNBC). We conducted evaluations on these subtypes, where netHD and netNE performed well compared to other models, with the netNE gene signature showing the best performance (Table 3).

**Functional co-occurrence analysis of gene signatures**
Although gene signatures identified by different methods often vary in quantity, they typically represent biological functions. Among the nine gene signatures, we aimed to further quantify their functional co-occurrence. To this end, we utilized the METABRIC and TCGA datasets to evaluate the functional co-occurrence of these gene signatures. Each sample from these two cohorts was assigned an enrichment score for each gene signature. Subsequently, we evaluated the functional correlations between gene signatures using these scores.

The results indicate that these gene signatures can be grouped into three clusters based on their functional correlations. Notably, netHD and netNE showed the highest similarity, with scores of 0.97 in the TCGA dataset (Fig. 4a) and 0.95 in the METABRIC dataset (Fig. 4b). However, these two signatures exhibit low similarity with others, except for the LM gene signature, which has a moderate similarity score of approximately 0.5 with both netHD and netNE. CGI97, scP.W, Mamma, and PAM50 formed a distinct cluster, demonstrating high similarity and significant functional co-occurrence. Furthermore, all these gene signatures were enriched in the biological process of cell cycle checkpoints (Supplementary Table S3). Similarly, Endo and RS showed a notable similarity, each achieving a score of 0.5 in both the TCGA and METABRIC datasets. This analysis of functional co-occurrence suggested that multiple functions likely contribute to breast cancer prognosis.

**Gene signature function and tumor microenvironment analysis**
Multiple functions derived from various gene signatures may collectively influence breast cancer prognosis. The netHD and netNE exhibit high similarity in terms of functional co-occurrence, with netHD

**Table 1 | Performance comparison for cancer prognosis using benchmark methods and network-based approaches (netHD and netNE)**

|  | Endo | GGI97 | LM | Mamma | Pam50 | RS | scP.W | netNE | netHD |
|---|---|---|---|---|---|---|---|---|---|
| TCGA (OS) | 0.58 | 0.52 | 0.62 | 0.50 | 0.52 | 0.60 | 0.57 | 0.64 | 0.68 |
| TCGA (RFS) | 0.59 | 0.50 | 0.57 | 0.54 | 0.50 | 0.57 | 0.61 | 0.59 | 0.61 |
| METABRIC (OS) | 0.65 | 0.62 | 0.64 | 0.63 | 0.64 | 0.64 | 0.63 | 0.65 | 0.65 |
| METABRIC (RFS) | 0.63 | 0.59 | 0.60 | 0.61 | 0.61 | 0.61 | 0.62 | 0.63 | 0.63 |
| GEO (GPL6098) | 0.65 | 0.60 | 0.61 | 0.63 | 0.57 | 0.62 | 0.63 | 0.66 | 0.65 |
| Mean rank | 2.4 | 8.2 | 5.4 | 6.4 | 6.6 | 4.8 | 4.4 | 1.6 | 1.2 |

**Table 2 | HR evaluation of the network-based method with seven benchmark methods**

|  | Endo | GGI97 | LM | Mamma | Pam50 | RS | scP.W | netHD | netNE |
|---|---|---|---|---|---|---|---|---|---|
| TCGA (OS) | 1.87 | 0.89 | 0.95 | 0.89 | 0.85 | 1.06 | 2.47 | 1.13 | 1.30 |
| TCGA (RFS) | 2.72 | 0.90 | 1.05 | 1.11 | 0.88 | 1.75 | 2.27 | 1.31 | 1.31 |
| METABRIC (OS) | 2.16 | 1.31 | 1.47 | 1.54 | 1.45 | 2.00 | 1.96 | 1.71 | 1.78 |
| METABRIC (RFS) | 2.17 | 1.24 | 1.58 | 1.57 | 1.54 | 2.02 | 2.11 | 1.74 | 1.84 |
| GEO (GPL6098) | 1.95 | 1.00 | 1.03 | 1.17 | 1.04 | 1.49 | 1.83 | 1.32 | 1.70 |
| Mean rank | 1.2 | 8.4 | 6.8 | 6.4 | 8.2 | 3.4 | 2 | 4.6 | 3.6 |

achieving the best performance according to the aforementioned benchmarks. We further analyzed the functions of netHD and examined the tumor microenvironment to gain deeper insights into breast cancer prognosis.

To thoroughly assess the function of netHD, we utilized two distinct databases for gene set enrichment analysis: the human-curated REACTOME[29] database and the computational database GO[30]. Our findings indicate that netHD predominantly captures immune-related annotation terms. Specifically, gene set enrichment analysis using REACTOME revealed that the top three enriched terms were interferon gamma signaling, interferon signaling, and chemokine receptor binding to chemokines (Fig. 5a). Similarly, analysis conducted with the GO database highlighted that the terms associated with cell–cell interactions, particularly regulation of cell-cell adhesion, regulation of leukocyte cell-cell adhesion, and leukocyte cell-cell adhesion, ranked among the top three (Fig. 5b). Furthermore, terms associated with interferon-gamma were among the top ten enriched terms in both the Reactome and GO databases.

The netHD gene signature, which represents immune-related functions, plays a crucial role in breast cancer prognosis. We delved deeper into the tumor microenvironment (TME) across various prognostic risk groups by examining the abundance of immune cells and the activity of cancer hallmarks. Specifically, CIBERSORT[31,32] was used for cell deconvolution, and the 50 gene signatures[33] associated with cancer hallmarks were used to evaluate cancer activity. Patients in the TCGA cohort were categorized into high-risk and low-risk groups based on their risk scores. Significantly, by setting the p-value cutoff below 0.0001, we observed distinct immunological profiles between the high-risk and low-risk breast cancer groups. The high-risk group demonstrated reduced levels of CD4 + T cells and CD8 + T cells but showed an increased abundance of M2 and M0 macrophages (Fig. 5c). Conversely, the low-risk group exhibited the opposite trend, underscoring the critical influence of these cells on breast cancer prognosis. Furthermore, with the same stringent p-value cutoff, notable differences in cancer hallmarks were evident between the two groups, with 20 out of 50 hallmarks showing significant variation (Fig. 5d). Among these, the INTERFERON_GAMMA_RESPONSE hallmark aligned consistently with gene set enrichment terms related to interferon-gamma. Additionally, the hallmarks IL2_STAT5_SIGNALING and IL6_JAK_STAT3_SIGNALING were identified as cytokine-mediated processes integral to the immune response, highlighting their potential roles in modulating cancer risk profiles.

## Discussion

To accurately predict breast cancer prognosis, we introduced a novel approach by transferring observations from pan-cancer immunotherapy responses to breast cancer prognosis. This approach employs a network-based method that utilizes both heat diffusion and node embedding to identify gene signatures. Compared with other established methods, these gene signatures have demonstrated promising performance in various evaluations.

Current methods for identifying gene signatures in breast cancer prognosis are cancer-specific and often miss insights into tumor homogeneity across multiple cancer types, an area of growing research interest. Our approach extends this focus to include pan-cancer immunotherapy responses. Existing methods, such as LM[15] for distant metastasis, GGI97[21] for a histologic grade, and Endo[18] for survival, primarily concentrate on the survival outcomes and malignant clinical phenotypes of breast cancer. Unlike methods that directly use clinical phenotypes, scP.W[23] focuses on the EMT phenotype of breast cancer cells, highlighting EMT as a key factor in tumor homogeneity and prognosis. However, this method is specifically applied to identify gene signatures within a breast cancer-specific single-cell atlas. Transferring clinical phenotypes across cancer types provides an alternative approach to breast cancer prognosis, reflecting tumor homogeneity through interferon gamma-related processes. interferon-gamma (IFNγ) plays a critical role in the anti-tumor immune response during immunotherapy. Immune checkpoint blockade has been shown to upregulate IFNγ, which in turn facilitates the clearance of tumor cells across various cancers[34]. Moreover, IFNγ has been identified as a prognostic marker in melanoma[35] and the majority of breast cancer cases[36]. This shared feature between immunotherapy and prognosis underscores the consistent role of IFNγ in influencing clinical outcomes.

The biological network serves as the bridge between genes and clinical phenotypes across various cancer types. Tumor heterogeneity often arises from the gap between genes and clinical phenotypes[37]. A biological network was introduced, revealing that only a minority of established clinical phenotype-specific genes can be extended to gain deeper insights[38]. Additionally, these insights can be transferred across cancer types to guide other clinical predictions. Based on this, the tumor homogeneity characteristics hidden in immunotherapy responses can be transferred to breast cancer prognosis. Similarly, the deep neural network framework has successfully made predictions from pixels to objects[39,40], effectively bridging the gap between lower pixel features and higher object features. Unlike deep neural
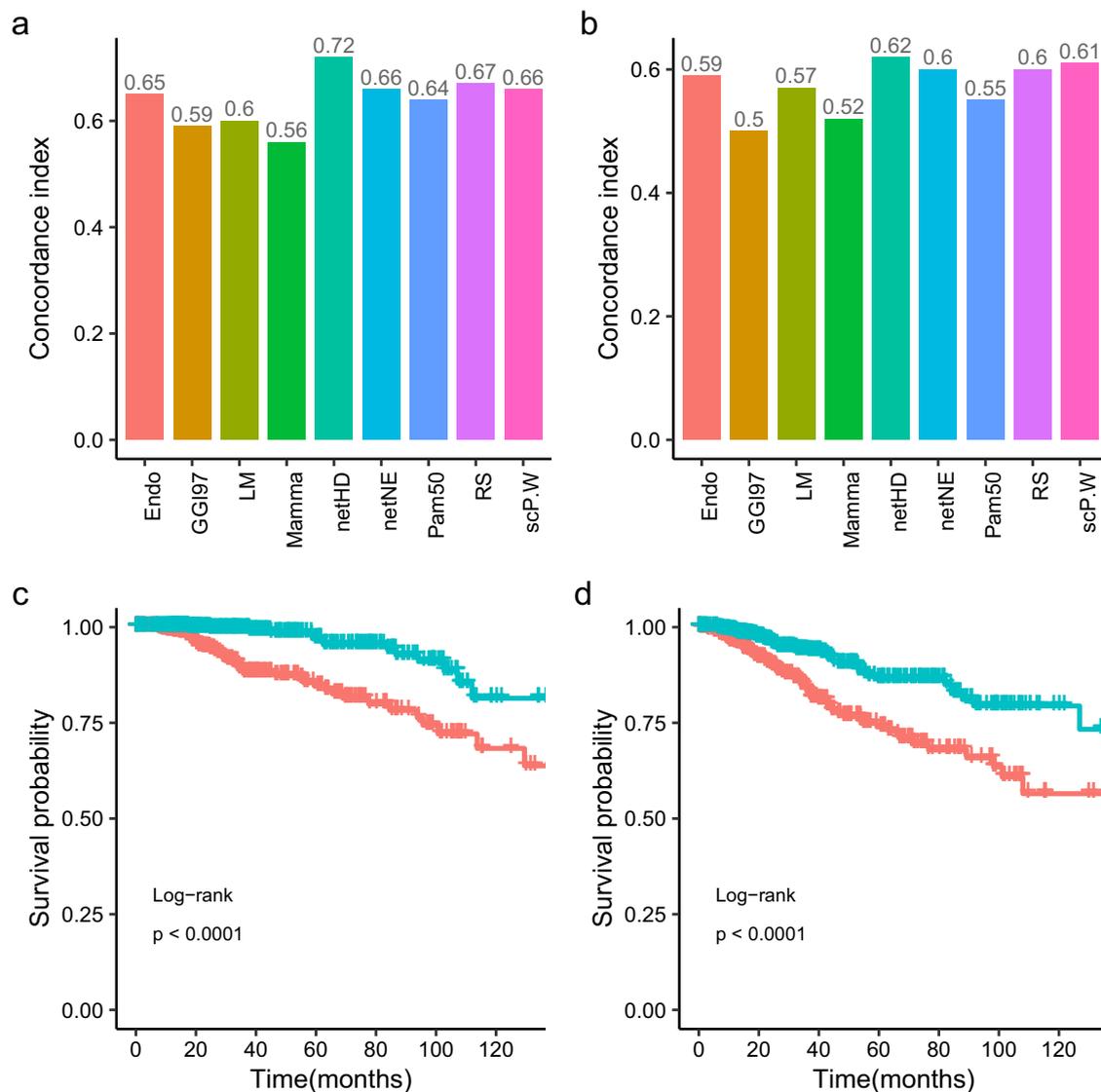
**Fig. 3 | Independent test for nine gene signatures. a** Concordance index in TCGA (OS): prognosis models developed with nine distinct gene signatures, initially trained on the METABRIC dataset, were evaluated on the TCGA dataset to assess their concordance indices for overall survival (OS). **b** Concordance index in TCGA (RFS): similar to (**a**), models were tested for relapse-free survival (RFS) in TCGA. **c** Survival analysis in TCGA (OS): Kaplan–Meier curves for overall survival, categorizing patients into prognostic risk groups based on netHD model predictions. **d** Survival analysis in TCGA (RFS): similar to (**c**), Kaplan–Meier curves for relapse-free survival.

networks, whose parameters require large amounts of data to learn, biological networks are inherently meaningful and incorporate a wealth of prior knowledge. This makes it feasible to robustly identify gene signatures for breast cancer prognosis.

Many biological processes can influence breast cancer prognosis, with immune-related biological processes being a significant factor in determining clinical phenotypes across cancer types. We collected nine gene signatures representing different functional co-occurrence patterns. Some of these are enriched in the cell cycle checkpoints biological pathway, indicating that the cancer cell proliferation phenotype influences breast cancer prognosis[41]. The gene signatures netHD and netNE, which we proposed, do not co-occur with others. They are enriched in interferon-gamma and chemokine biological processes. Notably, the shared genes within netHD and netNE include key markers associated with CD8 + T cells, underscoring the critical role of CD8 + T cells in tumor immunity. These genes not only indicate CD8 + T cell infiltration (e.g., CD8A, CXCL9[42], and CXCL10[43]) but also reflect CD8 + T cell activity, with positive regulators such as IFNG[44] and IRF1[45]. Interestingly, LAG3 appears

as an inhibitory marker[46]. Collectively, these genes provide a comprehensive view of the CD8 + T cell landscape within the tumor microenvironment across cancer types. CD8 + T cells are the most powerful effectors in the anticancer immune response and form the backbone of current successful cancer immunotherapies[47]. In terms of cancer prognosis, CD8 + T cell infiltration has also been identified as an independent prognostic marker for glioma[48], melanoma[49], kidney renal cell carcinoma[50], and the majority of breast cancer cases[51]. This finding further suggests that immune-related gene signatures provide valuable insights into breast cancer prognosis. Nevertheless, HER2+ and TNBC, as subtypes of breast cancer, do not fully support this conclusion, as they exhibit distinct immune-related pathway activities compared to other subtypes. This suggests the need for subtype-specific approaches when using immune-related gene signatures for prognosis prediction.

In summary, we adopted a new perspective on breast cancer prognosis by leveraging responses from pan-cancer immunotherapy, utilizing a network-based method to transfer clinical phenotypes across different cancer types. The gene signature we identified not only performed well in

**Table 3 | Independent testing on the TCGA breast cancer subtype datasets**

|  | Endo | GGI97 | LM | Mamma | Pam50 | RS | scP.W | netHD | netNE |
|---|---|---|---|---|---|---|---|---|---|
| ER+ (OS) | 0.6 | 0.55 | 0.59 | 0.55 | 0.65 | 0.65 | 0.62 | 0.62 | 0.67 |
| ER+ (RFS) | 0.57 | 0.52 | 0.59 | 0.51 | 0.53 | 0.6 | 0.61 | 0.6 | 0.64 |
| PR+ (OS) | 0.57 | 0.56 | 0.63 | 0.62 | 0.69 | 0.64 | 0.67 | 0.62 | 0.66 |
| PR + (RFS) | 0.58 | 0.56 | 0.59 | 0.53 | 0.54 | 0.57 | 0.64 | 0.62 | 0.66 |
| HER2+ (OS) | 0.85 | 0.65 | 0.81 | 0.65 | 0.71 | 0.82 | 0.73 | 0.74 | 0.73 |
| HER2+ (RFS) | 0.67 | 0.54 | 0.74 | 0.53 | 0.57 | 0.68 | 0.7 | 0.67 | 0.71 |
| TNBC (OS) | 0.57 | 0.7 | 0.58 | 0.59 | 0.61 | 0.59 | 0.55 | 0.71 | 0.55 |
| TNBC (RFS) | 0.52 | 0.66 | 0.52 | 0.53 | 0.52 | 0.61 | 0.56 | 0.64 | 0.53 |
| Mean rank | 5.8 | 6.5 | 4.9 | 7.5 | 5.5 | 3.6 | 3.9 | 3.6 | 3.4 |



**Fig. 4 | Functional co-occurrence analysis. a** Functional co-occurrence of nine gene signatures within the TCGA cohort, with color intensity indicating the degree of co-occurrence; darker colors denote greater degrees of co-occurrence. The different text colors denote different clusters. **b** Similar to (**a**), the functional co-occurrence of the same nine gene signatures within the METABRIC cohort.

cancer prognosis but also uniquely showed enrichment in immune-related pathways. Looking ahead, a predictive model for cancer immunotherapy could be developed by applying single-sample gene set enrichment analysis (ssGSEA)[52,53], machine learning[54,55], or neural networks[56,57] based on our gene signatures. These approaches could also be adapted for predicting breast cancer immunotherapy responses. Additionally, the gene signatures we identified could be used to develop a prognosis model specifically tailored for breast cancer subtypes.

## Methods
### Pan-cancer immunotherapy cohort
Pan-cancer immunotherapy responses play a crucial role in filtering anchor genes and transferring clinical phenotypes. We systematically collected a total of 948 samples, including 11 cohorts from immune checkpoint therapy studies. Among these patients, 275 responded to immunotherapy, and 673 were non-responders. These cohorts included five melanoma cohorts (SKCM): Liu2019[58], Riaz2017[59], Gide2019[60], Van Allen2015[61], and HugoLo2016[62]; three renal cell carcinoma cohorts (KIRC): IMmotion150[63], Miao2018[64], and Choueiri2016[65]; one gastric cancer cohort (STAD)[66]; one bladder cancer cohort (BLCA): IMvigor210[67]; and one glioma cohort

(GBM): Prins2019[68]. The entire set of immunotherapy responses, was evaluated according to mRECIST (Supplementary Table S4). To facilitate subsequent analyses, we converted all transcriptome data into TPM format.

### Breast cancer cohorts
To train and test the identified prognostic gene signatures, we aggregated breast cancer cohort data from various sources and molecular profiling techniques. Specifically, we collected breast cancer data from 987 patients from the TCGA cohort, 1420 patients from the METABRIC, and 216 patients from the GEO data platform GPL6098, with accession number GSE22219[69]. The TCGA and METABRIC cohorts provided data on overall survival (OS) and relapse-free survival (RFS), while the GPL6098 cohort provided only RFS data (Supplementary Table S5). Finally, to address batch effects across these datasets, we first applied a $\log2(x+1)$ transformation to stabilize variance and normalize the data, followed by batch effect correction using Combat[70].

### Identification of anchor genes
The anchor genes are phenotype-specific genes, with expression anchored to variations in clinical phenotypes across samples. These genes were identified from the pan-cancer immunotherapy cohort. According to
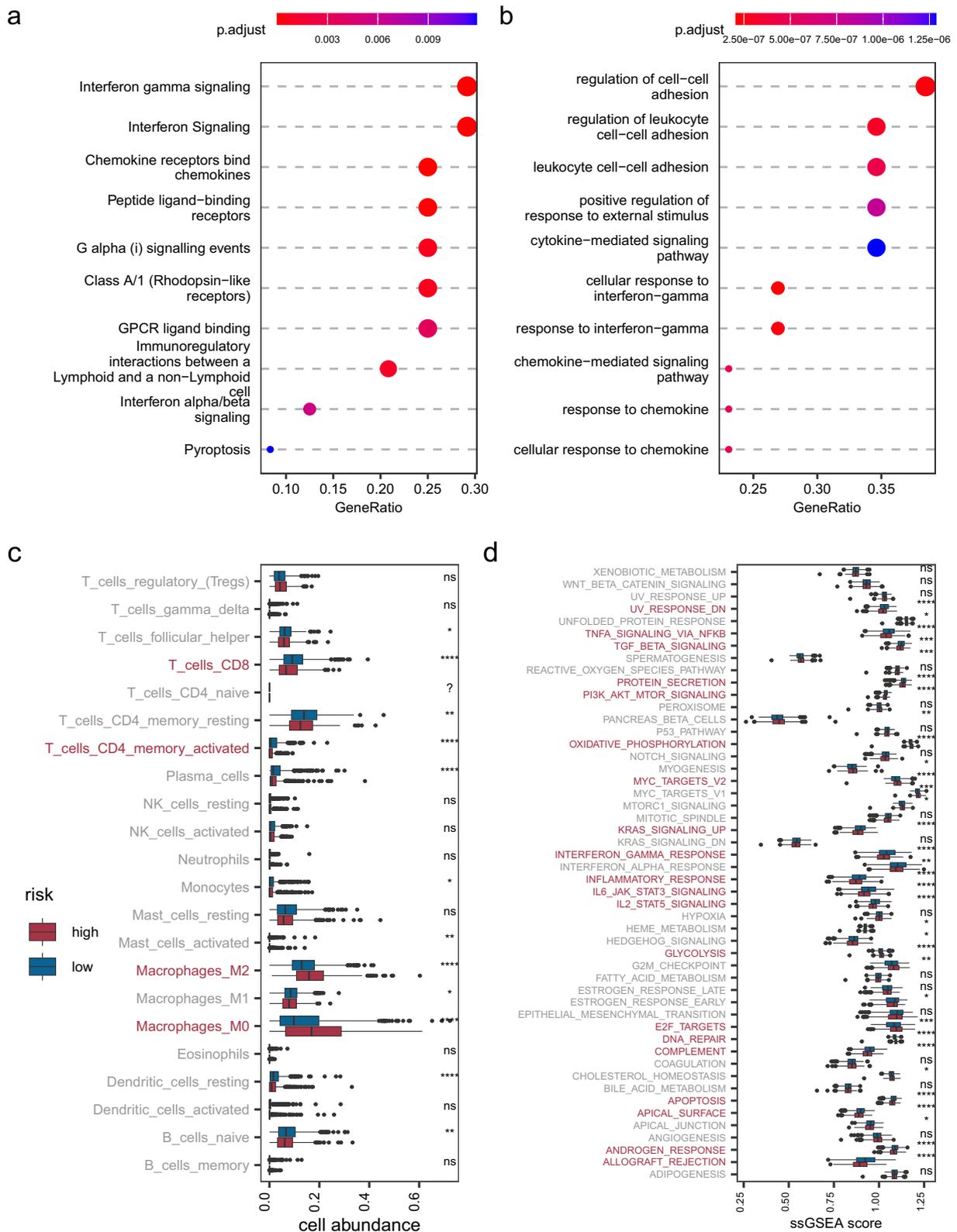
**Fig. 5 | Functional enrichment and tumor microenvironment (TME) analysis.**
**a** Gene set enrichment analysis for netHD using the Reactome database. **b** Similar to
(**a**), enrichment analysis using the Gene Ontology (GO) database. **c** Cellular
abundance in the TME between high and low prognosis risk groups. **d** Cancer
hallmark enrichment analysis for high and low prognostic risk groups.

mRECIST, immunotherapy responses are categorized into four groups: complete response (CR), partial response (PR), stable disease (SD), and progressive disease (PD).

To identify anchor genes, we quantified the immunotherapy responses using a digital scale (CR = 4, PR = 3, SD = 2, and PD = 1) based on the degree of response. We designed a user-defined consistent index to measure the consistency between gene expression and phenotypic variation across samples.

$$C_g = \frac{\sum[(g_{p_i} - g_{p_j}) \cdot (r_{p_i} - r_{p_j}) > 0]}{C_N^2} \qquad (1)$$

The formula (1) defines $C_g$ the consistency value for a specific gene $g$ within patients from the pan-cancer immunotherapy cohort. For each pair of patients $p_i$ and $p_j$, the product of the difference in gene $g$'s expression level $(g_{p_i} - g_{p_j})$ and the difference in immunotherapy responses grade $(r_{p_i} - r_{p_j})$ is calculated. We count only those pairs where the product of differences is positive, as indicated by the specified indicator function $[\cdot > 0]$. The denominator $C_N^2$, representing the total count of all possible unique patient pairs within the immunotherapy cohort, effectively normalized the sum. Ultimately, genes that exhibited a $C_g$ exceeding 0.6 were identified as anchor genes. In most prognostic models, a C-index of around 0.6 or higher is generally considered acceptable for predictive performance, which guided our decision to adopt this threshold directly[71] (see Supplementary Table S6 for a list of all anchor genes).

### Node embedding for PPI nodes

To identify gene signatures, a ranked list of genes related to clinical phenotypes is crucial for gene filtering. We mapped the anchor genes onto PPIs and added a virtual node connected to all anchor genes. The PPIs were extracted from STRING with a threshold score of 700. Each node in this new network was vectorized using the DeepWalk algorithm. Each node was sampled 100 times with a sampling length of six genes. We employed the Word2Vec method from gensim, setting a vector size of 256 for rich semantic detail and a window size of 4 to focus on close contextual relationships. The Skip-gram model (sg = 1) with negative sampling (negative = 10) was used to enhance robustness, while an initial learning rate of 0.03, decreasing to 0.0007, ensured stable convergence. The network was trained for 50 epochs to obtain embedding vectors for each node. These embedding vectors facilitate the generation of a gene ranking list for clinical phenotypes.

### Heat diffusion in PPI networks

Unlike the node embedding strategy, heat diffusion offers an alternative approach for obtaining a ranked gene list based on the heat value of each node. It is an important and widely used algorithm in systems biology, with applications in protein function prediction, disease gene prioritization, and patient stratification. Due to its algorithmic advantages in terms of memory usage, the heat diffusion model is much faster to compute, which is why we have chosen it for our analysis. The calculation is as follows (2):

$$d = h * exp(-Lt) \qquad (2)$$

where $h$ represents a vector of the original query, and $d$ is the resultant vector. $L$, the graph Laplacian, is defined as $D - A$, where $D$ is a diagonal matrix containing the degree of each node, and $A$ is the graph adjacency matrix of the input network. The scalar parameter $t$ represents the total diffusion time and controls the extent to which the original signal can spread across the network. We use a default value of $t = 0.1$. The term $exp(*)$ denotes the matrix exponential. Given the sensitivity of the model's performance to the network propagation parameter $t$, we conducted a sensitivity analysis (Supplementary Fig. S2). While $t = 0.1$ may not be the optimal choice, it is the default value commonly used in the Cytoscape heat diffusion plugin[72], and we have adopted it.

### Methods for evaluating prognostic models

To evaluate the performance of prognostic models, several widely used methods are often combined. Here is a brief overview of each:

**Concordance index.** The Concordance index (C-index)[73] measures the predictive accuracy of prognosis models by calculating the proportion of concordant pairs among all possible pairs in the study data. It is particularly useful for models predicting patient survival times or other time-related events. To compute the C-index, the input data included patient survival times or event occurrence times, observed status (with '1' indicating an event occurrence and '0' otherwise), and model-predicted scores for each patient. A pair is considered concordant if the patient with a longer actual survival has a higher predicted survival probability. Pairs, where neither patient experienced the event nor where one patient had not yet reached the event endpoint, were excluded from the calculation. The C-index is calculated as $K/M$, where $K$ is the number of concordant pairs, and $M$ is the total number of valid pairs evaluated.

**Hazard ratio (HR).** The hazard ratio was used to assess and compare the accuracy of various prognostic methods. To facilitate personalized treatment planning by clinicians, patients are often divided into high-risk and low-risk groups based on the median of their predicted risk scores. These risk scores are binarized, and each patient is assigned to a respective risk group, R. A Cox proportional hazards model is then utilized to estimate the differences in risk between these two survival groups, as described below (3):

$$h(t|R) = h_0(t)\exp(\beta R) \qquad (3)$$

where $h_0(t)$ is the baseline hazard function. The term $\exp(\beta)$, defined as the hazard ratio (HR), quantifies the difference in risk between two patient groups. The HR quantifies the difference in risk, with a higher HR indicating a greater disparity in risk between the low-risk and high-risk groups, thus reflecting the superior performance of the predictive method.

**KM survival curve.** The Kaplan–Meier (KM) survival curves[74] and the log-rank tests[75] were used to determine whether the two risk groups had significantly different survival patterns. In a KM plot, the $Y$-axis shows the probability of survival over time, represented on the $X$-axis. Distinct, non-overlapping KM curves indicate clear differences in survival outcomes. The log-rank test was used to determine if the survival curves were statistically equivalent, with a $p$-value less than 0.05 indicating a significant difference, suggesting effective differentiation between groups.

### Improved stepwise algorithm for gene signature identification

Both node embedding and heat diffusion algorithms generate a ranked gene list based on either the similarity between gene expression and clinical phenotype variation or the heat value of each gene. Given that our training dataset comprises approximately 1500 samples, we selected the top 150 candidate genes, representing the top 10% of the ranked genes. To construct a robust prognostic model, we needed to filter a minority of genes from the ranked list. Traditional stepwise strategies—forward, backward, and both—do not adequately filter a moderate size for gene signatures. Consequently, we obtained 139, 43, and 41 genes, respectively, without a significant increase in the C-index (supplementary Table S7). Instead, we devised an improved greedy strategy to identify the gene signature.

Using the ranked gene list, we employed a multi-iteration approach to progressively reduce the number of genes until the size stabilized. In this process, if adding a new gene decreases the model's performance (as measured by the C-index), we label that gene 'hindering'. Conversely, if adding a gene increases the C-index, it is labeled 'helpful'. Before the next iteration, we removed all 'hindering' genes and retained those labeled 'helpful'. This operation is repeated in subsequent iterations until the size of the gene set no longer changes. The final set of genes forms the gene signature we require. The pseudocode outlining the algorithm for identifying gene signatures

from a ranked gene list is provided (Supplementary Table S8). Iterations for gene signature identification for netHD are listed (Supplementary Fig. S3).

## Gene signature functional co-occurrence analysis

To assess the functional co-occurrence among breast cancer gene signatures within a specific cohort, we calculated the functional activity of each gene signature in each sample using the single-sample gene set enrichment analysis (ssGSEA) algorithm[53]. This process provided a set of functional activity scores for each gene signature across all patients. We then computed the Pearson correlation coefficient between these gene signatures. The resulting correlation values reflect the degree of functional co-occurrence among the signatures. To further analyze the clustering of gene signatures, a hierarchical clustering method was employed to extract cluster information.

## Data availability

All data utilized in this research are accessible publicly, with details provided in the Methods section. The paper includes web links and unique identifiers for the public cohorts and datasets used.

## Code availability

Source codes and data used to generate the results were deposited on GitHub https://github.com/immbal/NBT.

## References

1. Sung, H. et al. Global cancer statistics 2020: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA Cancer J. Clin.* **71**, 209–249 (2021).
2. Young, R. H. & Louis, D. N. The Warrens and other pioneering clinician pathologists of the Massachusetts General Hospital during its early years: an appreciation on the 200th anniversary of the hospital founding. *Mod. Pathol.* **24**, 1285–1294 (2011).
3. Goldhirsch, A. et al. Meeting highlights: international expert consensus on the primary therapy of early breast cancer 2005. *Ann. Oncol.* **16**, 1569–1583 (2005).
4. Goldhirsch, A. et al. Progress and promise: highlights of the international expert consensus on the primary therapy of early breast cancer. *Ann. Oncol.* **18**, 1133–1144 (2007).
5. Michor, F. & Polyak, K. The origins and implications of intratumor heterogeneity. *Cancer Prev. Res.* **3**, 1361–1364 (2010).
6. Sommer, S. & Fuqua, S. A. Estrogen receptor and breast cancer. *Semin. Cancer Biol.* **11**, 339–352 (2001).
7. Trabert, B., Sherman, M. E., Kannan, N. & Stanczyk, F. Z. Progesterone and breast cancer. *Endocr. Rev.* **41**, 320–344 (2020).
8. Loibl, S. & Gianni, L. HER2-positive breast cancer. *Lancet* **389**, 2415–2429 (2017).
9. Nielsen, T. O. et al. Assessment of Ki67 in breast cancer: updated recommendations from the international Ki67 in breast cancer working group. *J. Natl. Cancer Inst.* **113**, 808–819 (2021).
10. Abubakar, M. et al. Combined quantitative measures of ER, PR, HER2, and KI67 provide more prognostic information than categorical combinations in luminal breast cancer. *Mod. Pathol.* **32**, 1244–1256 (2019).
11. Rakha, E. A. & Pareja, F. G. New advances in molecular breast cancer pathology. *Semin. Cancer Biol.* **72**, 102–113 (2021).
12. Tsang, J. Y. S. & Tse, G. M. Molecular classification of breast cancer. *Adv. Anat. Pathol.* **27**, 27–35 (2020).
13. Latha, N. R. et al. Gene expression signatures: a tool for analysis of breast cancer prognosis and therapy. *Crit. Rev. Oncol. Hematol.* **151**, 102964 (2020).
14. Liu, R. et al. The prognostic role of a gene signature from tumorigenic breast-cancer cells. *N. Engl. J. Med.* **356**, 217–226 (2007).
15. Minn, A. J. et al. Genes that mediate breast cancer metastasis to lung. *Nature* **436**, 518–524 (2005).
16. Huang, S., Yang, J., Fong, S. & Zhao, Q. Artificial intelligence in cancer diagnosis and prognosis: opportunities and challenges. *Cancer Lett.* **471**, 61–71 (2020).
17. Parker, J. S. et al. Supervised risk predictor of breast cancer based on intrinsic subtypes. *J. Clin. Oncol.* **27**, 1160–1167 (2009).
18. Filipits, M. et al. A new molecular predictor of distant recurrence in ER-positive, HER2-negative breast cancer adds independent information to conventional clinical risk factors. *Clin. Cancer Res.* **17**, 6012–6020 (2011).
19. Paik, S. et al. A multigene assay to predict recurrence of tamoxifen-treated, node-negative breast cancer. *N. Engl. J. Med.* **351**, 2817–2826 (2004).
20. van 't Veer et al. Gene expression profiling predicts clinical outcome of breast cancer. *Nature* **415**, 530–536 (2002).
21. Sotiriou, C. et al. Gene expression profiling in breast cancer: understanding the molecular basis of histologic grade to improve prognosis. *J. Natl. Cancer Inst.* **98**, 262–272 (2006).
22. Loo, J. F., Ho, H. P., Kong, S. K., Wang, T. H. & Ho, Y. P. Technological advances in multiscale analysis of single cells in biomedicine. *Adv. Biosyst.* **3**, e1900138 (2019).
23. Li, X. et al. A novel single-cell based method for breast cancer prognosis. *PLoS Comput. Biol.* **16**, e1008133 (2020).
24. Duffy, M. J. et al. Clinical use of biomarkers in breast cancer: updated guidelines from the European Group on Tumor Markers (EGTM). *Eur. J. Cancer* **75**, 284–298 (2017).
25. Waldman, A. D., Fritz, J. M. & Lenardo, M. J. A guide to cancer immunotherapy: from T cell basic science to clinical practice. *Nat. Rev. Immunol.* **20**, 651–668 (2020).
26. Lencioni, R. & Llovet, J. M. Modified RECIST (mRECIST) assessment for hepatocellular carcinoma. *Semin Liver Dis.* **30**, 52–60 (2010).
27. Bagchi, S., Yuan, R. & Engleman, E. G. Immune checkpoint inhibitors for the treatment of cancer: clinical impact and mechanisms of response and resistance. *Annu. Rev. Pathol.* **16**, 223–249 (2021).
28. Szklarczyk, D. et al. The STRING database in 2023: protein–protein association networks and functional enrichment analyses for any sequenced genome of interest. *Nucleic Acids Res.* **51**, D638–D646 (2023).
29. Milacic, M. et al. The reactome pathway knowledgebase 2024. *Nucleic Acids Res.* **52**, D672–D678 (2024).
30. Gene Ontology Consortium et al. The gene ontology knowledgebase in 2023. *Genetics* **224**, iyad031 (2023).
31. Newman, A. M. et al. Robust enumeration of cell subsets from tissue expression profiles. *Nat. Methods* **12**, 453–457 (2015).
32. Chen, B., Khodadoust, M. S., Liu, C. L., Newman, A. M. & Alizadeh, A. A. Profiling tumor infiltrating immune cells with CIBERSORT. *Methods Mol. Biol.* **1711**, 243–259 (2018).
33. Liberzon, A. et al. The molecular signatures database (MSigDB) hallmark gene set collection. *Cell Syst.* **1**, 417–425 (2015).
34. Ni, L. & Lu, J. Interferon gamma in cancer immunotherapy. *Cancer Med.* **7**, 4509–4516 (2018).
35. Versluis, J. M. et al. Interferon-gamma signature as prognostic and predictive marker in macroscopic stage III melanoma. *J. Immunother. Cancer* **12**, e008125 (2024).
36. Todorovic-Rakovic, N., Milovanovic, J., Greenman, J. & Radulovic, M. The prognostic significance of serum interferon-gamma (IFN-gamma) in hormonally dependent breast cancer. *Cytokine* **152**, 155836 (2022).
37. Lenz, G. et al. The origins of phenotypic heterogeneity in cancer. *Cancer Res.* **82**, 3–11 (2022).
38. Barrio-Hernandez, I. et al. Network expansion of genetic associations defines a pleiotropy map of human cell biology. *Nat. Genet.* **55**, 389–398 (2023).
39. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).

40. Ciregan, D., Meier, U., Schmidhuber, J. Multi-column deep neural networks for image classification, 2012 IEEE conference on computer vision and pattern recognition, 3642–3649 (IEEE, 2012).

41. Dai, H. et al. A cell proliferation signature is a marker of extremely poor outcome in a subpopulation of breast cancer patients. *Cancer Res.* **65**, 4059–4066 (2005).

42. Tokunaga, R. et al. CXCL9, CXCL10, CXCL11/CXCR3 axis for immune activation—a target for novel cancer therapy. *Cancer Treat. Rev.* **63**, 40–47 (2018).

43. Kohli, K., Pillarisetty, V. G. & Kim, T. S. Key chemokines direct migration of immune cells in solid tumors. *Cancer Gene Ther.* **29**, 10–21 (2022).

44. Bhat, P., Leggatt, G., Waterhouse, N. & Frazer, I. H. Interferon-gamma derived from cytotoxic lymphocytes directly enhances their motility and cytotoxicity. *Cell Death Dis.* **8**, e2836 (2017).

45. Wang, L. et al. The multiple roles of interferon regulatory factor family in health and disease. *Signal Transduct. Target Ther.* **9**, 282 (2024).

46. Aggarwal, V., Workman, C. J. & Vignali, D. A. A. LAG-3 as the third checkpoint inhibitor. *Nat. Immunol.* **24**, 1415–1422 (2023).

47. Raskov, H., Orhan, A., Christensen, J. P. & Gogenur, I. Cytotoxic CD8(+) T cells in cancer and cancer immunotherapy,. *Br. J. Cancer* **124**, 359–367 (2021).

48. Yang, I. et al. CD8+ T-cell infiltrate in newly diagnosed glioblastoma is associated with long-term survival. *J. Clin. Neurosci.* **17**, 1381–1385 (2010).

49. Fu, Q. et al. Prognostic value of tumor-infiltrating lymphocytes in melanoma: a systematic review and meta-analysis. *Oncoimmunology* **8**, 1593806 (2019).

50. Qiu, Y. et al. NKG2A(+)CD8(+) T cells infiltration determines immunosuppressive contexture and inferior response to immunotherapy in clear cell renal cell carcinoma. *J. Immunother. Cancer* **12**, e008368 (2024).

51. Ali, H. R. et al. Association between CD8+ T-cell infiltration and breast cancer survival in 12,439 patients. *Ann. Oncol.* **25**, 1536–1543 (2014).

52. Wu, C. C., Wang, Y. A., Livingston, J. A., Zhang, J. & Futreal, P. A. Prediction of biomarkers and therapeutic combinations for anti-PD-1 immunotherapy using the global gene network association. *Nat. Commun.* **13**, 42 (2022).

53. Barbie, D. A. et al. Systematic RNA interference reveals that oncogenic KRAS-driven cancers require TBK1. *Nature* **462**, 108–112 (2009).

54. Zhang, Z. et al. Integrated analysis of single-cell and bulk RNA sequencing data reveals a pan-cancer stemness signature predicting immunotherapy response. *Genome Med.* **14**, 45 (2022).

55. Kong, J. et al. Network-based machine learning approach to predict immunotherapy response in cancer patients. *Nat. Commun.* **13**, 3703 (2022).

56. Jiang, Y. et al. IRnet: immunotherapy response prediction using pathway knowledge-informed graph neural network. *J. Adv. Res.* **7**, S2090-1232(24)00320-5 (2024).

57. Ding, X., Zhang, L., Fan, M. & Li, L. TME-NET: an interpretable deep neural network for predicting pan-cancer immune checkpoint inhibitor responses. *Brief Bioinform.* **25**, bbae410 (2024).

58. Liu, D. et al. Integrative molecular and clinical modeling of clinical outcomes to PD1 blockade in patients with metastatic melanoma. *Nat. Med.* **25**, 1916–1927 (2019).

59. Riaz, N. et al. Tumor and microenvironment evolution during immunotherapy with nivolumab. *Cell* **171**, 934–949.e16 (2017).

60. Gide, T. N. et al. Distinct immune cell populations define response to anti-PD-1 monotherapy and anti-PD-1/anti-CTLA-4 combined therapy. *Cancer Cell* **35**, 238–255.e6 (2019).

61. Van Allen, E. M. et al. Genomic correlates of response to CTLA-4 blockade in metastatic melanoma. *Science* **350**, 207–211 (2015).

62. Hugo, W. et al. Genomic and transcriptomic features of response to anti-PD-1 therapy in metastatic melanoma. *Cell* **165**, 35–44 (2016).

63. Atkins, M. B. et al. IMmotion150: a phase II trial in untreated metastatic renal cell carcinoma (mRCC) patients (pts) of atezolizumab (atezo) and bevacizumab (bev) vs and following atezo or sunitinib (sun). *Am. Soc. Clin. Oncol.* **35**, 4505–4505 (2017).

64. Miao, D. et al. Genomic correlates of response to immune checkpoint therapies in clear cell renal cell carcinoma. *Science* **359**, 801–806 (2018).

65. Choueiri, T. K. et al. Immunomodulatory activity of nivolumab in metastatic renal cell carcinoma. *Clin. Cancer Res.* **22**, 5461–5471 (2016).

66. Kim, S. T. et al. Comprehensive molecular characterization of clinical responses to PD-1 inhibition in metastatic gastric cancer. *Nat. Med.* **24**, 1449–1458 (2018).

67. Balar, A. V. et al. Atezolizumab as first-line treatment in cisplatin-ineligible patients with locally advanced and metastatic urothelial carcinoma: a single-arm, multicentre, phase 2 trial. *Lancet* **389**, 67–76 (2017).

68. Cloughesy, T. F. et al. Neoadjuvant anti-PD-1 immunotherapy promotes a survival benefit with intratumoral and systemic immune responses in recurrent glioblastoma. *Nat. Med.* **25**, 477–486 (2019).

69. Buffa, F. M. et al. microRNA-associated progression pathways and potential therapeutic targets identified by integrated mRNA and microRNA expression profiling in breast cancer. *Cancer Res.* **71**, 5635–5645 (2011).

70. Leek, J. T., Johnson, W. E., Parker, H. S., Jaffe, A. E. & Storey, J. D. The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics* **28**, 882–883 (2012).

71. Zhang, Y., Wong, G., Mann, G., Muller, S., Yang, J. Y. H. SurvBenchmark: comprehensive benchmarking study of survival analysis methods using both omics data and clinical data. *Gigascience* **11**, giac071 (2022).

72. Carlin, D. E., Demchak, B., Pratt, D., Sage, E. & Ideker, T. Network propagation in the cytoscape cyberinfrastructure. *PLoS Comput. Biol.* **13**, e1005598 (2017).

73. Harrell, F. E. Jr., Lee, K. L. & Mark, D. B. Multivariable prognostic models: issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. *Stat. Med.* **15**, 361–387 (1996).

74. Rich, J. T. et al. A practical guide to understanding Kaplan–Meier curves. *Otolaryngol. Head. Neck Surg.* **143**, 331–336 (2010).

75. Bland, J. M. & Altman, D. G. The logrank test. *BMJ* **328**, 1073 (2004).

## Acknowledgements

## Author contributions

L.L. and M.F. collaborated as joint senior authors in designing and overseeing the study. L.Z. played a significant role in data collection and curation. X.D., L.Z., and M.F. conducted the analysis of the datasets and interpreted the findings. All authors contributed to the interpretation of the data, participated in writing and reviewing the manuscript, had complete access to all study data, and provided final consent for publication. The authors have thoroughly reviewed and approved the final manuscript.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1038/s41540-024-00486-7.

**Correspondence** and requests for materials should be addressed to Ming Fan or Lihua Li.

**Reprints and permissions information** is available at http://www.nature.com/reprints

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.