

Supplementary materials

A network medicine approach for identifying diagnostic and prognostic biomarkers and exploring drug repurposing in human cancer

Le Zhang^{1,#}, Shiwei Fan^{1,#}, Julio Vera^{2,3,4}, Xin Lai^{2,3,4,*}

1. College of Computer Science, Sichuan University, Chengdu, China
2. Laboratory of Systems Tumor Immunology, Department of Dermatology, Universitätsklinikum Erlangen and Friedrich-Alexander-Universität Erlangen-Nürnberg, Erlangen, Germany
3. Deutsches Zentrum Immuntherapie, Erlangen, Germany
4. Comprehensive Cancer Center Erlangen, Erlangen, Germany

#Equal contributors

*To whom correspondence should be addressed. E-mail: xin.lai@uk-erlangen.de

Supplementary table captions.....	3
Supplementary figure captions.....	5
Figure S1	8
Figure S2	9
Figure S3	10
Figure S4	11
Figure S5	12
Figure S6A	13
Figure S6B	14
Figure S7	15
Figure S8A	16
Figure S8B	17
Figure S8C	18
Figure S8D	19
Figure S9A	20
Figure S9B	21
Figure S10.....	22
Figure S11.....	23
Figure S12.....	24

Supplementary table captions.

Table S1 Differentially expressed genes in 18 cancer types. The file contains 18 sheets, each listing significantly differentially expressed genes with log₂ fold-change (tumor vs normal samples) and adjusted p-value for one cancer type.

Table S2 Gene set enrichment analyses in 18 cancer types. The file contains 18 sheets, each listing the hallmarks-of-cancer gene sets in which differentially expressed genes are significantly enriched for one cancer type. The columns from left to right are hallmark terms, the total number of genes in the hallmark term, normalized enrichment score, p-value, adjusted p-value, and the genes enriched in the hallmark term.

Table S3 Cancer type-specific networks. The file contains 18 sheets, each providing information for reconstructing a cancer type-specific network. The columns from left to right are source genes, target genes, interaction type (1: stimulation, 0: undirected, -1: inhibition), and edge weight quantified by Pearson correlation coefficients.

Table S4 Scores of nodes in cancer networks. The file contains 18 sheets, each listing scores of nodes computed by the network method for a cancer type. The columns from left to right are gene names, node degree, log₂ fold-change of gene expression (tumor vs normal samples), node score, and node rank according to node score. The genes highlighted in yellow are those showing the best performance for sample classification and thus regarded as cancer genes.

Table S5 Comparisons of performances of classifiers. The table shows the performance metrics of the classifiers trained in this study or from the literature. The pan-cancer classifiers from the literature include a regularized multi-task learning model in Hossain et al. (2021), a convolutional neural network model in Mostavi et al. (2020), and a variational autoencoder model in Withnell et al. (2021).

Table S6 Results of the Cox models. The file contains 18 sheets named after a cancer type, and each sheet contains detailed information about the Cox model using genes that are derived from the network method (network) or fold-changes in expression (log₂fc). The information includes the penalty parameter (λ) that gives the minimum error, the calculated error squares

and their estimated standard deviation, and the hazard ratio values of the corresponding genes (in capital letters).

Table S7 Drug information. The file contains drug annotations used in the drug repurposing analysis. The columns from left to right are identifiers of the DrugBank database, the drugs' names, targets, and FDA-approved treated cancer types by the drugs (the abbreviated names stand for the 18 cancer types investigated in this study).

Table S8 Networks for drug repurposing analysis. The file contains 14 sheets named after a cancer type plus a suffix '_node' or '_edge'. The node sheets provide information on the relevant genes, including gene names, log₂ fold-change of gene expression (tumor vs normal samples), adjusted p-value for expression change, and gene category (drug targets, cancer genes, first neighbor of drug targets, and other genes on the shortest path between drug targets and cancer genes). The edge sheets provide the interactions and their weights quantified by Pearson correlation coefficients. Drug-gene interactions do not have weights.

Supplementary figure captions

The supplementary figures are presented after the figure captions.

Figure S1 Flowchart of the computational methods used for conducting the analysis. The network medicine approach contains several modules (grey blocks). In each module, we showed the data (yellow rectangles), method and model used for analysis (purple rectangles), intermediate data (green rectangles), and final results (blue rectangles). The communication between different parts of a module is connected using thin arrowed lines. Cross communication between modules is connected using thin and thick arrowed lines. The texts in the parentheses are the corresponding source of the data, the packages and methods used by us, tables (Tab.), or figures (Fig.) presented in the article. DGE: differential gene expression. GSEA: gene set enrichment analysis.

Figure S2 Differential gene expression analysis. The heat map shows the hierarchical clustering of tumor samples (turquoise) and normal tissues (yellow) using significantly differentially expressed genes ($|\log_2 \text{fold-change}| \geq 2$ and adjusted p-value ≤ 0.05). We performed the hierarchical clustering using Euclidean distance with the average linkage algorithm.

Figure S3 Principal component analysis. The plot shows the clustering of tumor samples (blue) and normal tissues (yellow) in a reduced dimensional space. We performed the analysis using significantly differentially expressed genes.

Figure S4 Cancer networks. The plot shows the largest connected component for each of the 18 cancer-specific networks. The size of nodes is proportional to their node degree (in-degree plus out-degree). The node color visualizes log fold-changes of gene expression (tumor vs normal samples). The edge color visualizes the Pearson correlation between genes.

Figure S5 Node degree distribution of cancer networks. The plot shows the distribution of the node degree in cancer-specific networks. The line is fitted to the points using a linear function. The degree follows a power-law distribution in all 18 networks, i.e. few genes have many interactions and most genes have few interactions. The information that specifies FDA-

approved drugs for specific cancers was gathered from the NCI website (www.cancer.gov/about-cancer/treatment/drugs/cancer-type).

Figure S6 Hierarchical clustering of tumor and normal samples. The heat map shows the hierarchical clustering of tumor samples (turquoise) and normal tissues (yellow) using network-derived genes (A) or aberrantly expressed genes (B). We performed the hierarchical clustering using Euclidean distance with the average linkage algorithm. The genes were selected by the random forests that showed the best metrics in sample classification (Table 2).

Figure S7 Network topology analysis. The bar plot shows the ranking of four centrality values (such as degree, closeness, betweenness, and eigenvalue centrality) of the identified cancer genes in the networks. Degree centrality indicates how many connections each gene has to other genes in the network. Closeness centrality calculates the shortest paths between all genes, then assigns each gene a score based on its sum of shortest paths. Betweenness centrality shows which genes are 'bridges' between genes in a network. It does this by identifying all the shortest paths and then counting how many times each gene falls on one. Eigenvalue centrality shows that the centrality of a gene is proportional to the sum of the centralities of its direct neighbor genes. The y-axis is the percentile ranking of a gene using different centrality values (the higher the better ranking of a gene in the network), and the x-axis is the name of the cancer genes.

Figure S8 Density plots of p-values from the survival analyses using single genes. We used 10 (A), 20 (B), 50 (C) and 100 (D) genes to compare the performance of genes ranked by the network method (network) or fold-changes in expression (log2fc) in survival analyses. For each gene, we divided patients into high-expression (top 50%) and low-expression (remaining 50%) groups and performed Kaplan-Meier analysis (two-sided log-rank test). The p-values of the estimates were used to draw the density plots. To determine whether there were any differences in the distribution of p-values, we performed the Fisher-Pitman permutation test and added its calculated p-value at the top of the plots. The numbers in parentheses give the number of joint genes between the two ranking methods.

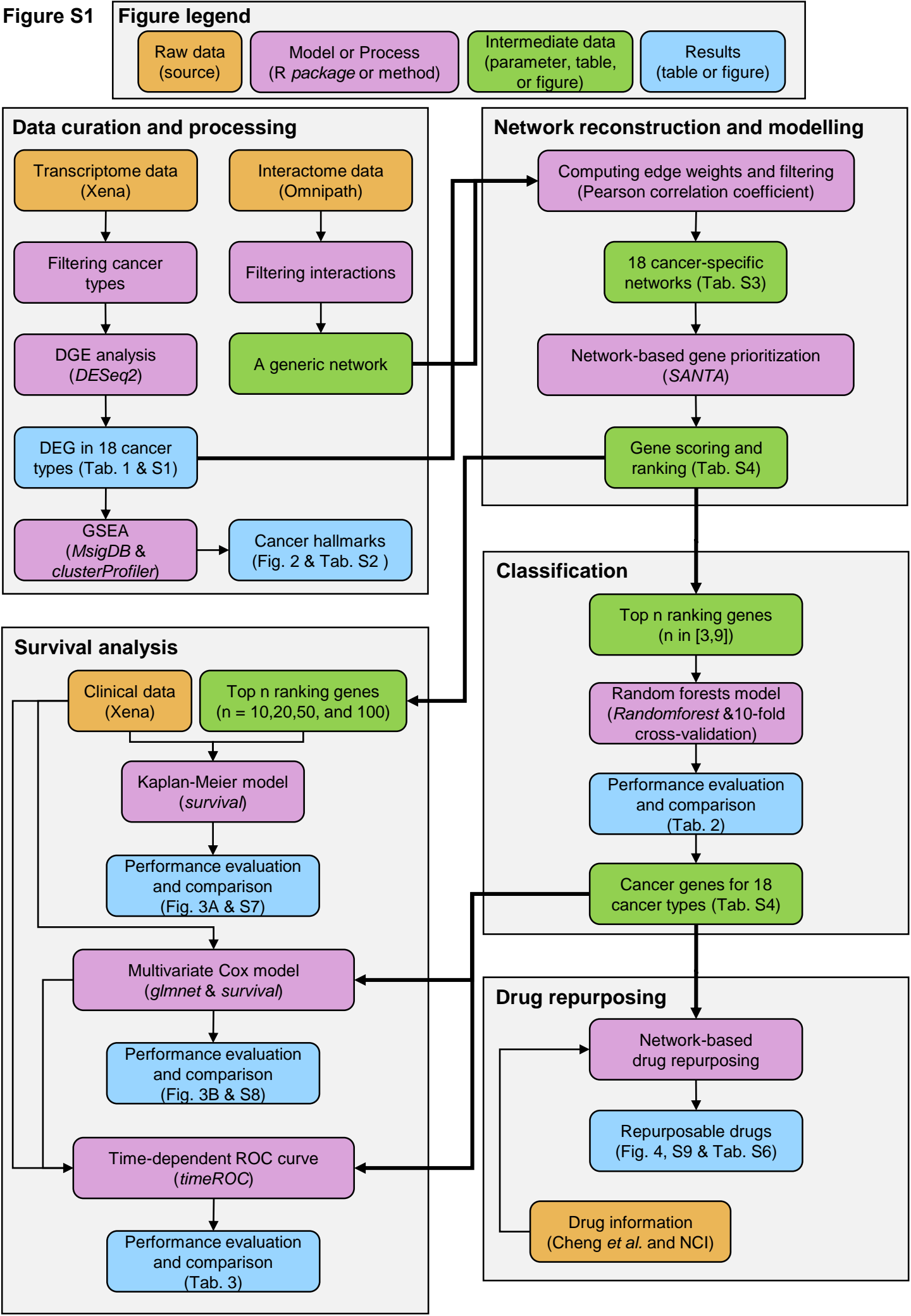
Figure S9 Comparison of survival analyses using combined genes. We used the combined genes identified by the network method (A) or fold-changes in expression (B) to compare their abilities in predicting cancer patients' survival rates. The patients were divided into high-risk and

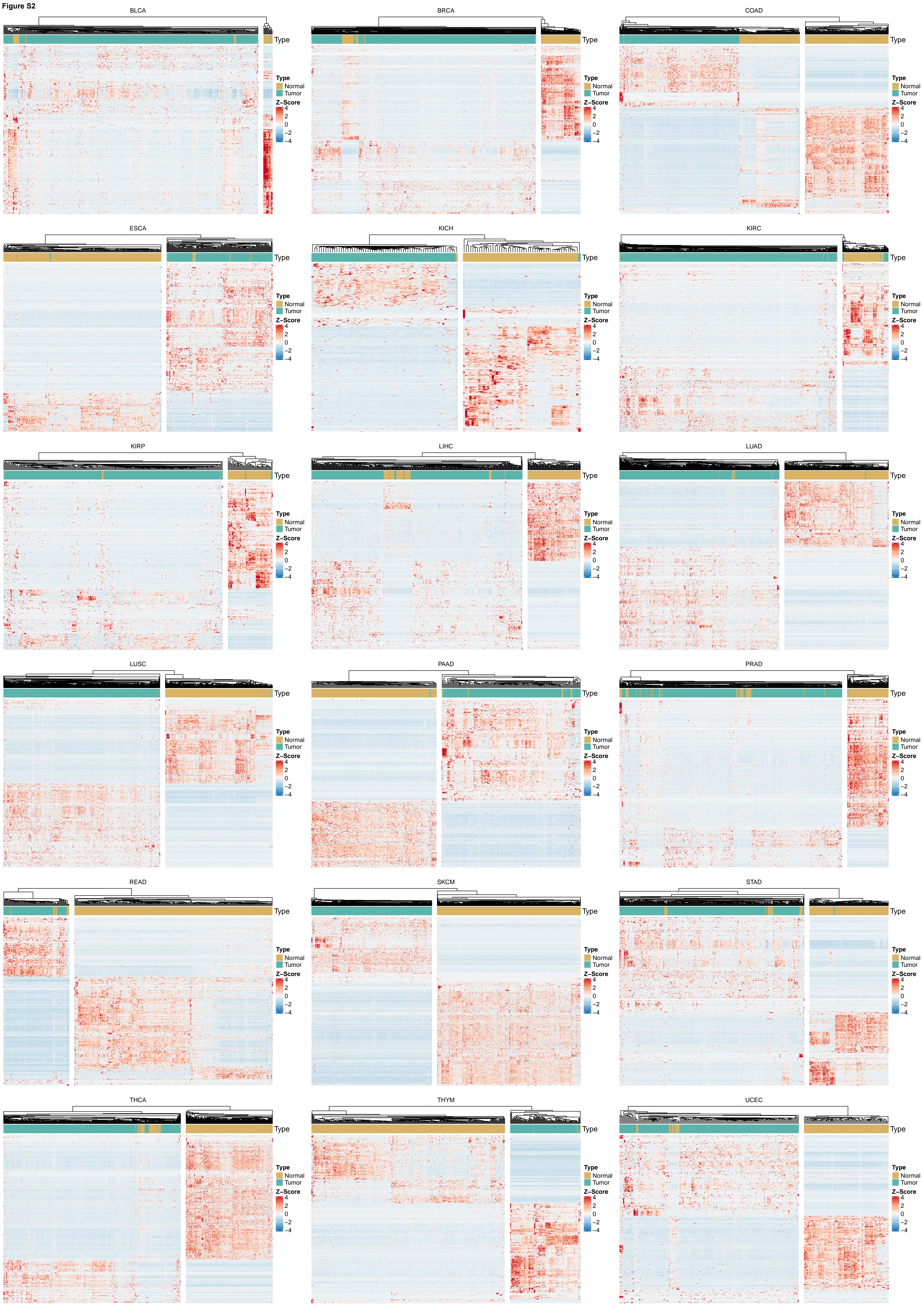
low-risk groups using the scores computed by a Cox regression model. The p-value of the survival curves was obtained using a two-sided log-rank test.

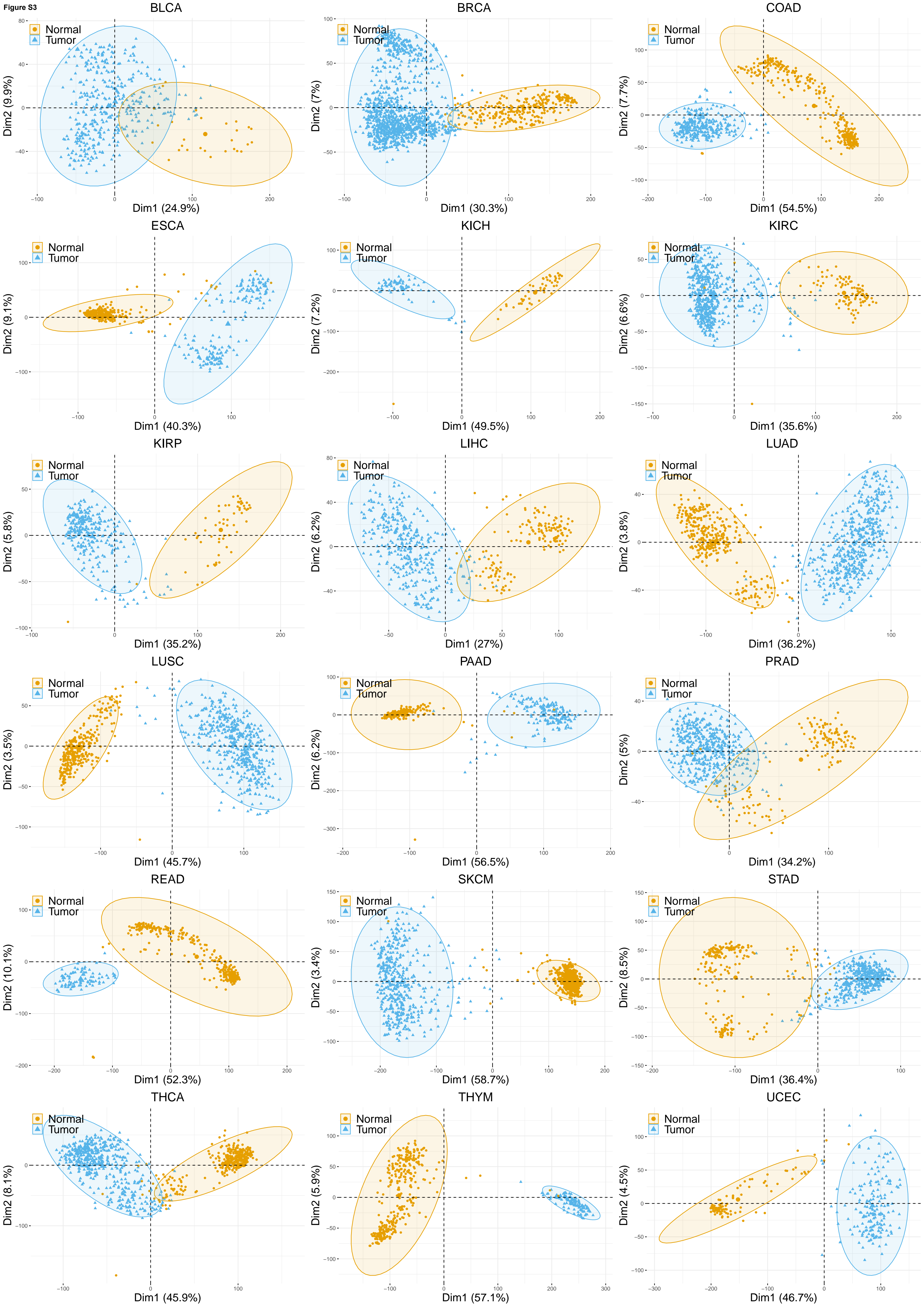
Figure S10 Drug distance to cancer gene and normal gene sets. The bar plot shows the distance of the identified repurposable drugs (x-axis) to cancer or normal genes in cancer networks. For each cancer type, its cancer gene set was identified using the network method while the counterpart is a normal gene set with high average expression in normal tissues and many interacting genes. The normal gene set contains top n genes (n equals to the number of genes in the corresponding cancer gene set). The genes were ranked by a score that multiplies the average expression of a gene in normal tissues with the gene's node degree in the network. The y-axis indicates the distance from the drug targets to the genes in both gene sets (see [Materials and Methods](#)).

Figure S11 Shortest paths between the repurposed drugs and cancer genes in KIRP (A) and LUAD (B). The paths were derived from the corresponding cancer-specific networks. Drugs and genes are shown as diamonds and circles, respectively. The node colors indicate the direction of gene expression change (red: upregulation; blue: downregulation). The label colors represent different gene categories: drug target genes (purple), cancer genes (orange), and other genes (grey) on the shortest path between drug targets and cancer genes. The edge colors visualize the interaction type between genes (pink: stimulation; green: inhibition); drug-target interactions are shown in black. The data illustrated in the network can be found in [Table S8](#). DB06616: bosutinib; DB08865: crizotinib; DB01254: dasatinib; DB00317: gefitinib; DB04868: nilotinib; DB05294: vandetanib.

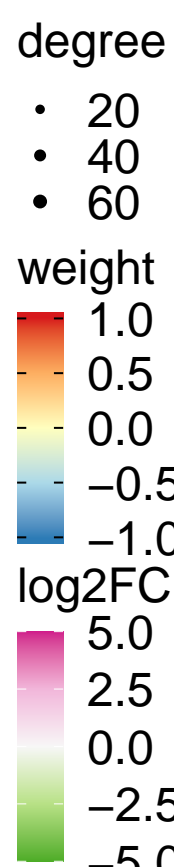
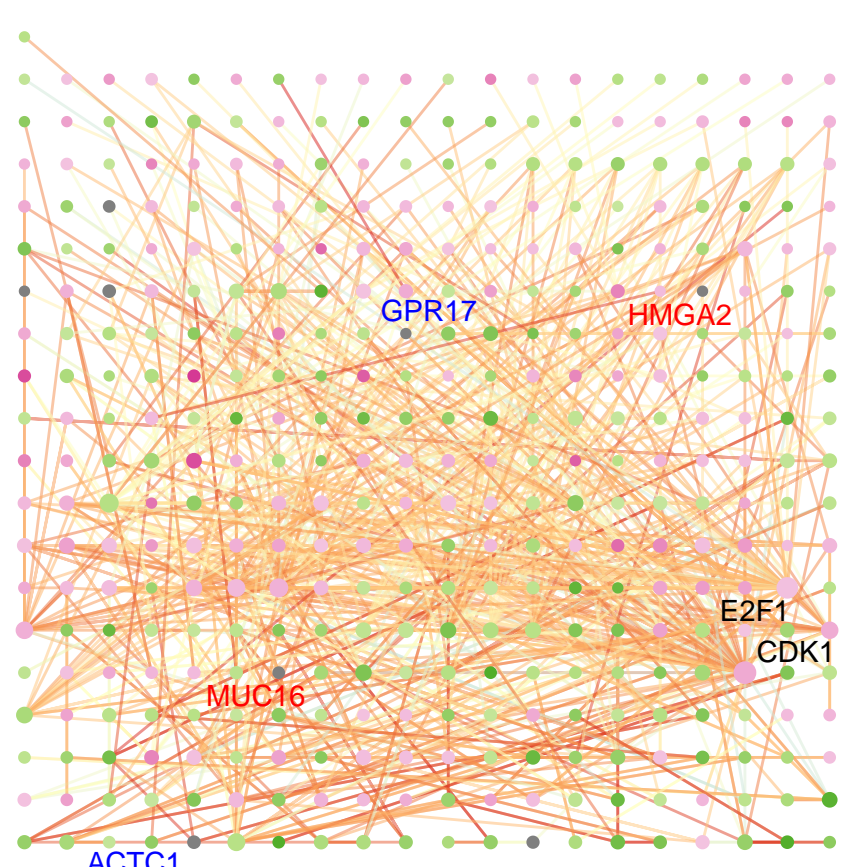
Figure S12 Cancer cell lines' sensitivity to the identified repurposable drugs. The plot shows the ranking (x-axis) of the drugs based on their measured IC₅₀ values (y-axis) in corresponding cancer cell lines (dots) for specific drugs. The higher ranking (more left on the x-axis) of a cell line the more sensitive it is to a specified drug. The IC₅₀ values of cell lines fall between the dashed black lines (i.e., maximum and minimum drug screening concentrations) indicating that the cell lines are sensitive (highlighted in green) to the drug otherwise resistant (highlighted in red). The figures were downloaded from the GDSC database and adapted.



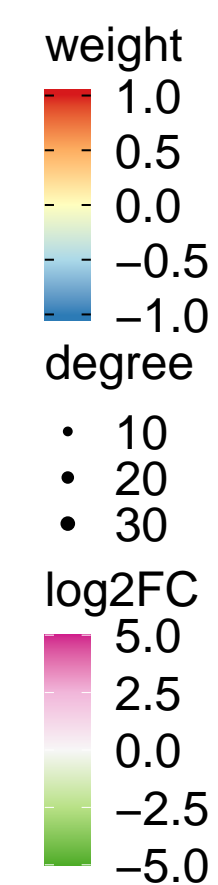
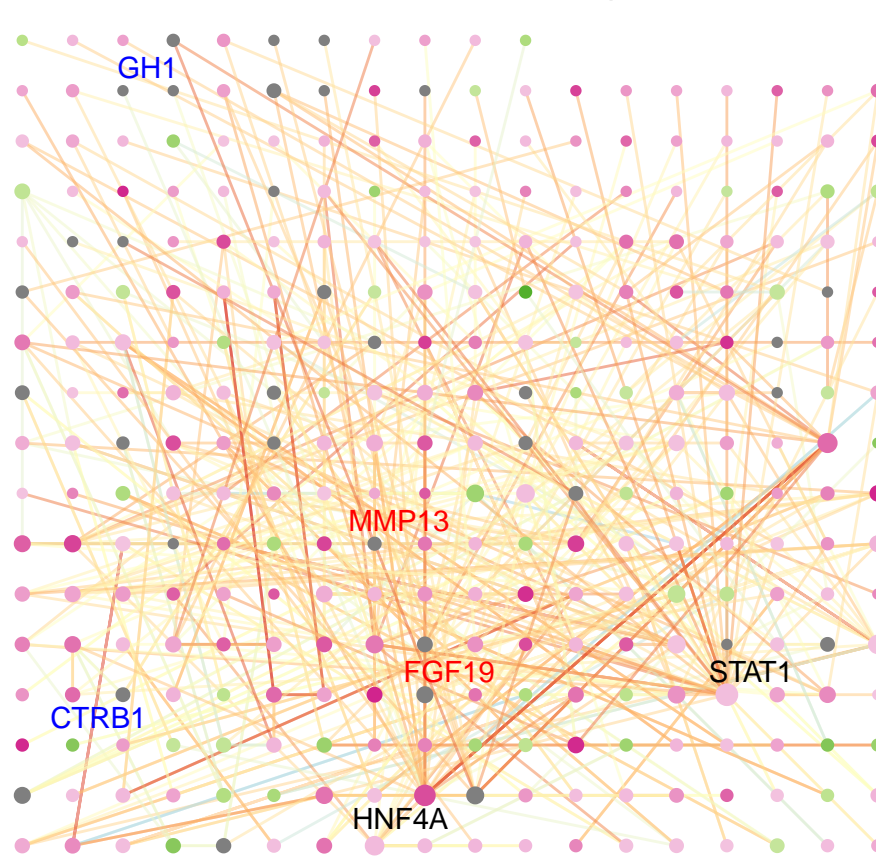




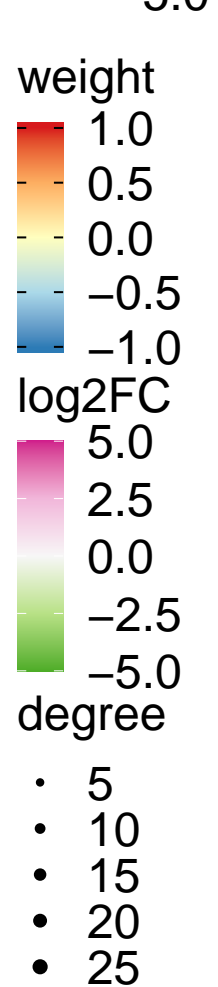
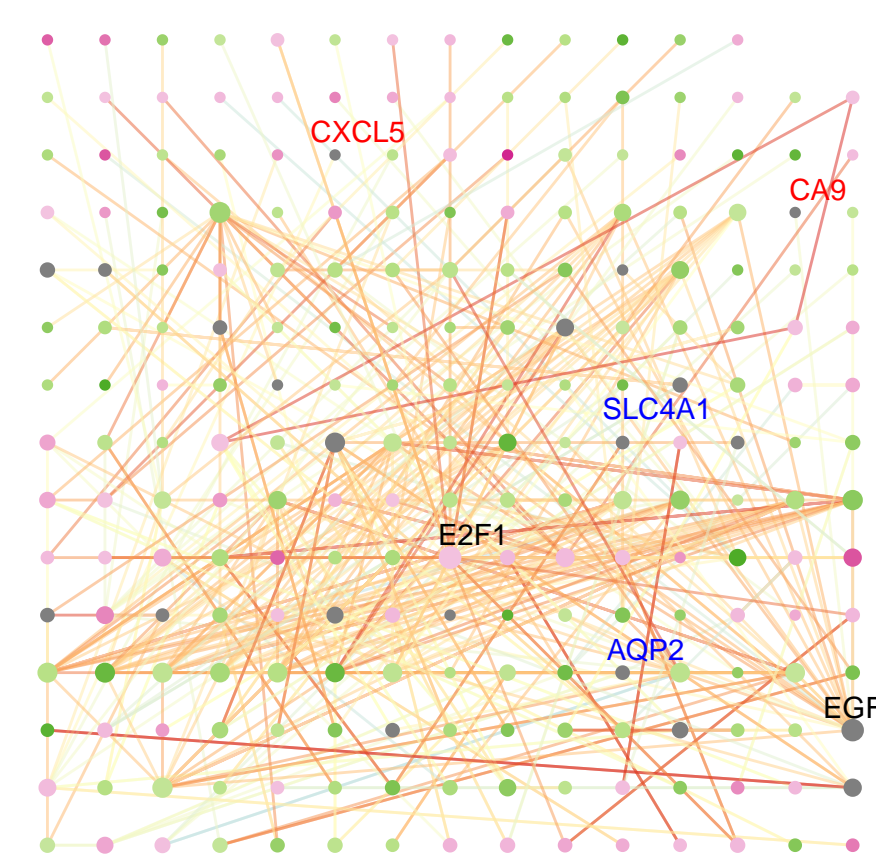
BLCA nodes:381 edges:851



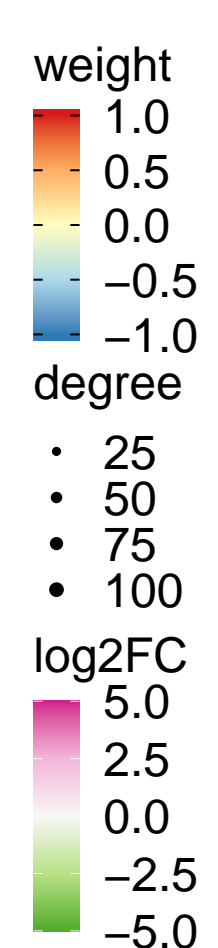
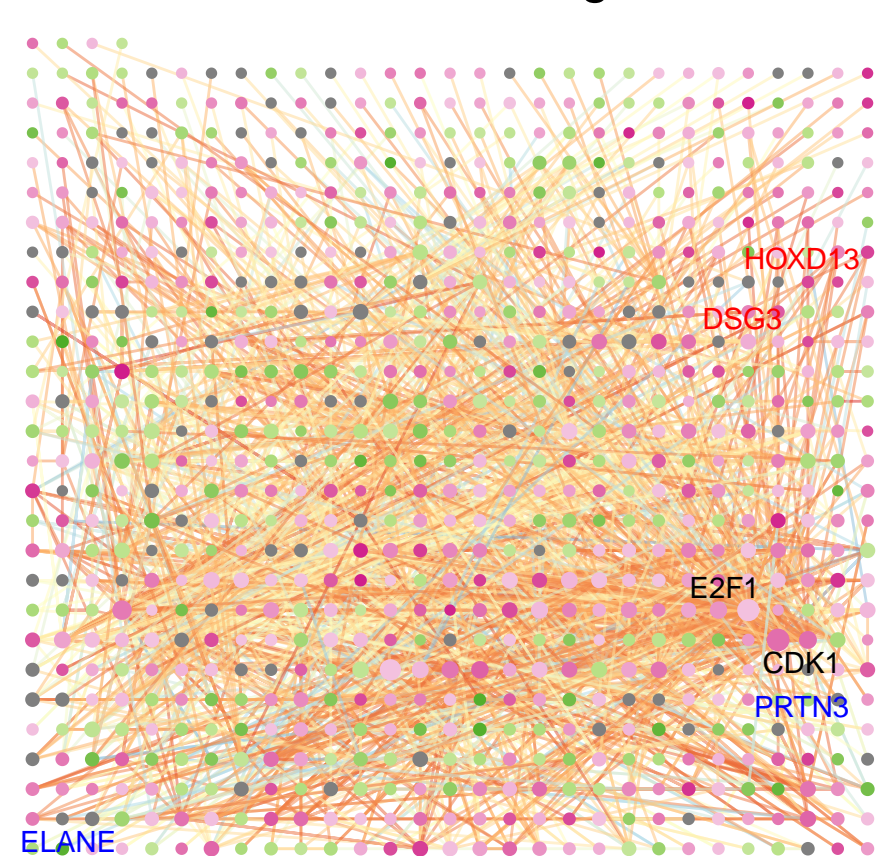
ESCA nodes:299 edges:515



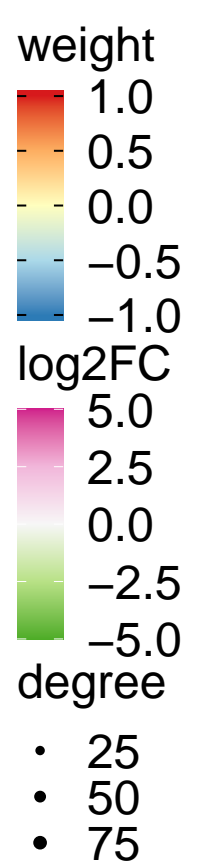
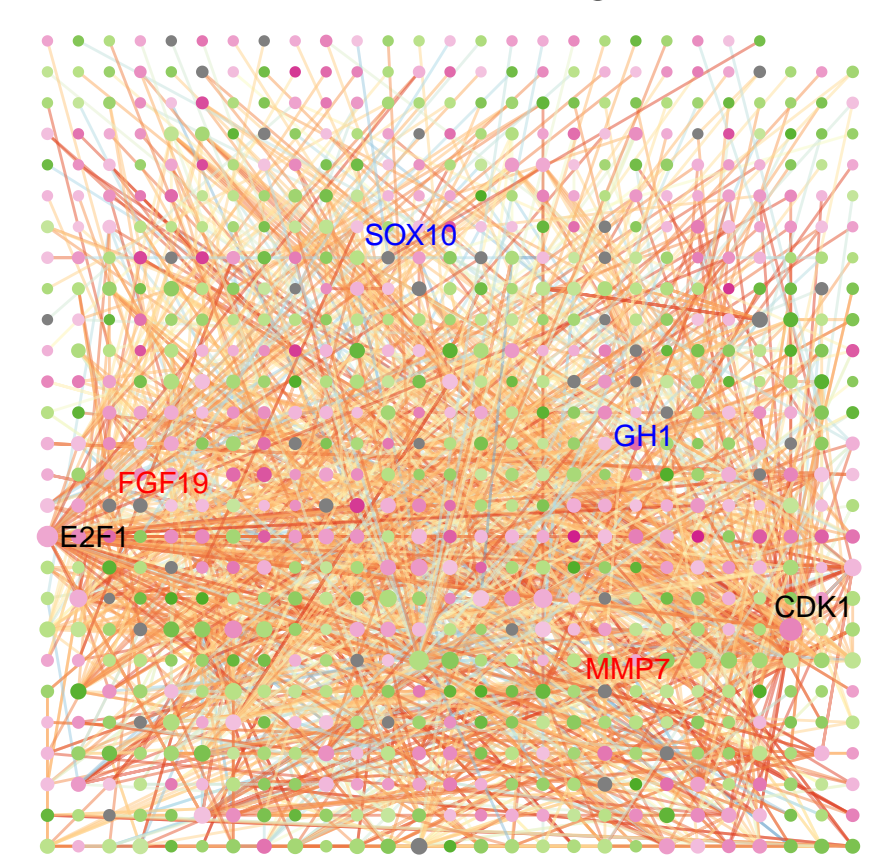
KIRP nodes:223 edges:424



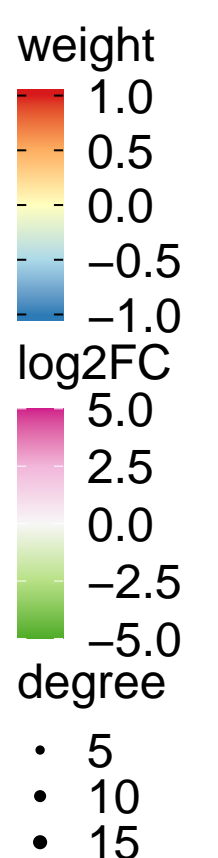
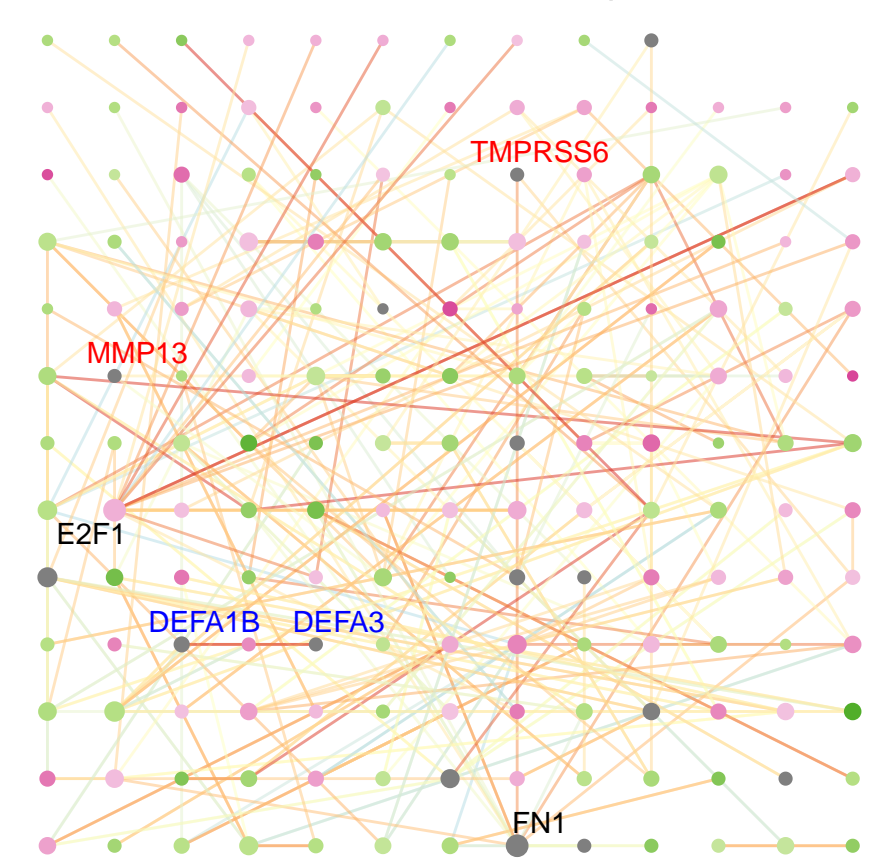
LUSC nodes:787 edges:1927



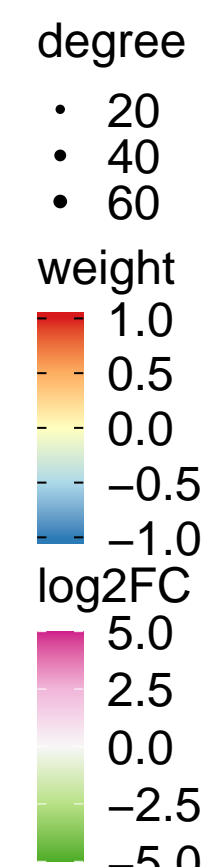
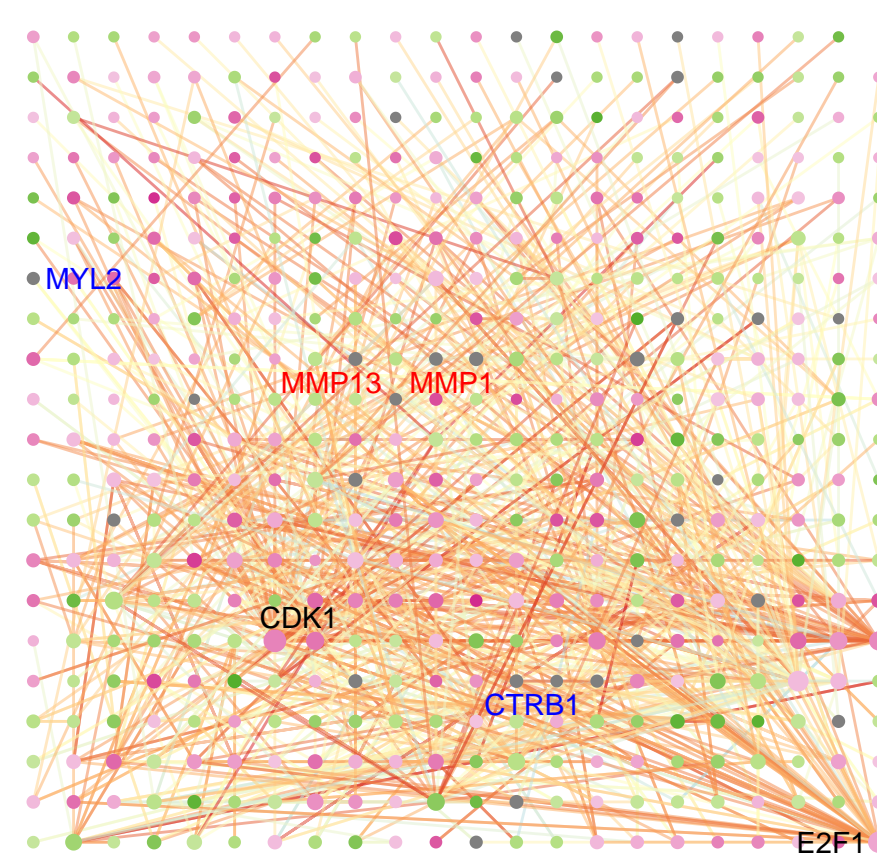
READ nodes:726 edges:1745



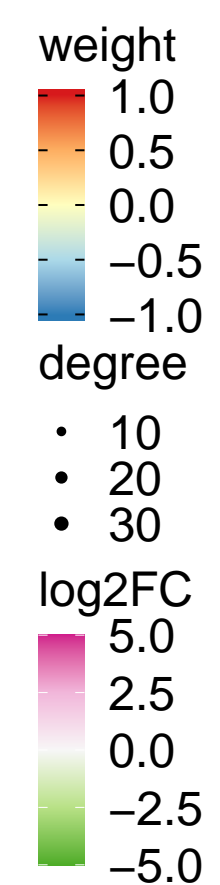
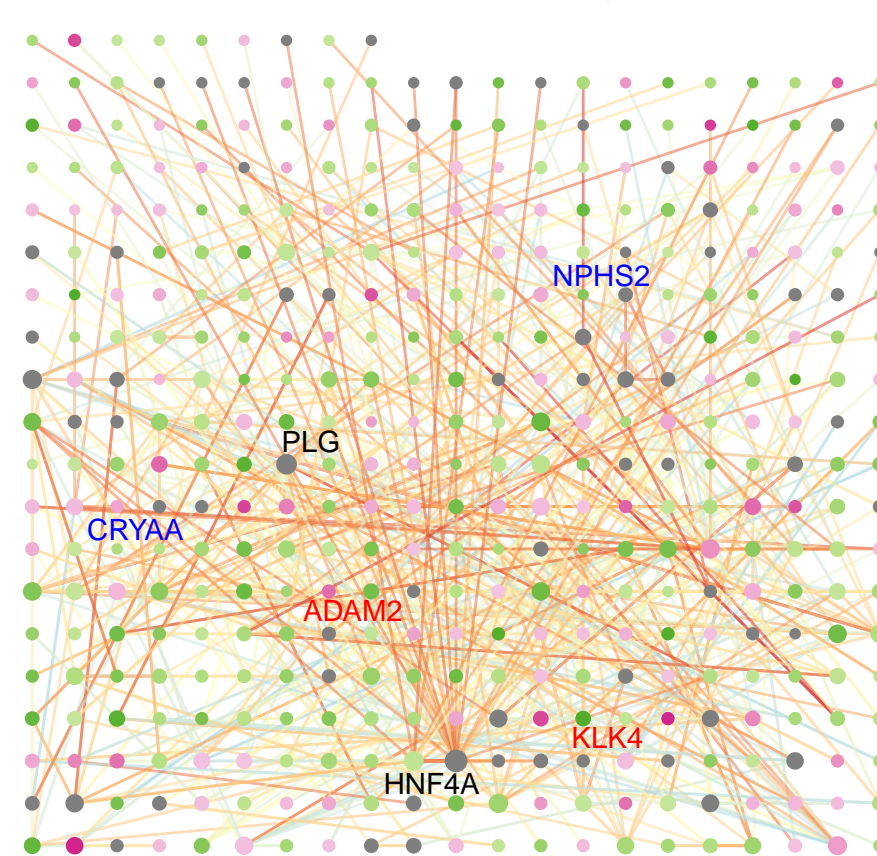
THCA nodes:166 edges:300



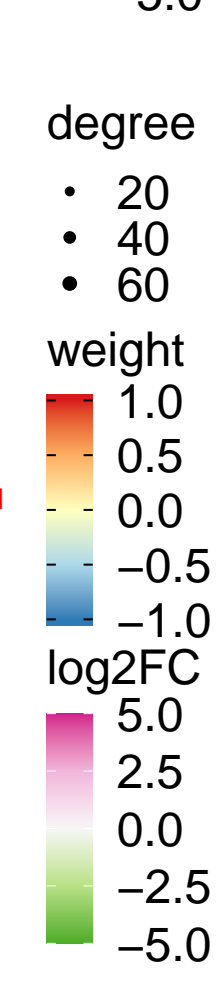
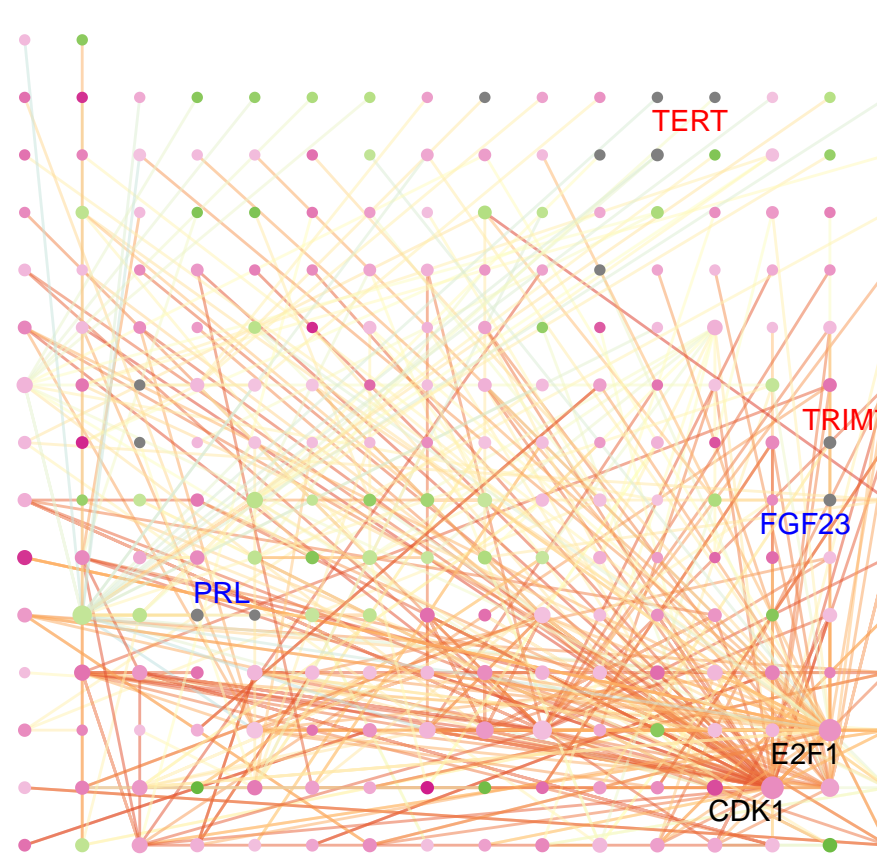
BRCA nodes:461 edges:962



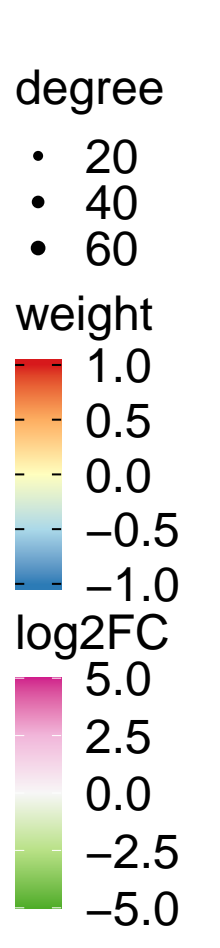
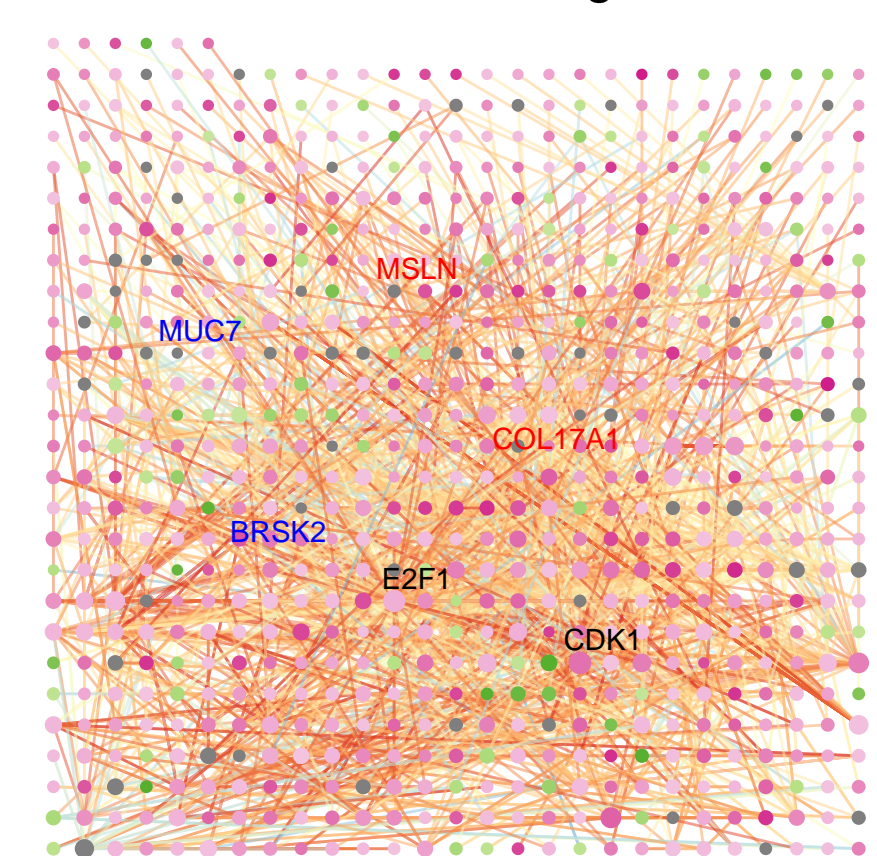
KICH nodes:408 edges:771



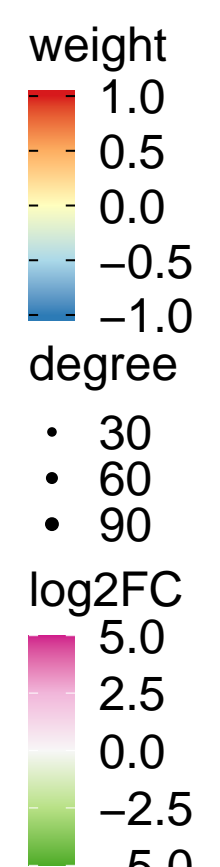
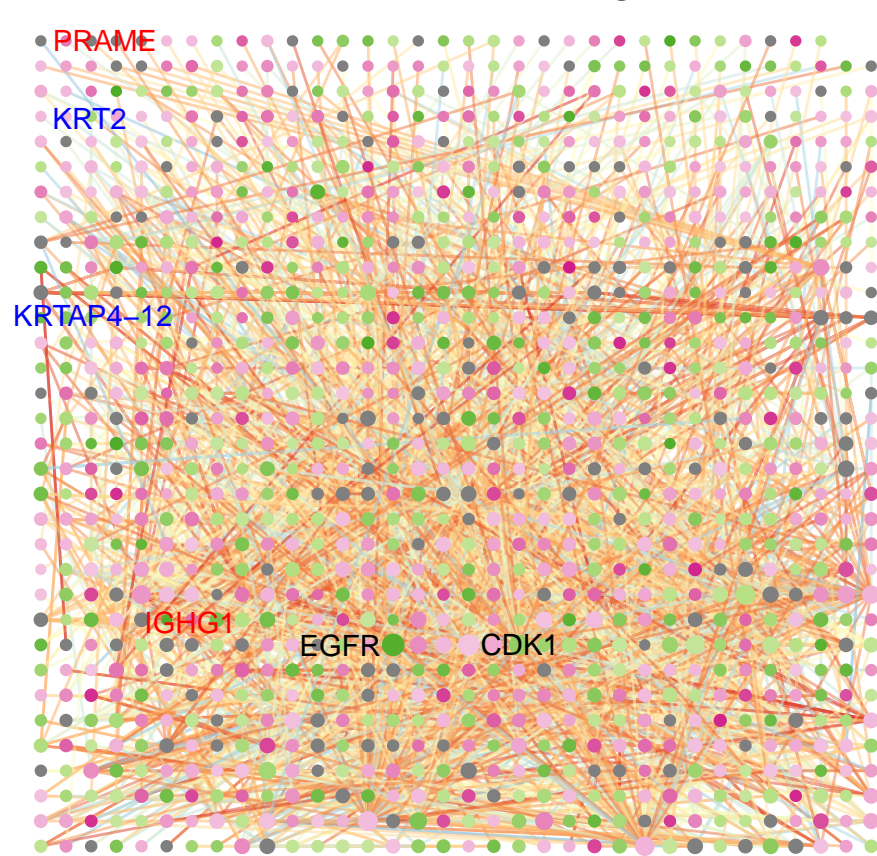
LIHC nodes:226 edges:472



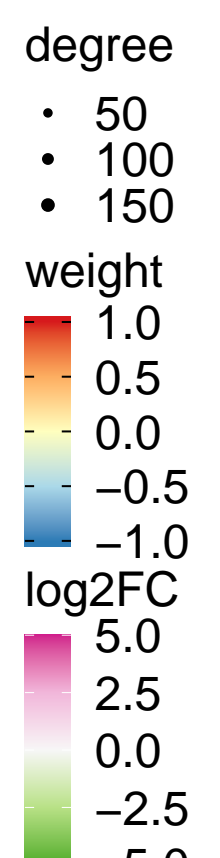
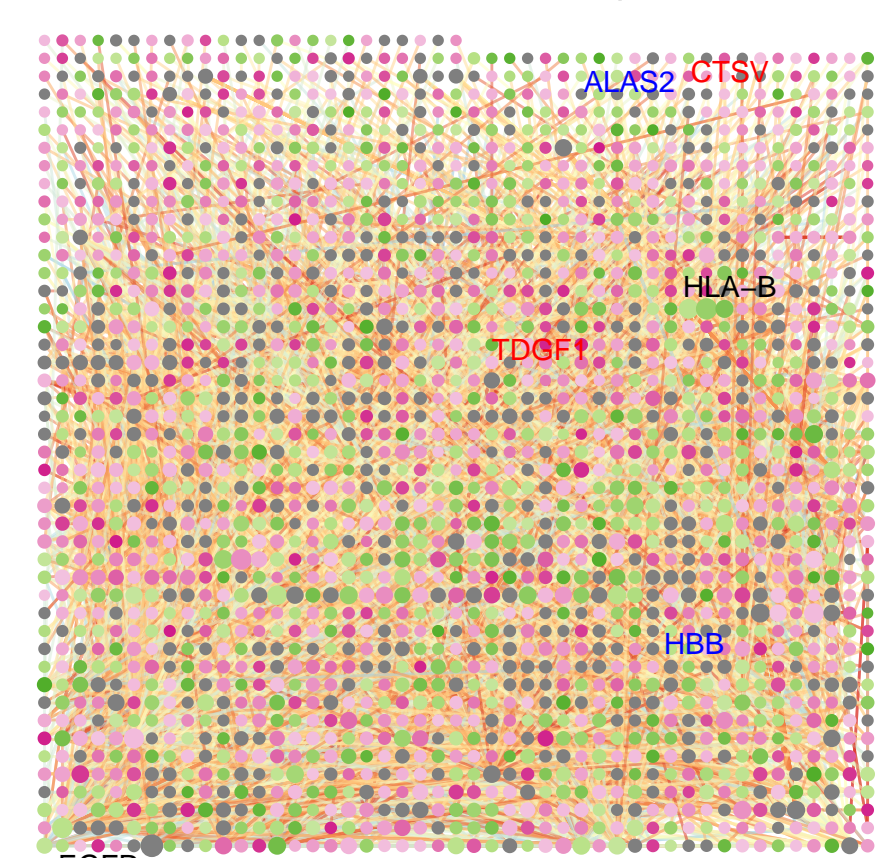
PAAD nodes:708 edges:1720



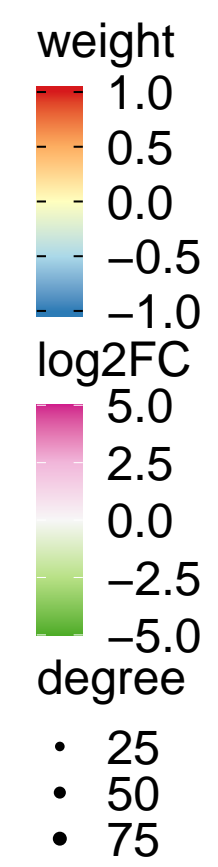
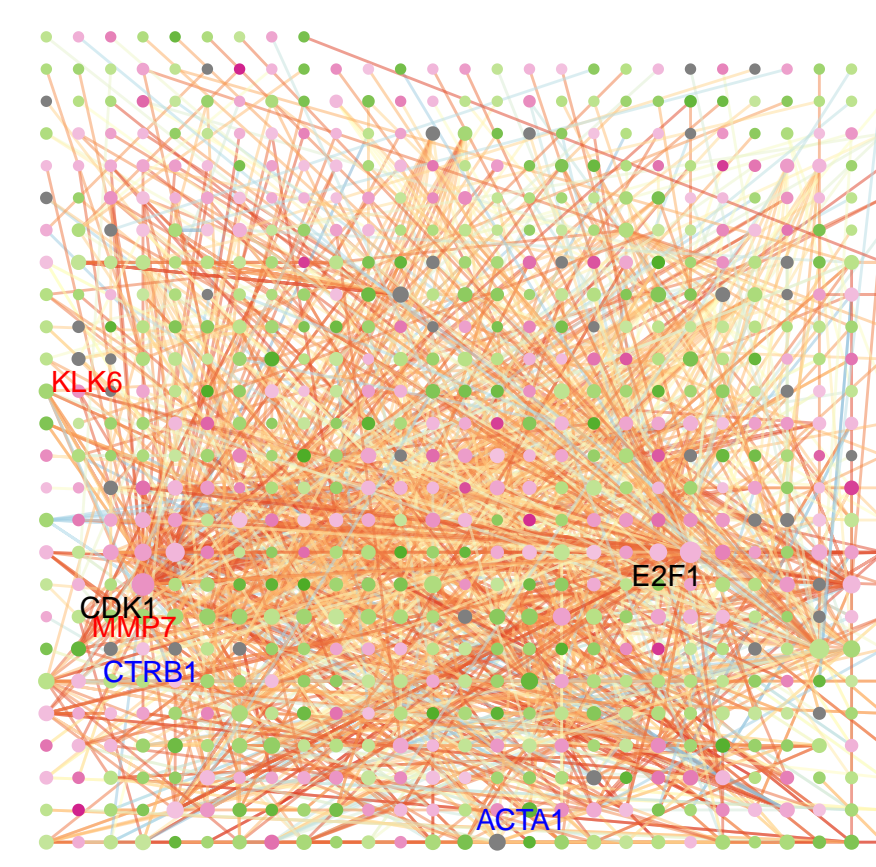
SKCM nodes:1120 edges:2641



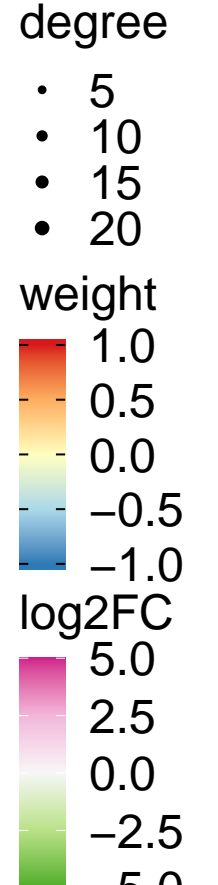
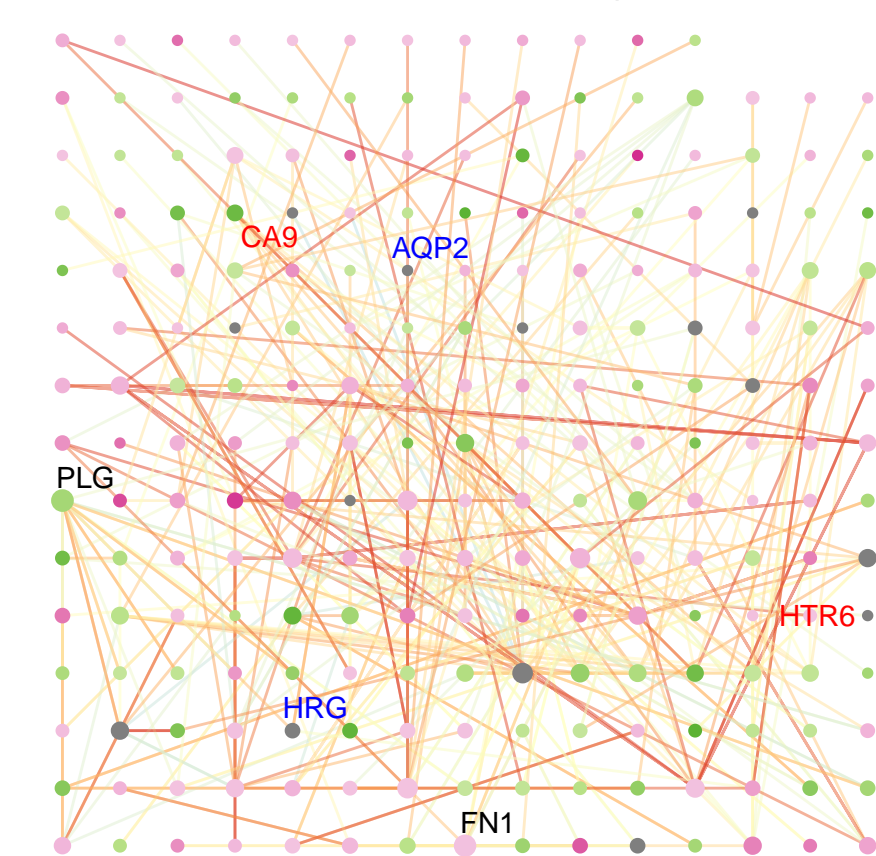
THYM nodes:2139 edges:7658



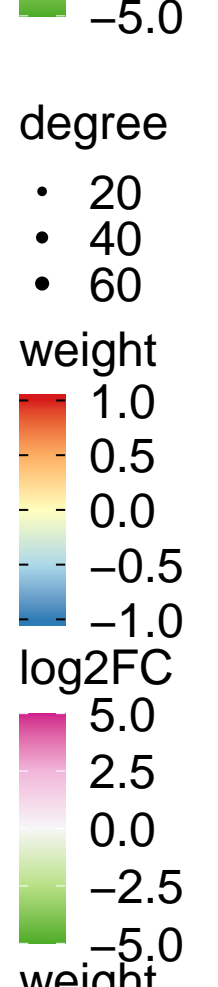
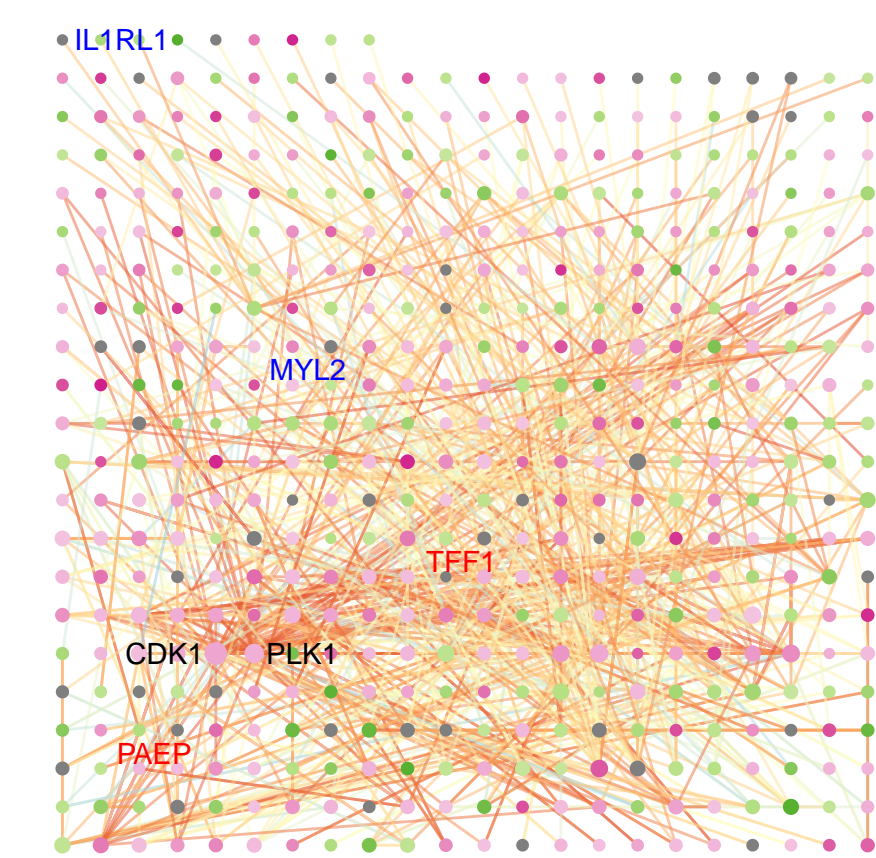
COAD nodes:684 edges:1616



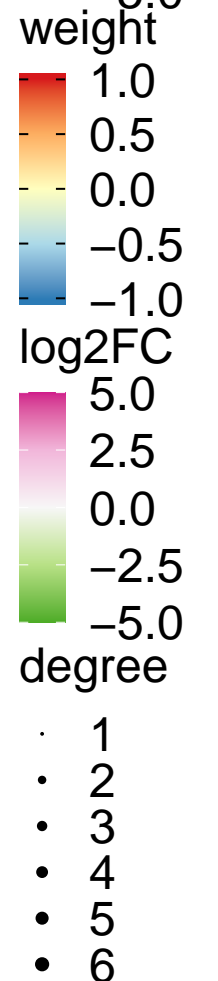
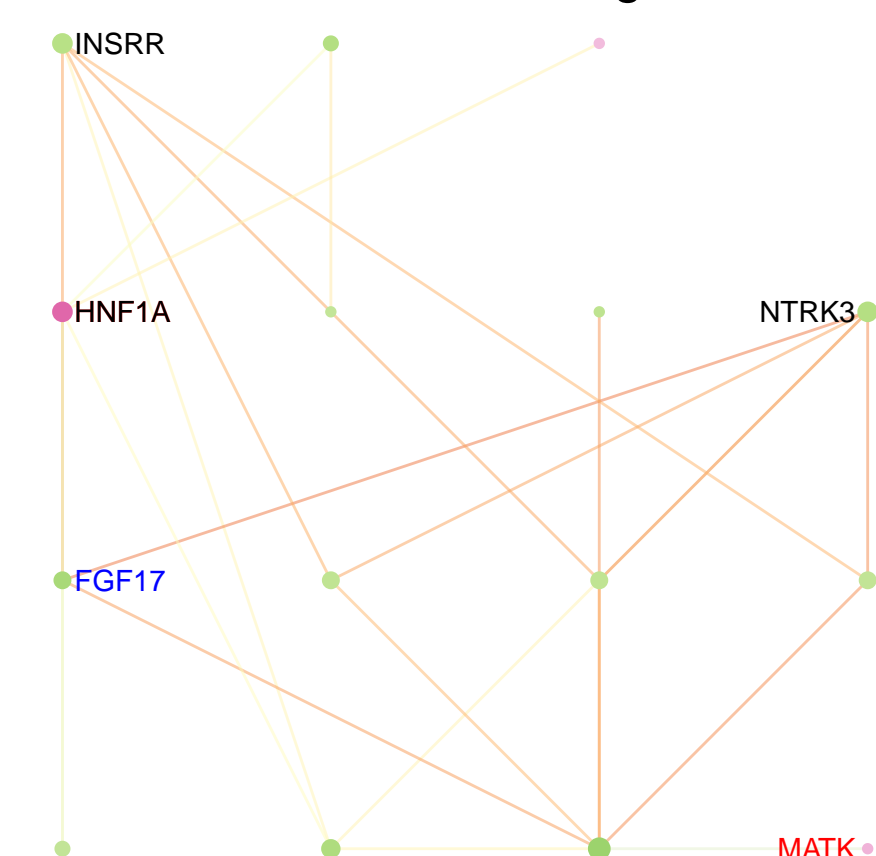
KIRC nodes:222 edges:374



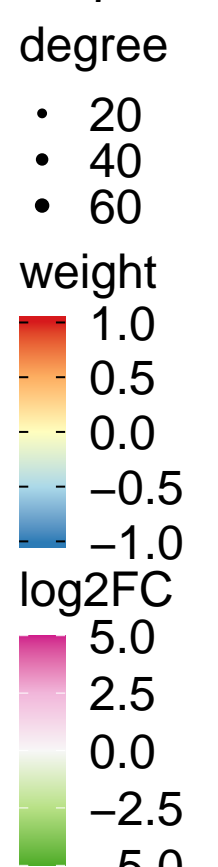
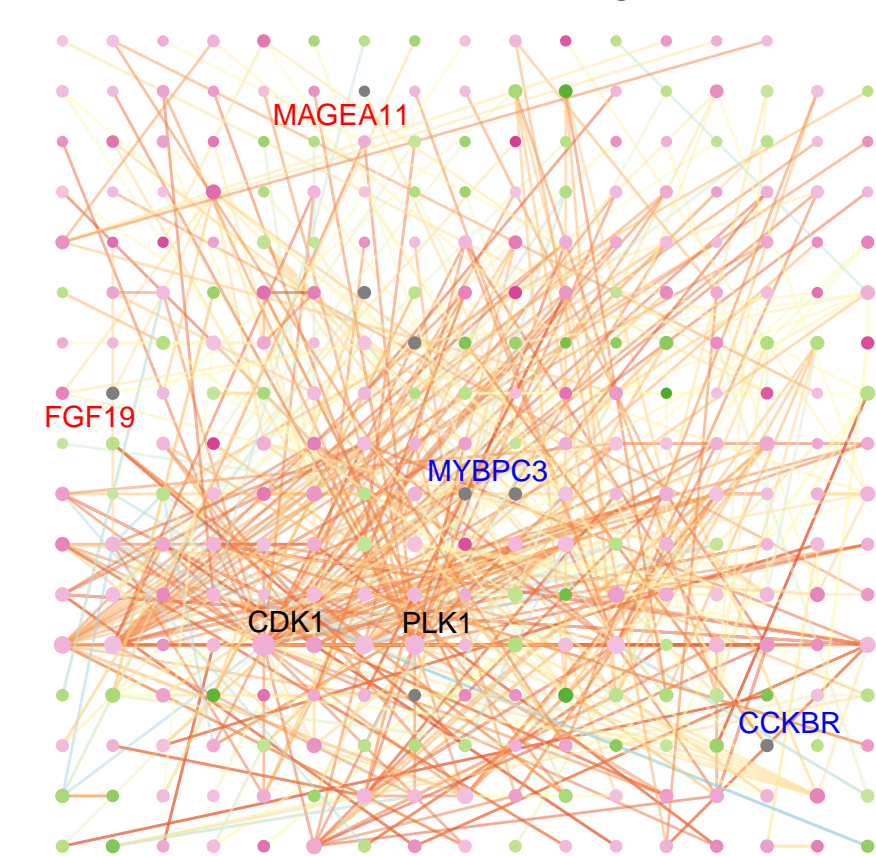
LUAD nodes:471 edges:950



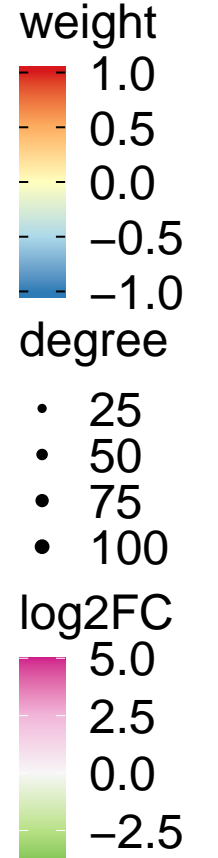
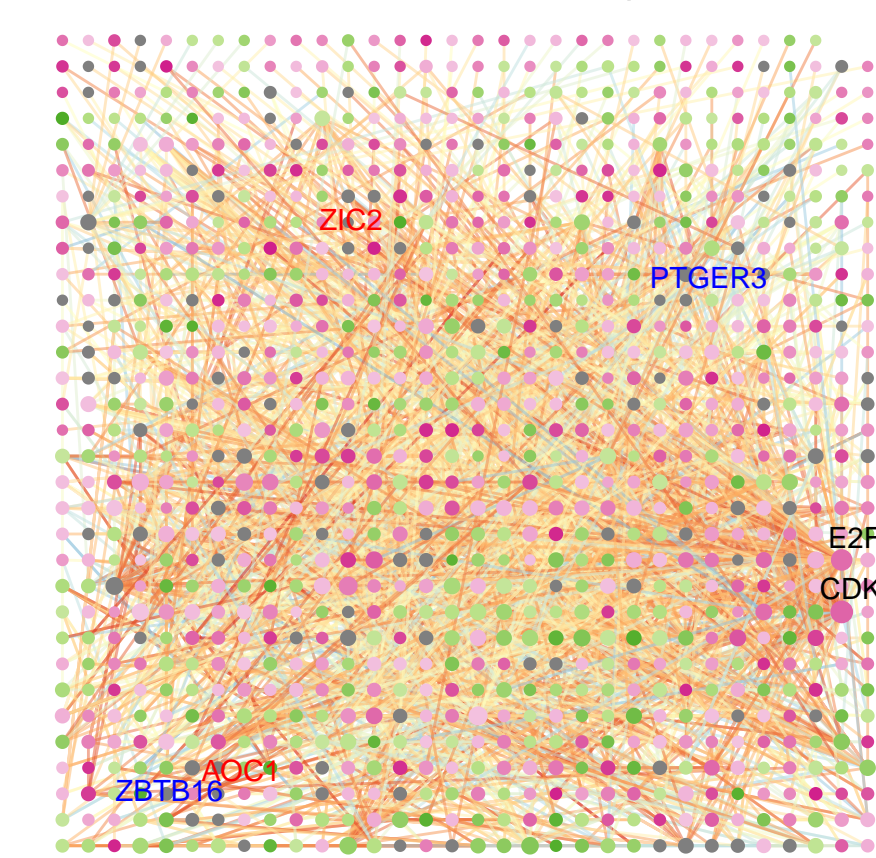
PRAD nodes:15 edges:23

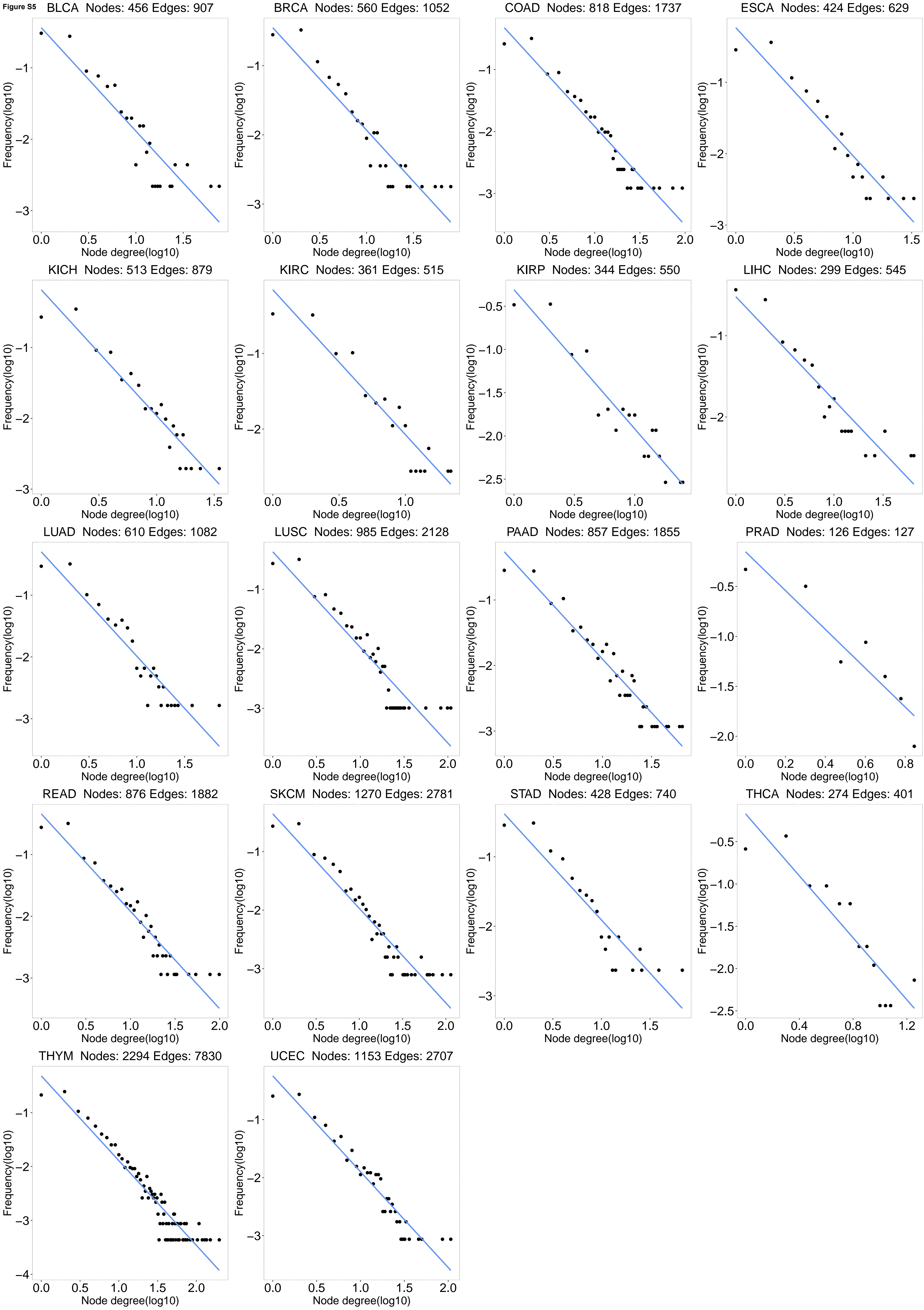


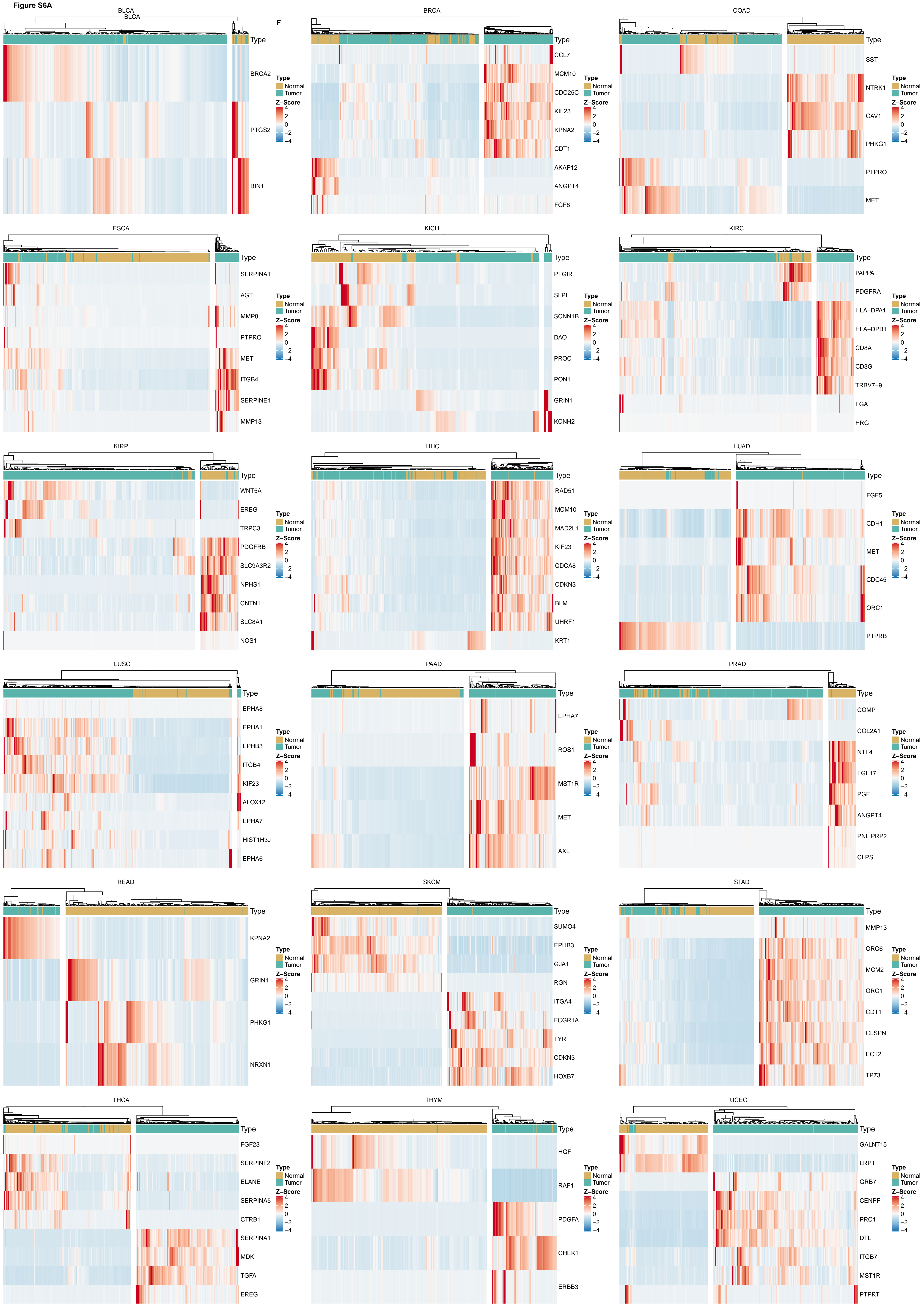
STAD nodes:287 edges:606

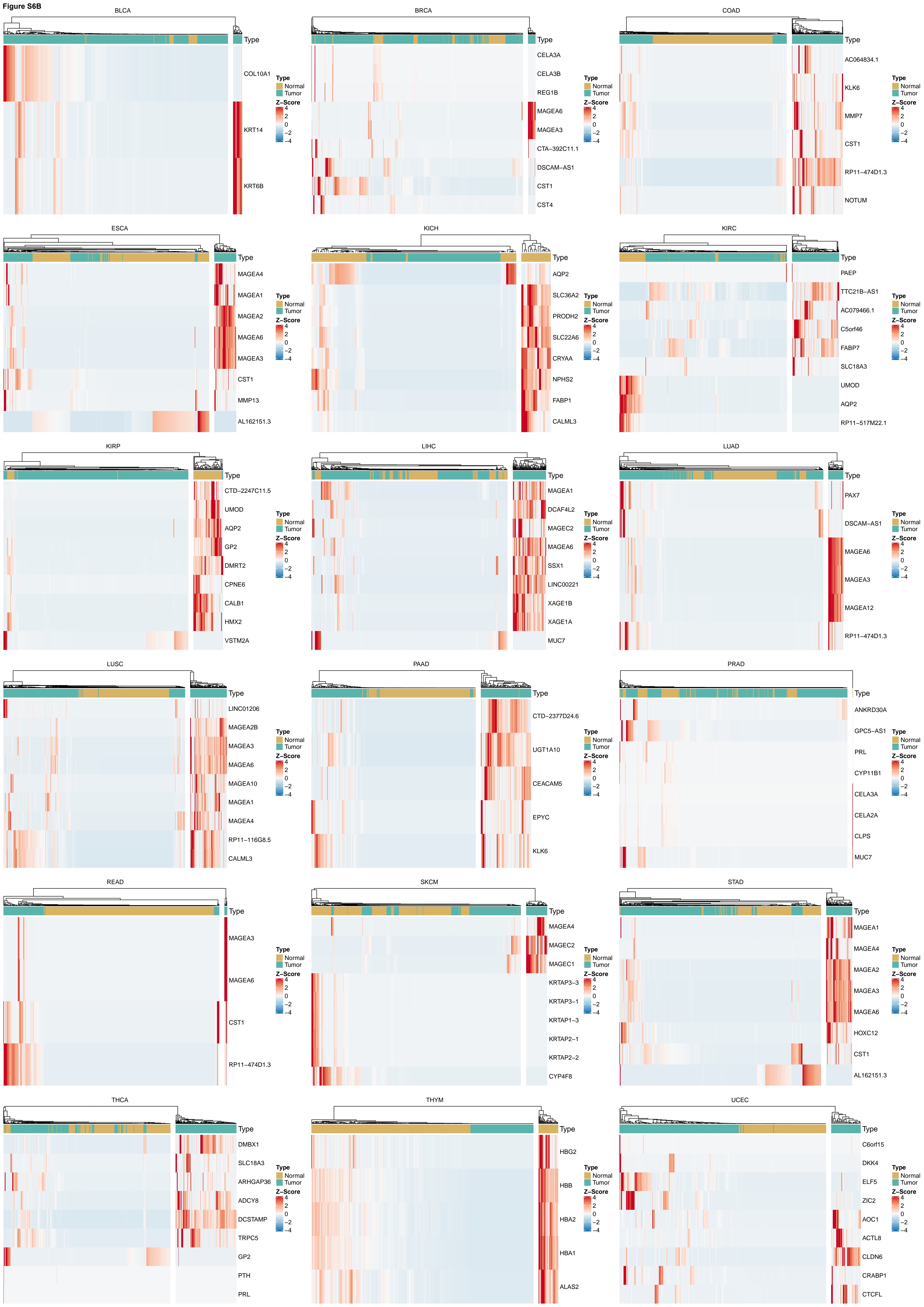


UCEC nodes:1022 edges:2582

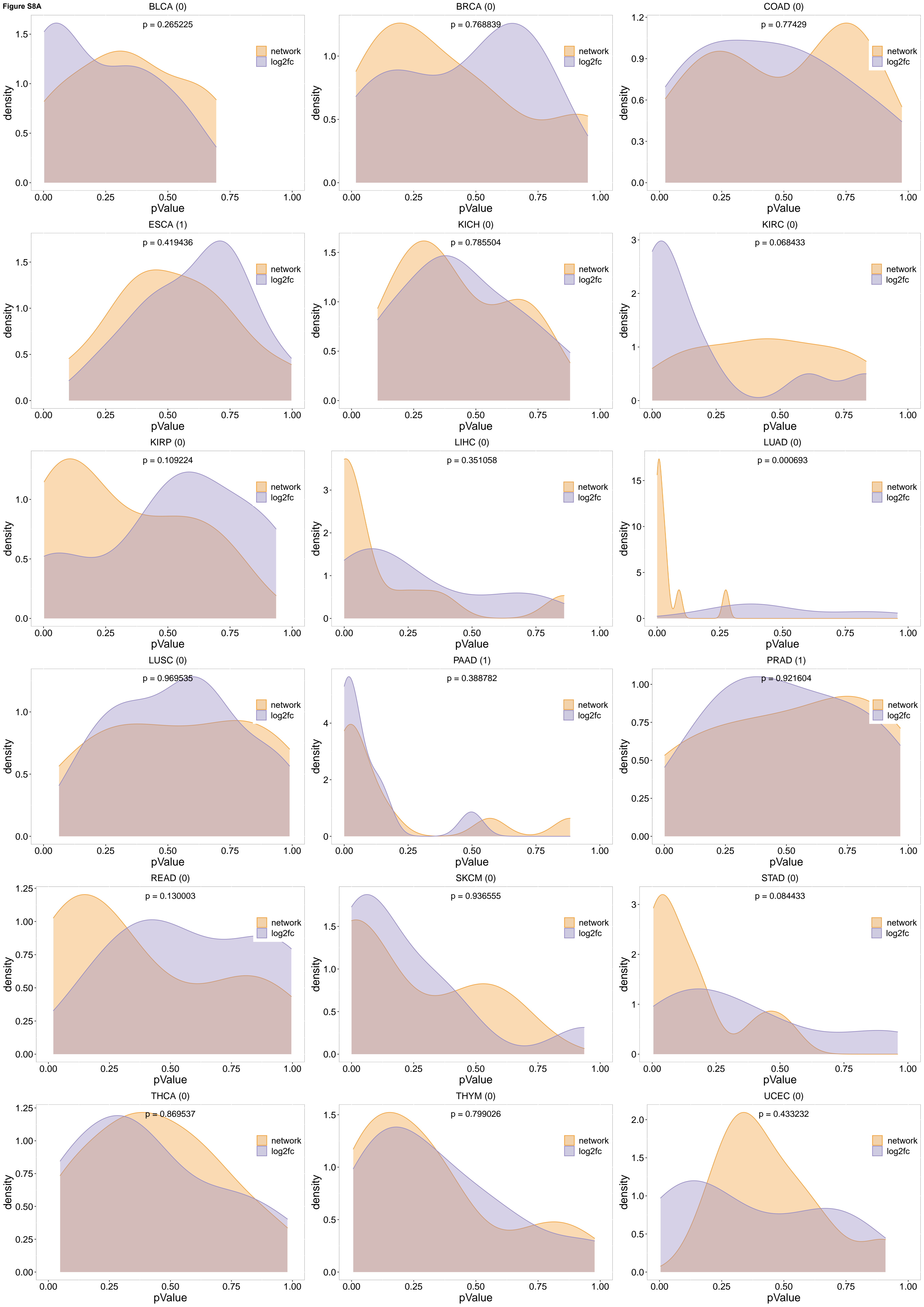


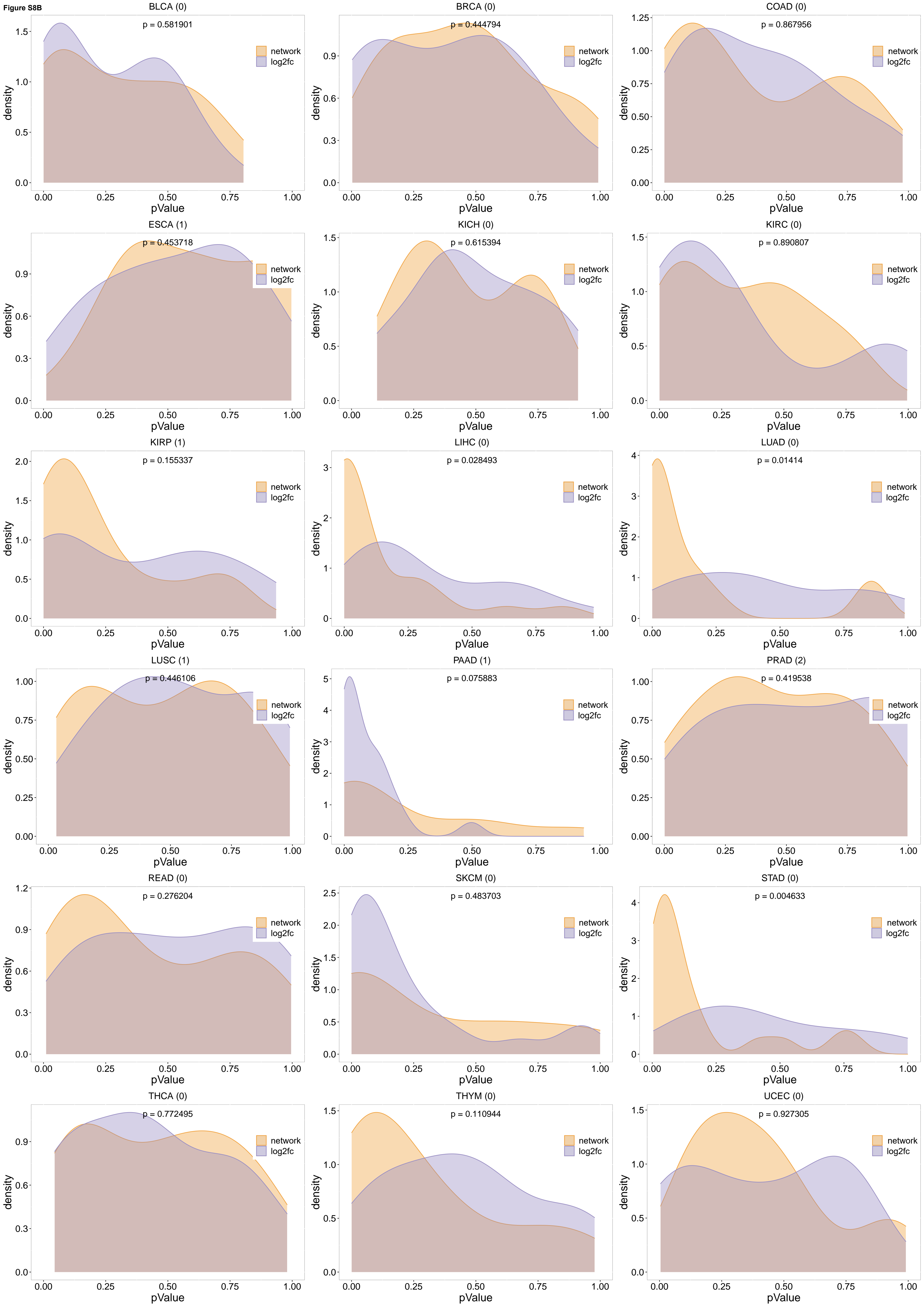


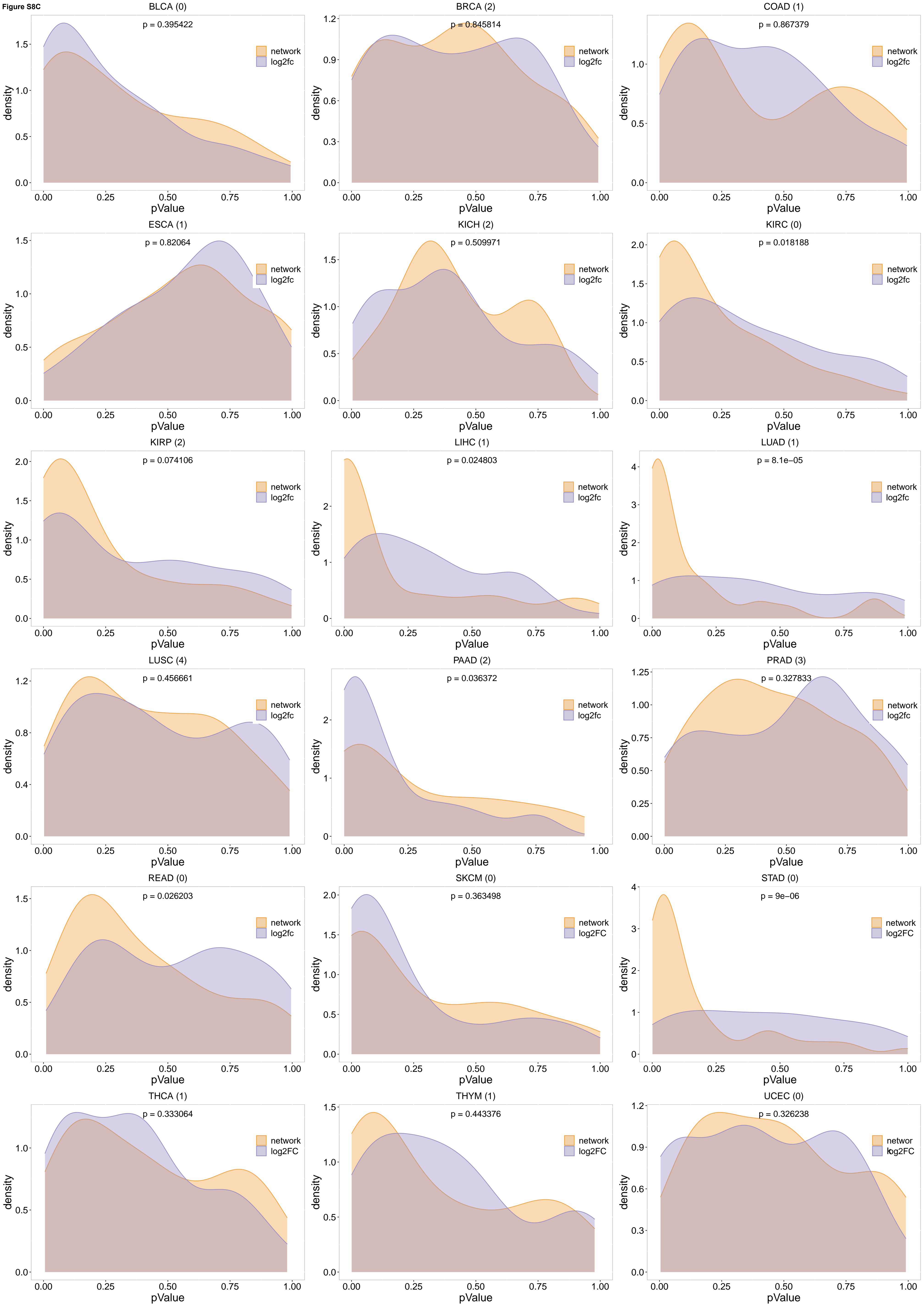


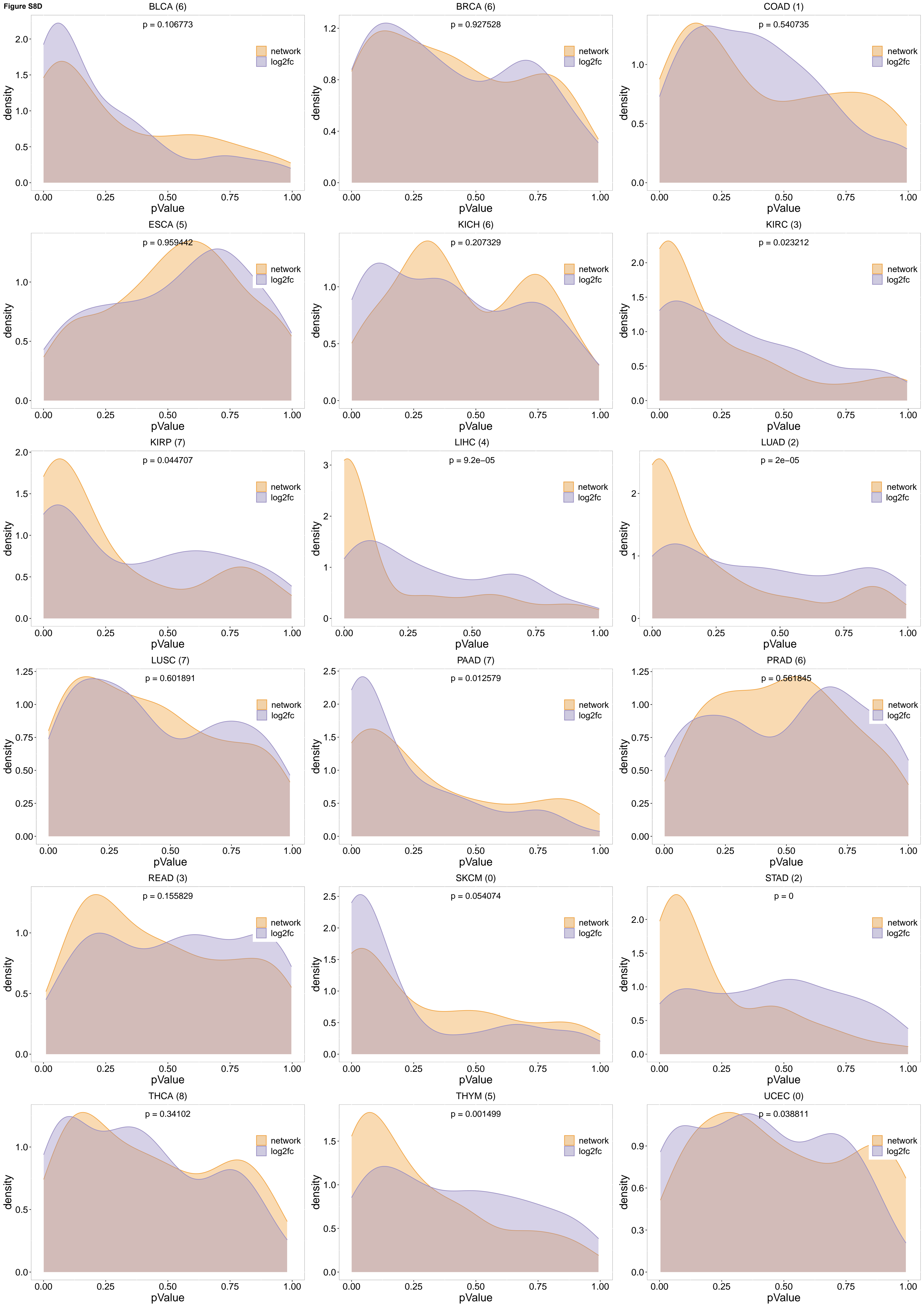


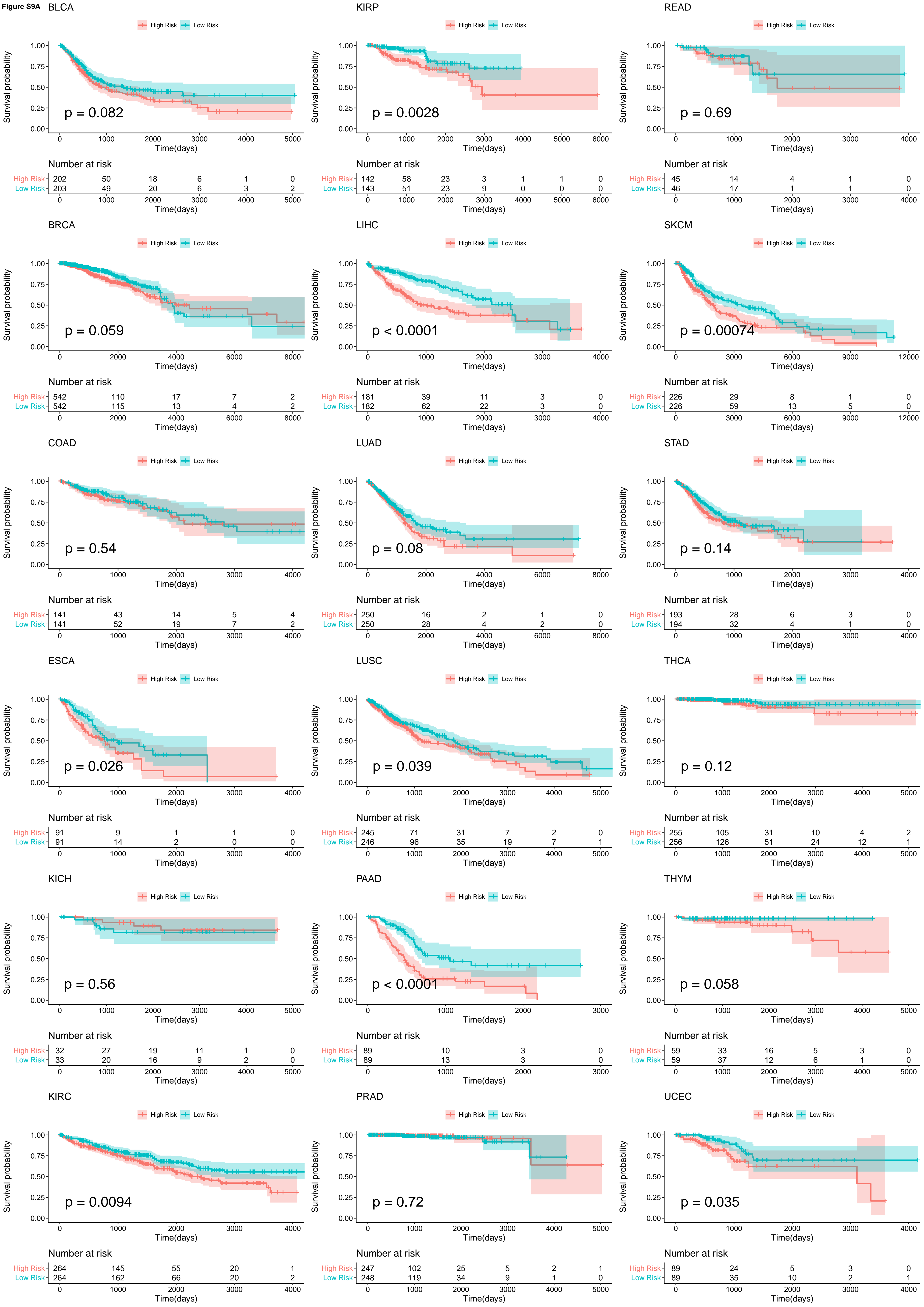












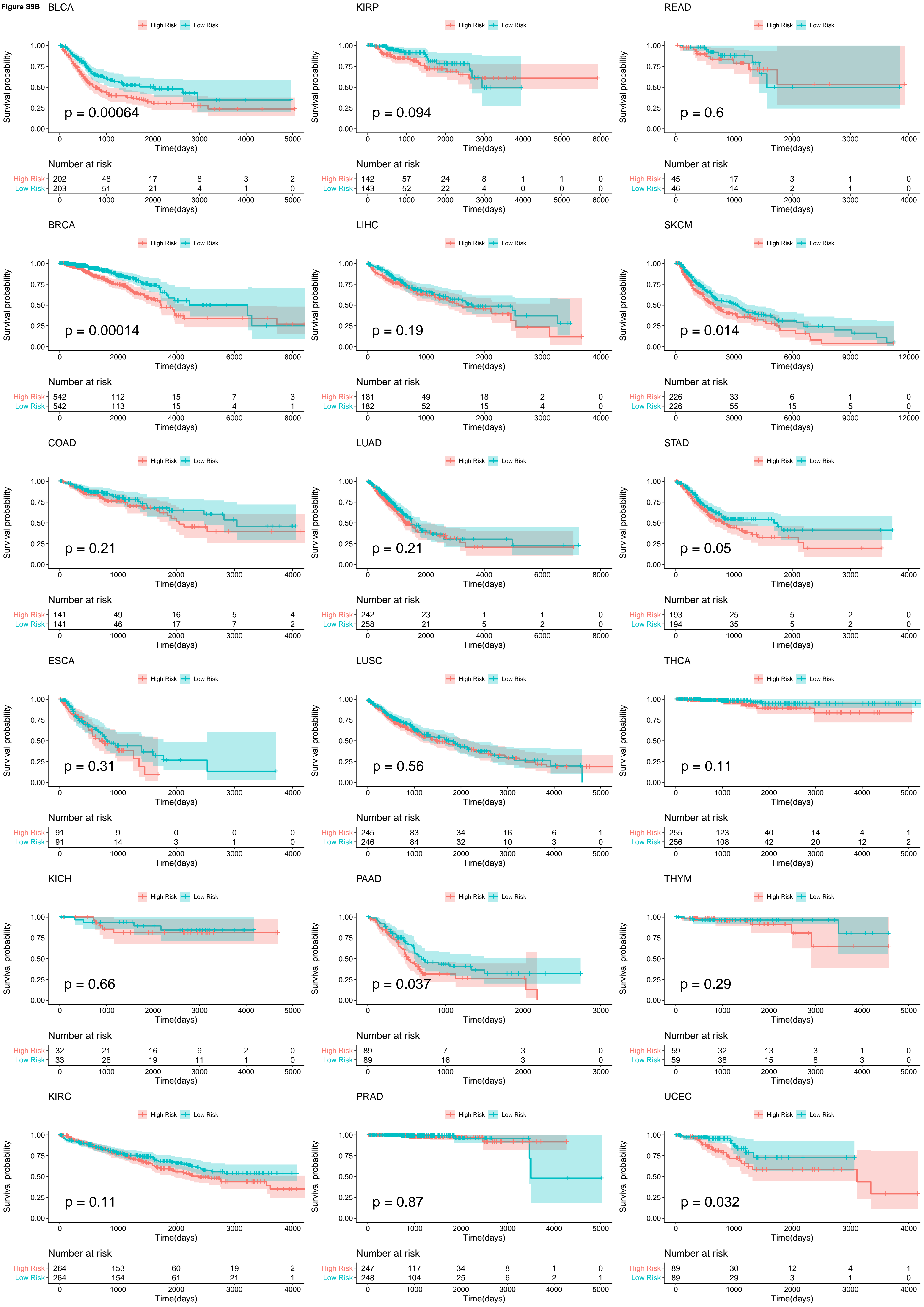


Figure S10

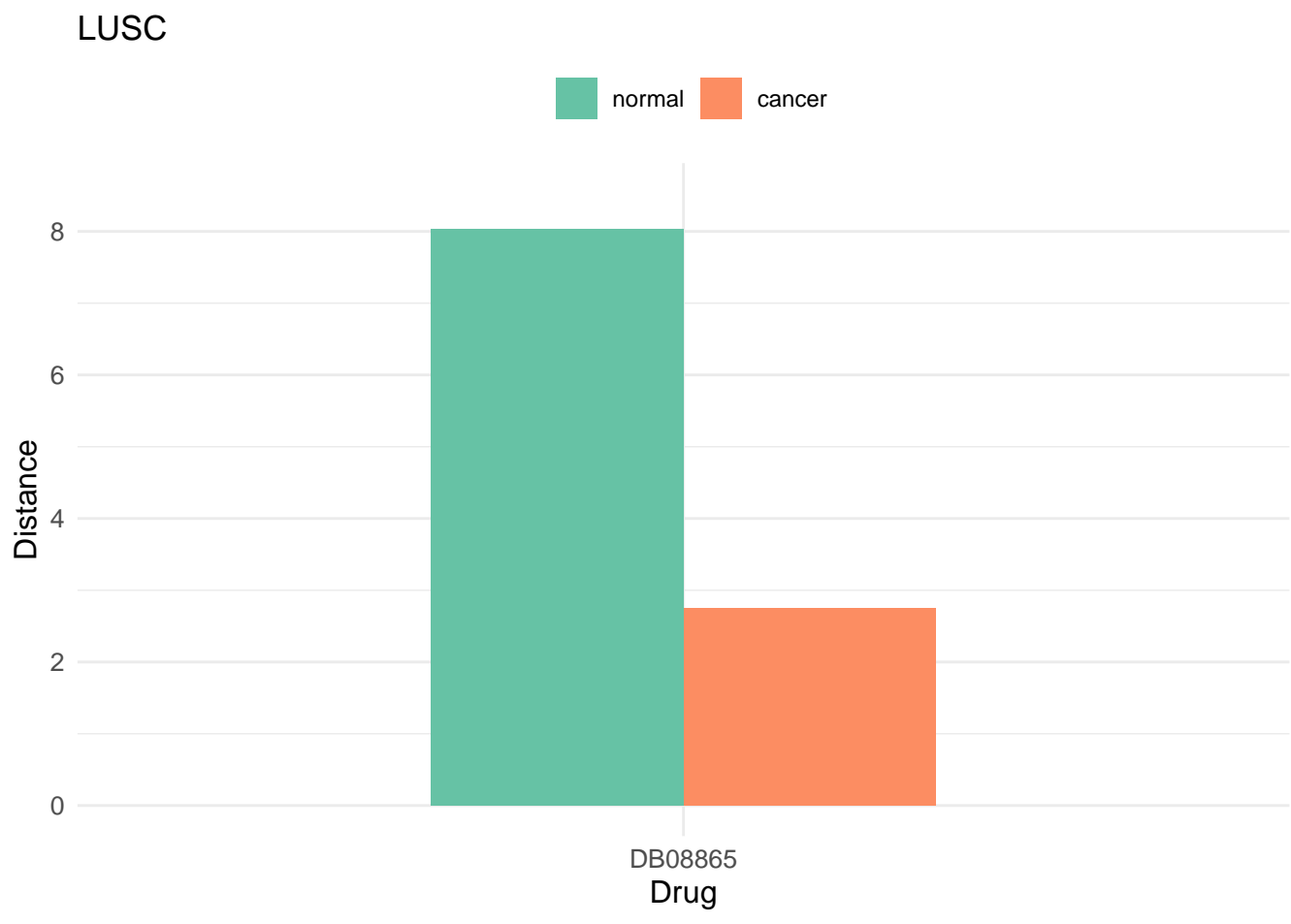
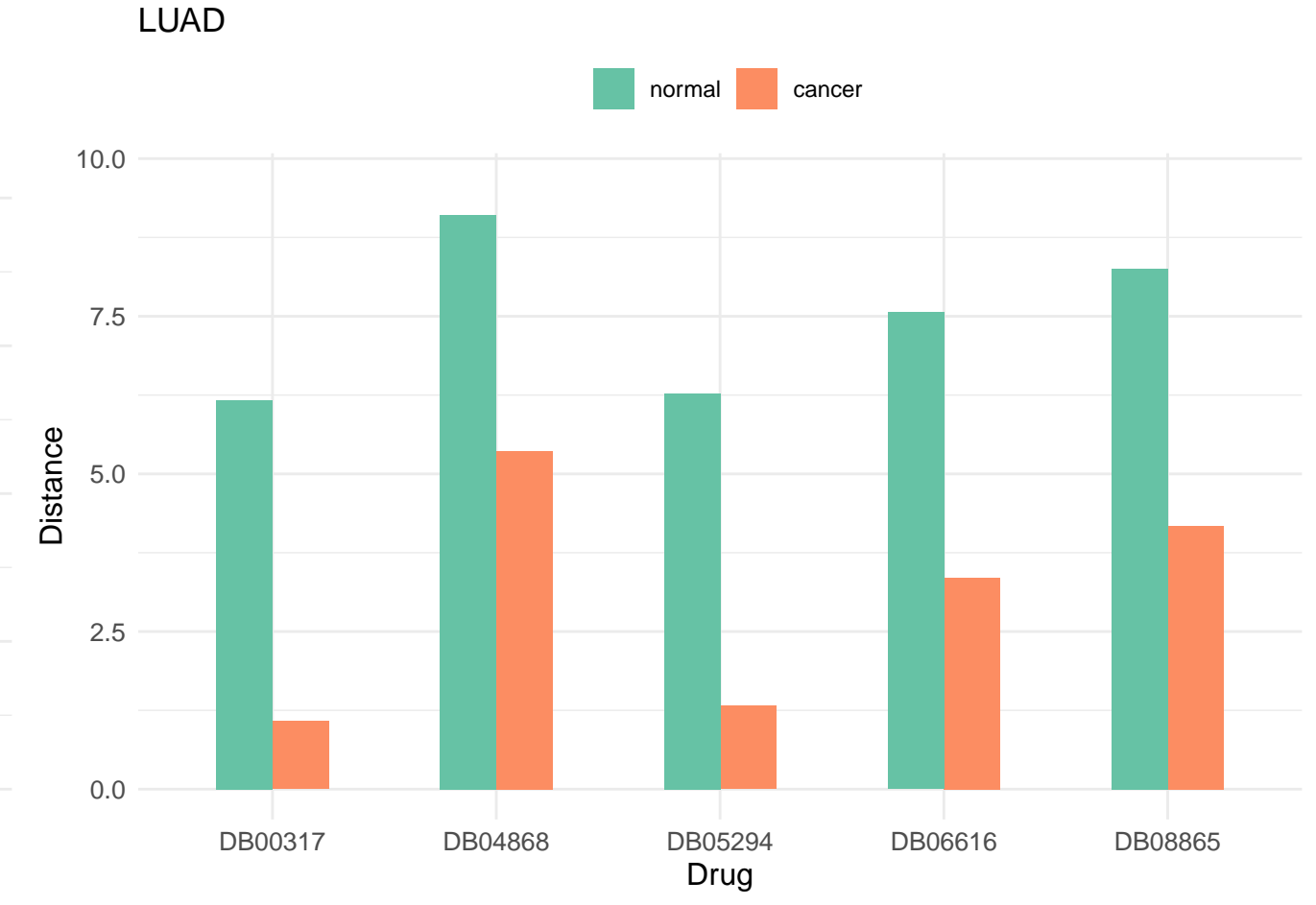
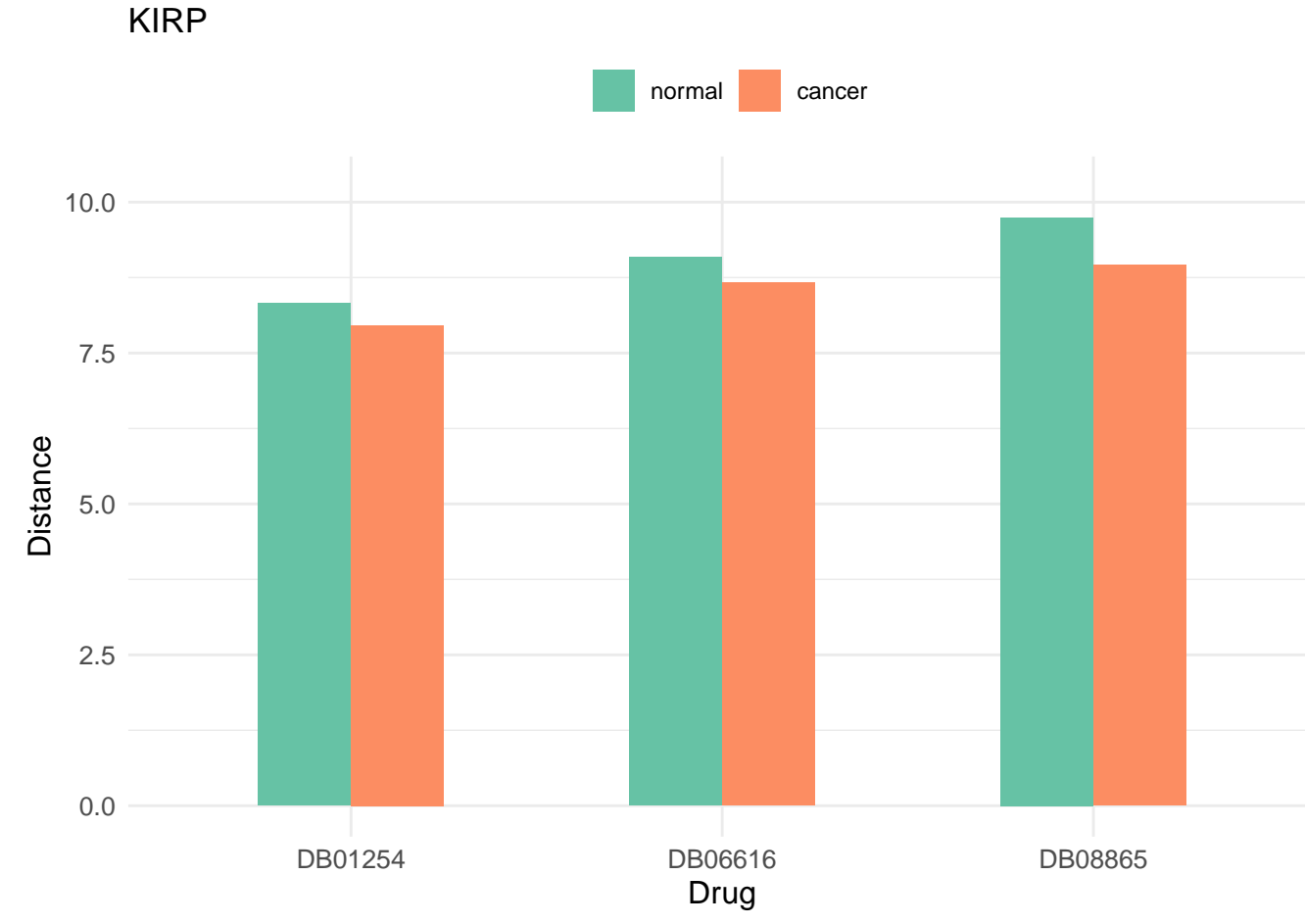
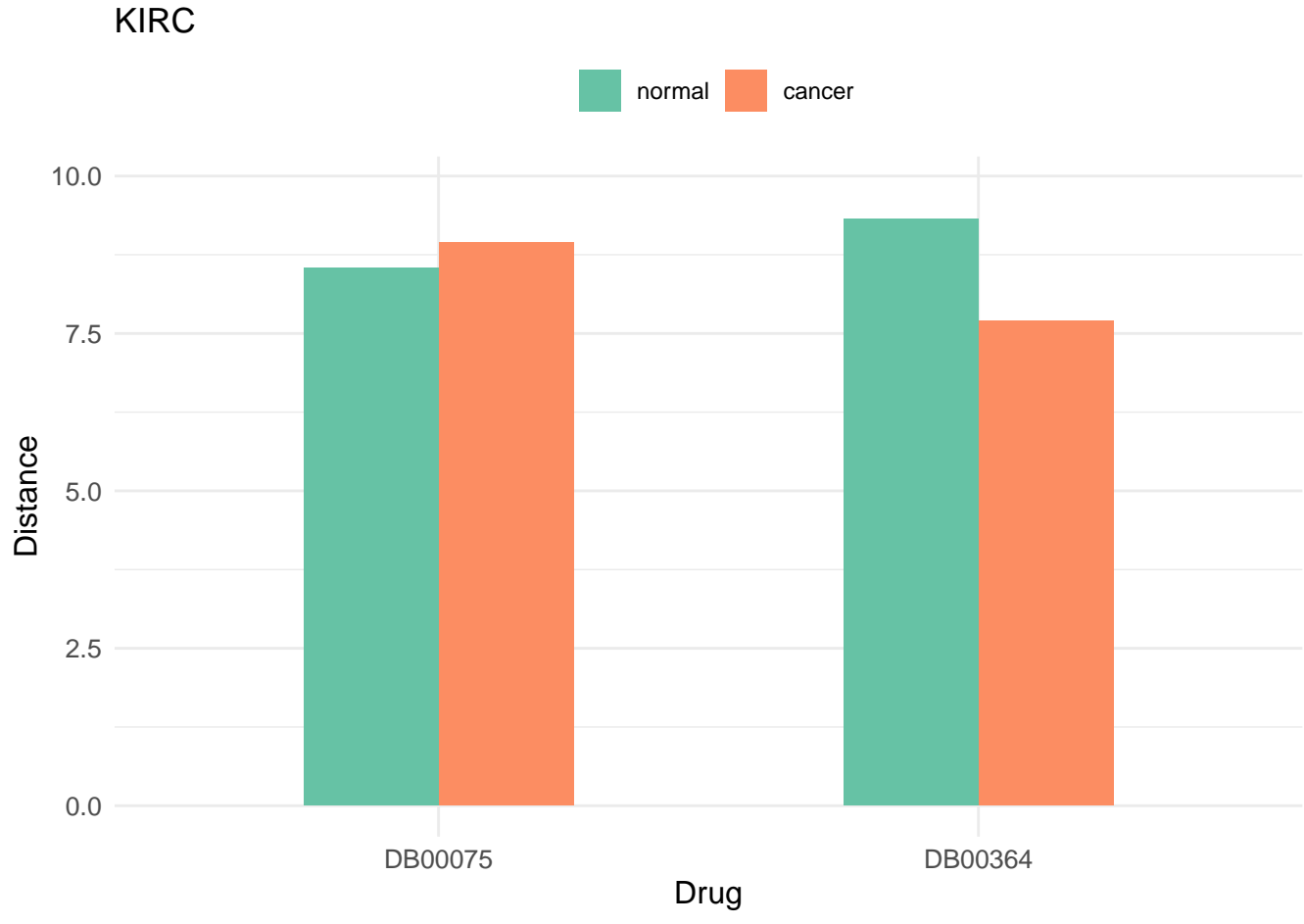
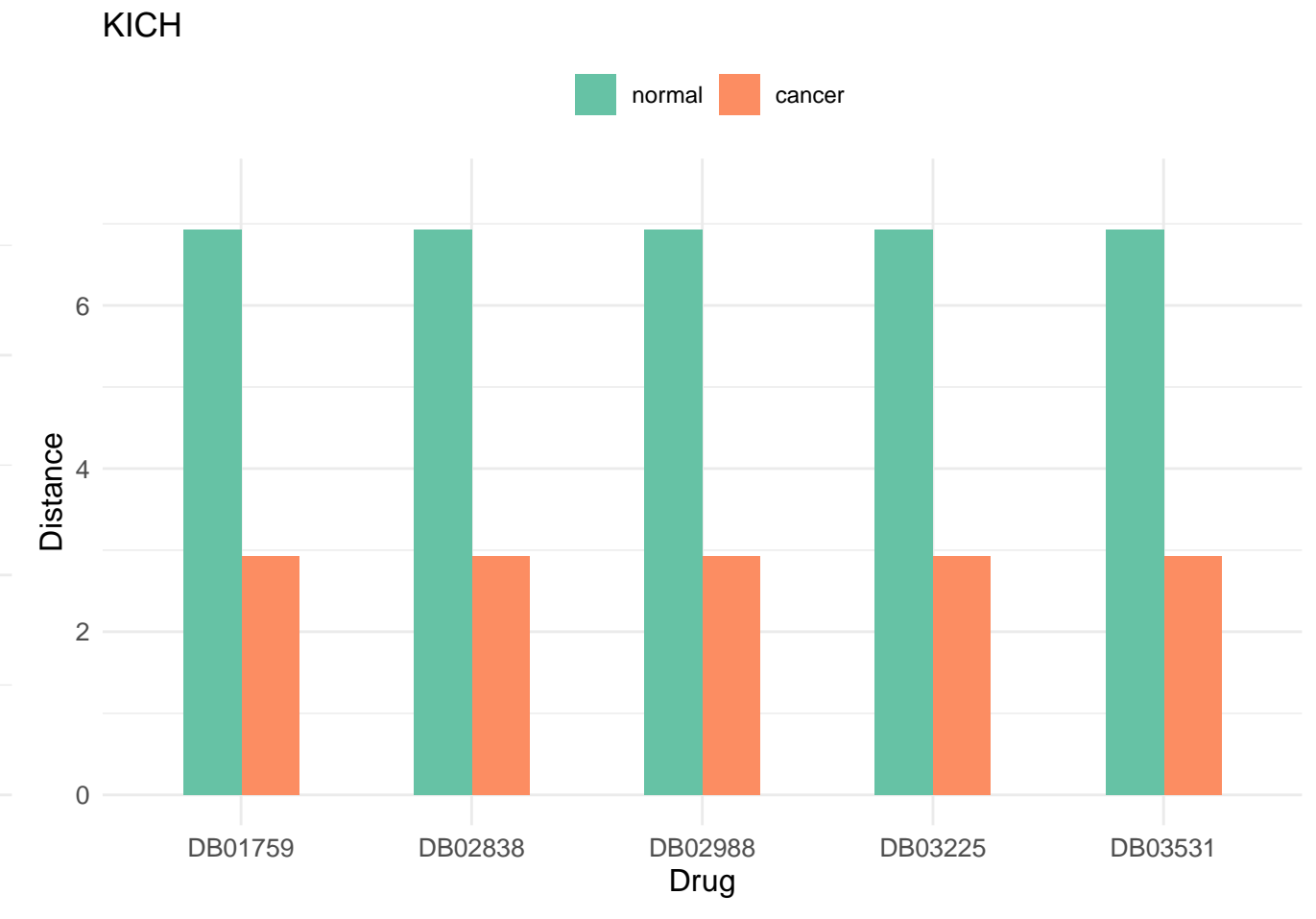
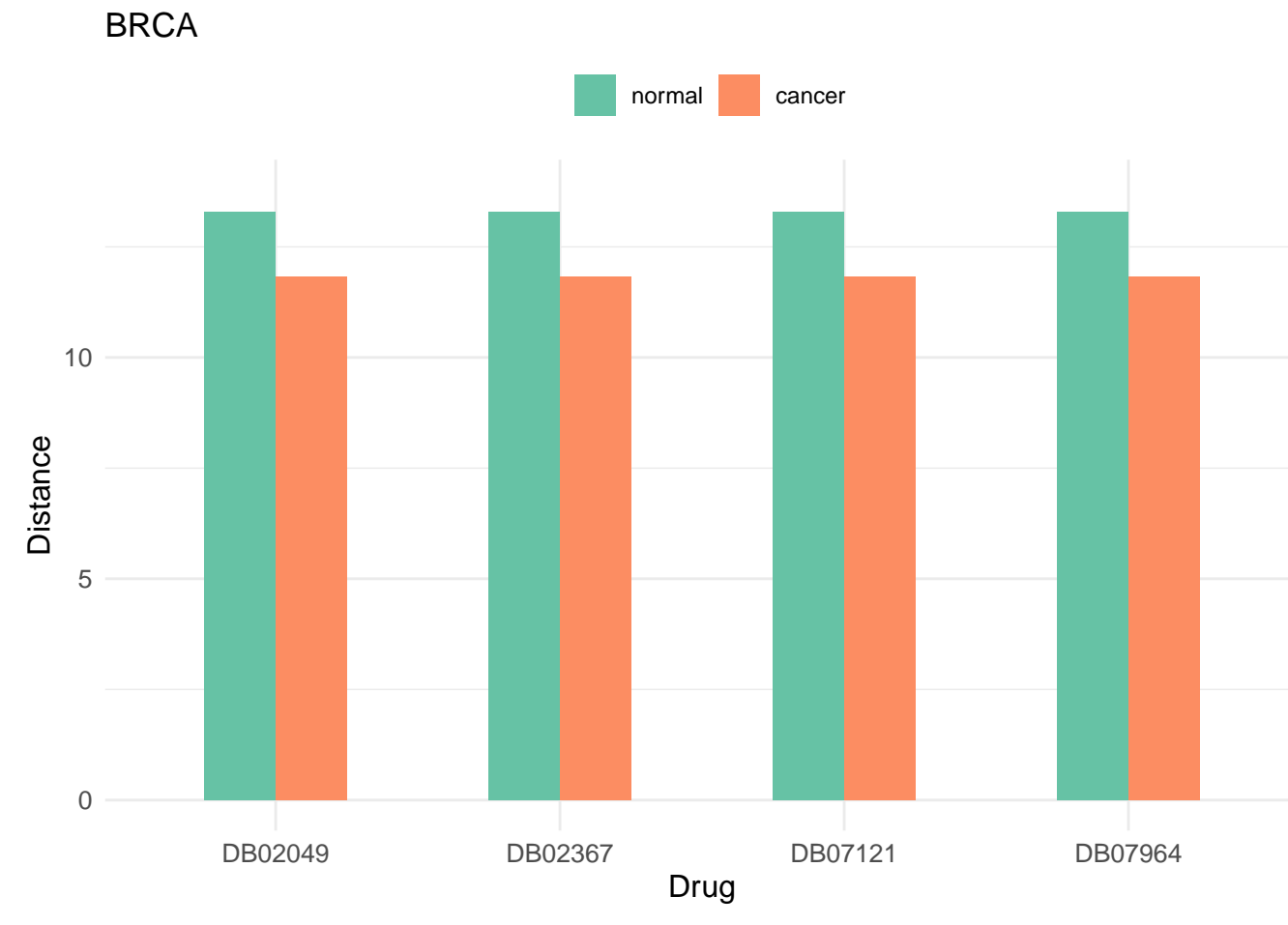
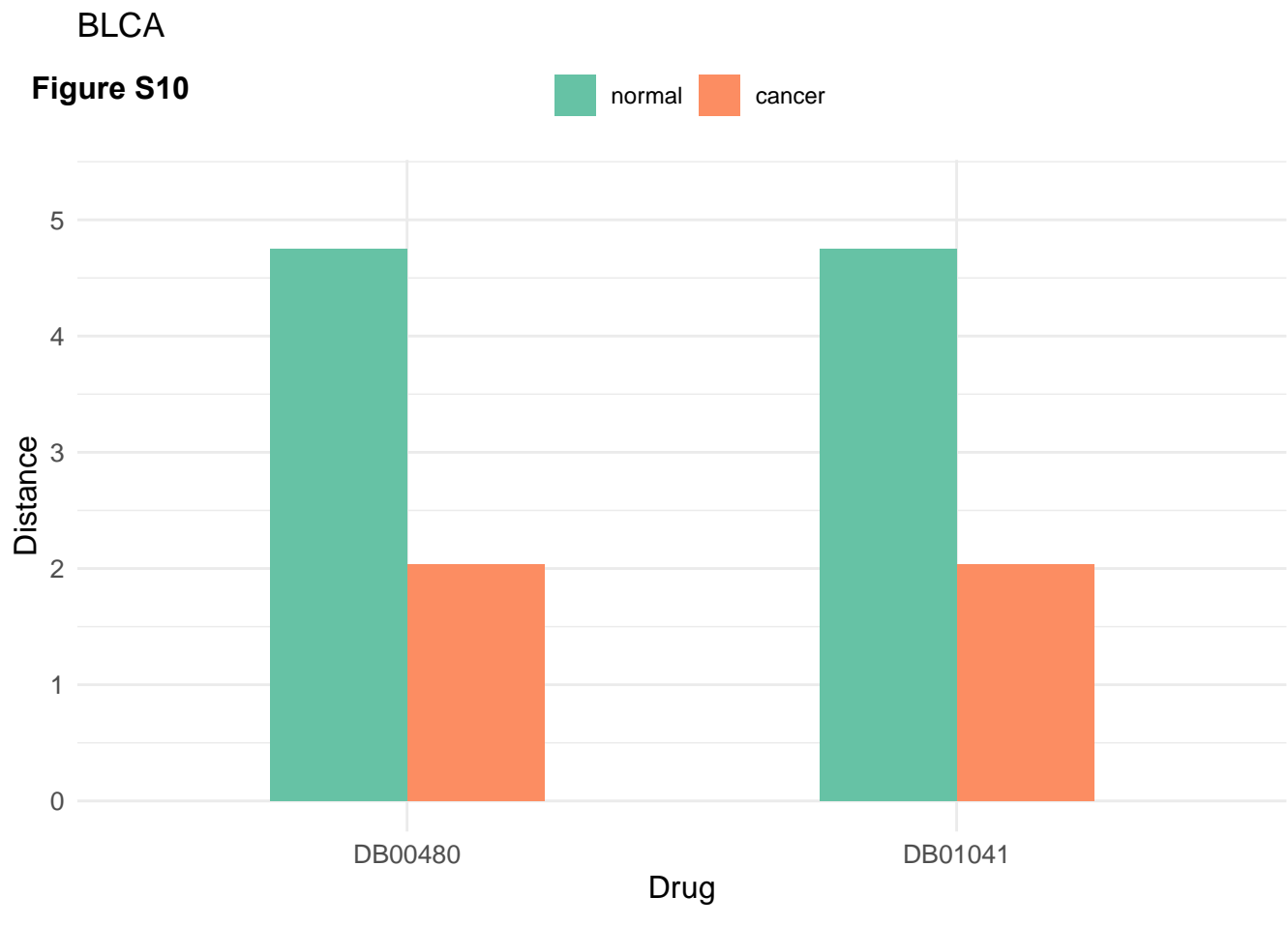


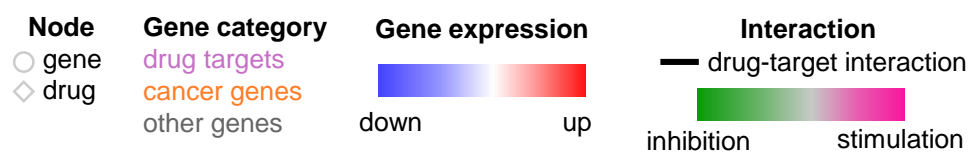
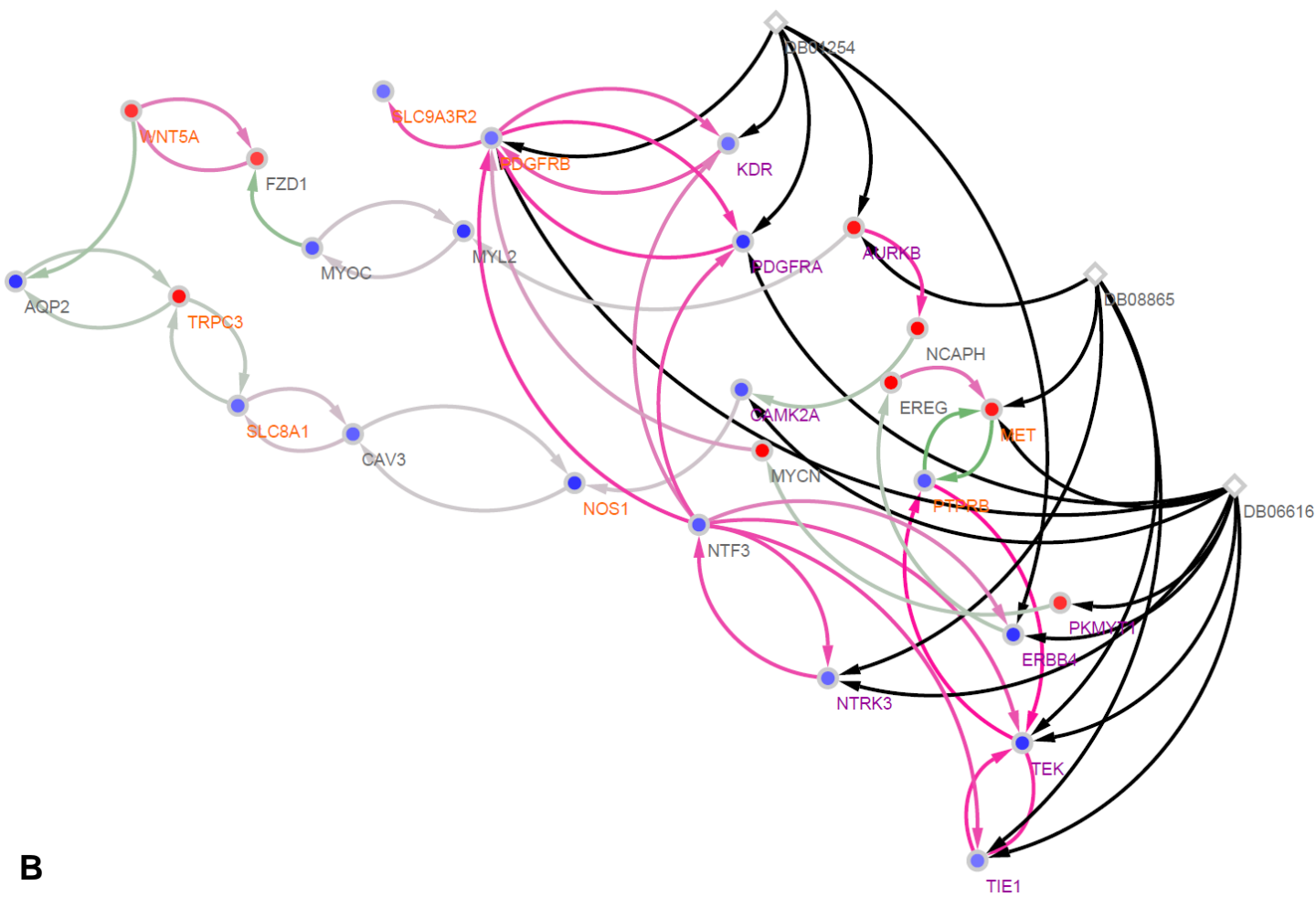
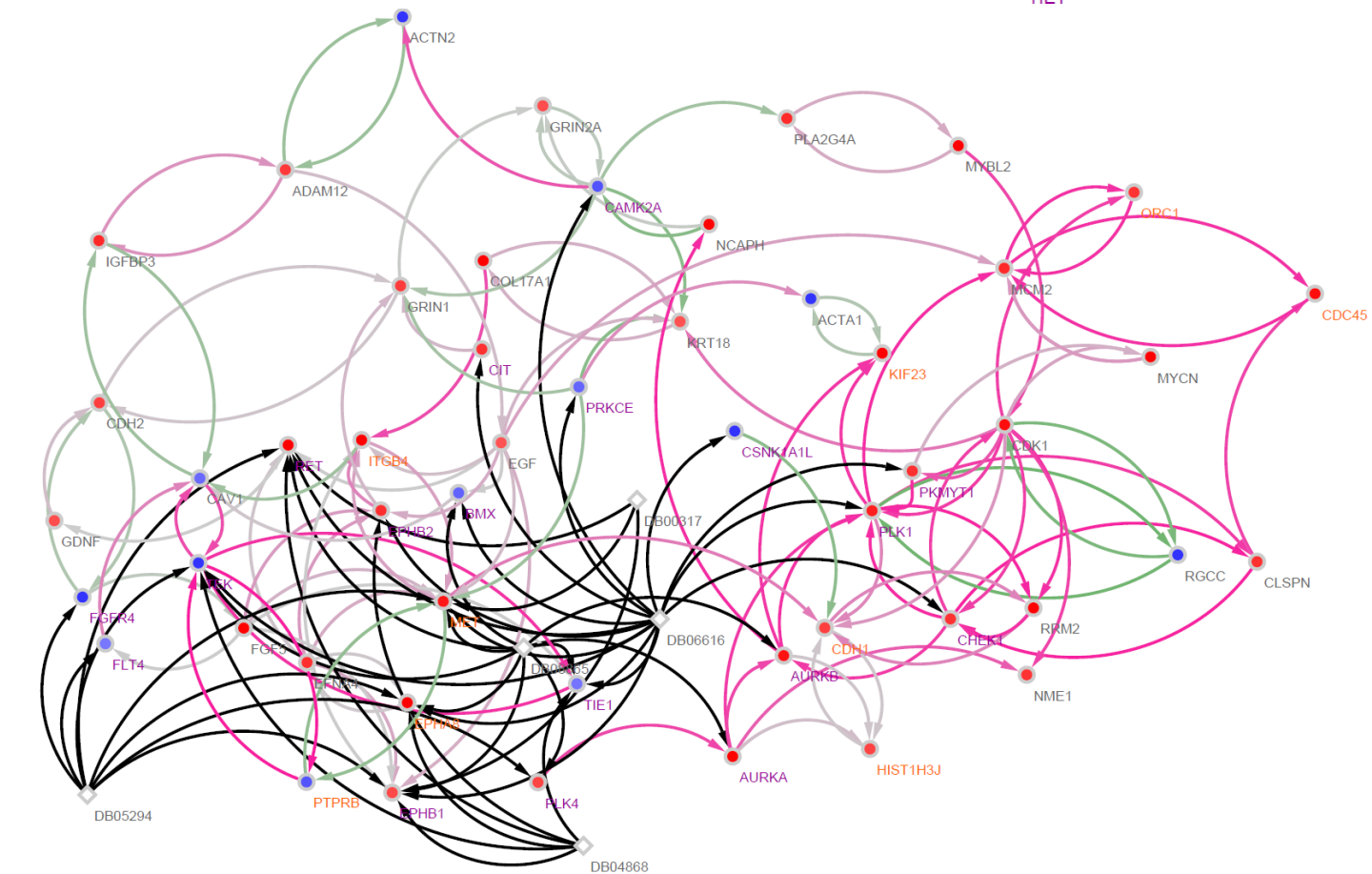
Figure S11**A****B**

Figure S12

