

# **An integrative somatic mutation analysis to identify pathways linked with survival outcomes across 19 cancer types**

**Sunho Park<sup>1</sup>, Seung-Jun Kim<sup>2</sup>, Donghyeon Yu<sup>1</sup>, Samuel Pena-Llopis<sup>3,4</sup>, Jianjiong Gao<sup>5</sup>, Jin Suk Park<sup>1</sup>, Hao Tang<sup>1</sup>, Beibei Chen<sup>1</sup>, Jiwoong Kim<sup>1</sup>, Jessie Norris<sup>1</sup>, Xinlei Wang<sup>6</sup>, Min Chen<sup>7</sup>, Minsoo Kim<sup>1</sup>, Jeongsik Yong<sup>8</sup>, Zabi WarDak<sup>4,9</sup>, Kevin S Choe<sup>4,9</sup>, Michael Story<sup>4,9</sup>, Timothy K. Starr<sup>10,11</sup>, Jaeho Cheong<sup>†12</sup>, Tae Hyun Hwang<sup>†1,4</sup>**

<sup>1</sup>Department of Clinical Sciences, University of Texas Southwestern Medical Center, Dallas, Texas, United States of America

<sup>2</sup>Department of Computer Science and Electrical Engineering, University of Maryland at Baltimore, Baltimore County, Maryland, United States of America

<sup>3</sup>Department of Internal Medicine, University of Texas Southwestern Medical Center, Dallas, Texas, United States of America

<sup>4</sup>Simmons Comprehensive Cancer Center, University of Texas Southwestern Medical Center, Dallas, Texas, USA.

<sup>5</sup>Center for Molecular Biology, Memorial Sloan Kettering Cancer Center, New York, New York, 10065, United States of America

<sup>6</sup>Department of Statistical Science, Southern Methodist University, Dallas, Texas, United States of America

<sup>7</sup>Department of Mathematical Sciences, University of Texas at Dallas, Dallas, Texas, United States of America

<sup>8</sup>Department of Biochemistry, Molecular Biology and Biophysics, University of Minnesota Twin Cities, Minneapolis, Minnesota, United States of America

<sup>9</sup>Department of Radiation Oncology, University of Texas Southwestern Medical Center, Dallas, Texas, United States of America

<sup>10</sup>Masonic Cancer Center, University of Minnesota – Twin Cities, Minneapolis, Minnesota, United States of America

<sup>11</sup>Department of Obstetrics, Gynecology & Women's Health and Department of Genetics, Cell Biology and Development, University of Minnesota, Minneapolis, Minnesota, United States of America

<sup>12</sup>Department of Surgery, Yonsei University College of Medicine, Seoul, Korea

<sup>†</sup>Corresponding author

Email addresses:

SP: [sunho.park@utsouthwestern.edu](mailto:sunho.park@utsouthwestern.edu)

SK: [sjkim@umbc.edu](mailto:sjkim@umbc.edu)

DY: [donghyeon.yu@utsouthwestern.edu](mailto:donghyeon.yu@utsouthwestern.edu)

SPL: [samuel.pena-llopis@utsouthwestern.edu](mailto:samuel.pena-llopis@utsouthwestern.edu)

JG: [jgao@cbio.mskcc.org](mailto:jgao@cbio.mskcc.org)

JP: [jinsuk.park@utsouthwestern.edu](mailto:jinsuk.park@utsouthwestern.edu)

HT: [hao.tang@utsouthwestern.edu](mailto:hao.tang@utsouthwestern.edu)

BC: [beibei.chen@utsouthwestern.edu](mailto:beibei.chen@utsouthwestern.edu)

JK: [jiwoong.kim@utsouthwestern.edu](mailto:jiwoong.kim@utsouthwestern.edu)

JN: [jessie.norris@utsouthwestern.edu](mailto:jessie.norris@utsouthwestern.edu)

XW: [swang@mail.smu.edu](mailto:swang@mail.smu.edu)

MC: [mchen@utdallas.edu](mailto:mchen@utdallas.edu)

MK: [mins.kim@utsouthwestern.edu](mailto:mins.kim@utsouthwestern.edu)

JY: [jyong@umn.edu](mailto: jyong@umn.edu)

ZW: [zabi.wardak@utsouthwestern.edu](mailto:zabi.wardak@utsouthwestern.edu)

KC: [kevin.choe@utsouthwestern.edu](mailto:kevin.choe@utsouthwestern.edu)

MS: [michael.story@utsouthwestern.edu](mailto:michael.story@utsouthwestern.edu)

TS: [star0044@umn.edu](mailto:star0044@umn.edu)

JC: [jhcheong@yuhs.ac](mailto:jhcheong@yuhs.ac)

THH: [taehyun.hwang@utsouthwestern.edu](mailto:taehyun.hwang@utsouthwestern.edu)

## **Abstract**

Identification of altered pathways that are clinically relevant across human cancers is a key challenge in cancer genomics. We developed a network-based algorithm to integrate somatic mutation data with gene networks and pathways, in order to identify pathways altered by somatic mutations across cancers. We applied our approach to The Cancer Genome Atlas (TCGA) dataset of somatic mutations in 4,790 cancer patients with 19 different types of malignancies. Our analysis identified cancer-type-specific altered pathways enriched with known cancer-relevant genes and drug targets. Consensus clustering using gene expression datasets that included 4,870 patients from TCGA and multiple independent cohorts confirmed that the altered pathways could be used to stratify patients into subgroups with significantly different clinical outcomes. Of particular significance, certain patient subpopulations with poor prognosis were identified because they had specific altered pathways for which there are available targeted therapies. These findings could be used to tailor and intensify therapy in these patients, for whom current therapy is suboptimal.

## **Background**

In the last few years, studies using high-throughput technologies have highlighted the fact that the development and progression of cancer hinges on somatic alterations. These somatic alterations may disrupt gene function, such as activating oncogenes or inactivating tumor suppressor genes, and thus, dysregulate critical pathways contributing to tumorigenesis. Therefore, precise identification and understanding of disrupted pathways may provide insights into therapeutic strategies and the development of novel agents. Many large-scale cancer genomics studies, such as The Cancer Genome Atlas (TCGA) and the International Cancer Genome Consortium (ICGC), have performed an integrated analysis to draft an overview of somatic alterations in the cancer genome<sup>1-4</sup>. Many of these studies have reported novel

candidate cancer genes mutated at high and intermediate frequencies in a specific cancer as well as across many cancer types<sup>4</sup>. However, it is a still challenge to translate somatic mutations in tumors into the pathway model to accurately predict patient clinical outcomes<sup>5,6</sup>. Recently, in order to improve the clinical relevance and utility of somatic mutation analyses, Hopfree et al<sup>7</sup> proposed integrating somatic mutation data with molecular interaction networks for patient stratification. They demonstrated that inclusion of prior knowledge, captured in molecular interaction networks, could improve identification of patient subgroups with significantly different histological, pathological or clinical outcomes and discover novel cancer-related pathways or subnetworks. In a similar manner, other network-based methods have demonstrated that incorporating molecular networks and/or biological pathways can improve accuracy in identifying cancer-related pathways<sup>8-11</sup>.

One limitation of these network-based methods is that they are not designed to fully utilize large-scale somatic mutation data from multiple cancer types to determine which particular pathways are altered by somatic mutations across a range of human cancers. In addition, due to the incomplete knowledge of existing gene set and/or pathway database, these methods are limited to detect pathways based on a number of altered genes annotated in existing gene set and pathway databases. Alternatively, the methods that build pathways *de novo* without incorporating biological prior knowledge can be applicable to detect altered pathways, but these methods were not designed to detect cancer-type specific or commonly altered pathways either.

To address these, we developed an algorithm named NTriPath (Network regularized non-negative TRI matrix factorization for PATHway identification) to integrate somatic mutation, gene-gene interaction networks and gene set or pathway databases to discover pathways altered by somatic mutations in 4,790 cancer patients with 19

different types of cancers. Incorporating existing gene set or pathway databases enables NTriPath to report a list of altered pathways across cancers, and make easy to determine/compare which particular pathways are altered in a particular cancer type(s). In particular, the use of large-scale genome-wide somatic mutations from 4,790 cancer patients and gene-gene interaction networks enables NTriPath to classify genes, which were not annotated in existing gene set or pathway databases, as new member genes of the identified altered pathways based on modular structures of mutational data within a cancer type and/or across multiple cancer types (using matrix factorization) and connectivity in the gene-gene interaction networks. The questions that we investigate here are: first, whether large-scale integrative somatic mutation analysis can reliably identify cancer-type-specific or commonly altered pathways by somatic mutations across cancers; second, whether the identified pathways can be used as a prognostic biomarker for patient stratification - with the assumption that the altered pathways contribute to cancer development and progression and, thus, impact survival.

In these experiments, we demonstrated that the cancer-type-specific and commonly altered pathways identified by NTriPath are biologically relevant to the corresponding cancer type and are associated with patient survival outcomes. We also showed that cancer-specific altered pathways are enriched with many known cancer-relevant genes and targets of available drugs including those already FDA-approved. These results imply that the cancer-specific altered pathways can guide therapeutic strategy to target the altered pathways that are pivotal in each cancer type.

## Results

### **NTriPath: An integrative somatic mutation analysis for discovering pathways altered by somatic mutations across multiple cancer types**

NTriPath integrates somatic mutations with gene-gene interaction networks and a pathway database to discover altered pathways across human cancers. We collected somatic mutation data from TCGA for 4,790 patients and 19 different cancer types (Table 1). A diagrammatic description of our algorithm is depicted in Figure 1. Four types of data were used as input for our algorithm. First, we generated a binary matrix ( $\mathbf{X}$ ) of patients x genes, with '1' indicating a mutation and '0' no mutation. Second, we constructed gene-gene interaction networks ( $\mathbf{A}$ ). Third, we incorporated a pathway database ( $\mathbf{V}_0$ ) (e.g., conserved 4,620 subnetworks across species<sup>12</sup>). Fourth, we included clinical data on the patient's tumor type ( $\mathbf{U}$ ). NTripPath produces two matrices as output; 1) altered pathways by mutated genes ( $\mathbf{V}$ ) and 2) altered pathways by cancer type matrix ( $\mathbf{S}$ ). The use of both large-scale somatic mutation profiles and gene-gene interaction networks enabled NTripPath to identify cancer-related pathways containing known cancer genes mutated at different frequencies across cancers with newly added member genes according to high network connectivity ( $\mathbf{V}$ )<sup>8, 13</sup>. Finally we use the altered pathways by cancer type matrix ( $\mathbf{S}$ ) to identify altered pathways that are specific for each cancer type. For further details, please see the Materials and Methods section. Our method is also available at [www.taehyunlab.org](http://www.taehyunlab.org).

### **NTriPath identifies cancer-type-specific altered pathways that are biologically and clinically relevant**

In each cancer type, we selected the top 3 ranked altered pathways by statistical significance from NTripPath with the 4,620 subnetwork modules to generate cancer-type-specific altered pathways (See Material and Method section and Supplementary Table 1). Interestingly, NTripPath was able to find altered pathways containing not only genes that were frequently mutated but also genes that were

mutated in a small subset of patients in each cancer type (Supplementary Table 2 and Supplementary Figure 1). Gene set enrichment analysis using the genes from the top 3 altered pathways showed that the altered pathways are significantly enriched with well-known cancer-related genes from COSMIC database<sup>14</sup> and known drug target genes as well as cancer-relevant biological processes (Supplementary Table 3 and 4).

Focusing on kidney renal clear cell carcinoma (KIRC) as a proof of concept, NTriPath identified the pathway consisting of *VHL*, *USP33*, *DIO2*, *TCEB1* and *TCEB2* as the top-ranked altered pathway in KIRC (Figure 2A). The *VHL* (*von-Hippel Lindau*) gene is a well-known tumor suppressor associated with KIRC, and is frequently mutated in patients with KIRC<sup>15-18</sup>. *VHL* was the most frequently mutated gene in TCGA KIRC with 55.7% of patients harboring mutations in the gene. *TCEB1* is mutated at very low frequency in TCGA KIRC cohort. A recent study found that *TCEB1* is mutated in about 3% of the KIRC patients without *VHL* inactivation, and found *TCEB1* preventing the binding of Elongin C to *VHL*, which inactivates the *VHL* pathway<sup>16</sup>. The second highest ranked pathway contained *EP300* and *TP53*. *EP300* and *TP53* were mutated in 8.1% and 5.2% of patients, respectively. *EP300* has been identified as a co-activator of hypoxia-inducible factor 1 alpha (*HIF1* $\alpha$ ), whose activation is a hallmark of KIRC tumors. *TP53* was previously found to be associated with poor outcome in TCGA KIRC<sup>19</sup>. The third highest ranked pathway contains *LRP1* and matrix metalloproteinases (*MMP1*, *MMP7*, *MMP9*, *MMP26*). *LRP1* is mutated in 10% of TCGA KIRC cohort, but matrix metalloproteinases (*MMPs*) were not mutated in TCGA KIRC cohort. Biological and clinical relevance of *LRP1* mutation in KIRC has not been previously reported. *MMPs* have been implicated in different types of cancer progression including the acquisition of invasive and metastatic properties in many cancer types. The aberrant expression of *MMPs* has

been associated with poor patient survival and prognosis in KIRC patients<sup>8, 20</sup>. Interestingly, recent studies suggested that *LRPI* induces the expression of matrix metalloproteinase (MMPs) and thus promotes cancer cell invasion and metastasis in many cancers including KIRC<sup>16,21-23</sup>.

NTriPath identified many new member genes in the top ranked pathways including *TCEB2*, *JUN*, and *SPI* as well as other tumor suppressors such as *CREBBP*, *SMAD3*, *BRCA1* and *RBI*. These newly identified member genes by NTriPath were mutated at a very low frequency or not mutated at all in TCGA KIRC patients. Instead, these genes interacted with many frequently mutated genes in the networks and were often dysregulated at the mRNA and protein levels in many KIRC patients (Figure 2B). For example, *TCEB2*, *SPI* and *JUN* were not mutated but yet their expression was dysregulated in 7%, 10% and 2% of TCGA KIRC patients, respectively. Previous studies have shown that dysregulation in *TCEB2* is expected to disrupt the protein complex that ubiquitinates *HIF1 $\alpha$* , resulting in the same phenotype as *VHL* inactivation by mutation or promoter hypermethylation<sup>24-26</sup>. In addition, *SPI* and *JUN* were previously identified as major transcriptional regulators associated with signaling circuit to promote tumor growth and invasion in KIRC<sup>18</sup>.

Taken together, these results demonstrate that NTriPath is an effective tool to accurately identify cancer-specific altered pathways including known cancer genes mutated at a high or intermediate frequency in the patients, as well as genes mutated at a very low frequency or not mutated at all yet may be fundamental role in development and/or progression of KIRC.

### **Cancer-type-specific altered pathways across cancer types correlate with patient survival outcomes**

We hypothesized that cancer-type-specific altered pathways reflect the molecular basis underlying the patient clinical outcomes. This would allow us to use member



genes in the altered pathways as gene signatures to stratify patients into subgroups with different clinical outcomes for each type of cancer. We first collected a dataset consisting of gene expression profiles from 3,656 patients with their survival information from TCGA cohorts. We then used member genes in the top 3 ranked cancer-type-specific altered pathways to perform consensus clustering for each cancer type (see Material and Method section). We generated Kaplan-Meier (KM) curves based on the groups produced by consensus clustering and found that patient survival was significantly different among the groups (Figure 3 and Supplementary Figure 2). In TCGA KIRC, we found three patient subgroups (A, B and C), with the Group C having the poorest survival. A log-rank test indicated that Groups A and C had significantly different survival outcomes (Log-rank test p-value =  $1.840 \times 10^{-8}$ , Hazard ratio = 2.94) with median survival times of 41.9 months for group A compared to 30.8 months for group C (Figure 3A). Other examples are Bladder Urothelial Carcinoma (BLCA), Head and Neck squamous cell carcinoma (HNSC), and Skin Cutaneous Melanoma (SKCM) patient subgroups identified by NTriPath pathway signatures. While the molecular classification of clinically relevant subtypes of these cancers is still challenging, we found patient subgroups having significantly different survival in these cancers (Log-rank test p-value = 0.0086, 0.0010, and 0.0210, respectively) (Figure 3B, 3C and 3D).

Experiments with other TCGA datasets, including those for Breast invasive carcinoma (BRCA), Glioblastoma Multiforme (GBM), Lung adenocarcinoma (LUAD), and Ovarian serous cystadenocarcinoma (OV) consistently showed that the use of member genes in cancer-type-specific altered pathways could serve as a prognostic biomarker for patient stratification (Supplementary Figure 2 and 3). For comparison, we also attempted to cluster patients using significant frequently mutated

genes previously identified by the TCGA Pan-Cancer study<sup>1</sup>. The results of consensus clustering using the NTriPath-derived pathway signatures and the TCGA Pan-Cancer-derived mutated gene signatures showed that the results from NTriPath-derived pathway signatures had higher significance levels for BLAC, BRCA, and KIRC, and comparable results for the GBM, HNSC, and LUAD cancer types (Figure 4). These findings suggested that NTriPath-derived altered pathways could be used as prognostic biomarkers for better patient stratification.

### **Independent cohorts for the validation of the cancer-type-specific altered pathways**

We performed multiple validations to evaluate the robustness and the reproducibility of NTriPath. First, we evaluated the robustness of the cancer-type-specific altered pathways identified in the TCGA cohort for prognostic stratification. We generated gene expression profiles of 102 HNSC patients from our institution and used the member genes of the top 3 HNSC cancer-type-specific altered pathways in the TCGA cohort for patient stratification. In addition, we also used publically available gene expression data from two ovarian cancer datasets, one lung cancer dataset, two colon cancer datasets for a total of 1,112 patients, and used the top 3 cancer-type specific altered pathways for corresponding cancer type for independent validation. In the HNSC cohorts, we found six patient subgroups (A through F), with the group F patients having the poorest survival times (Figure 5A). A log-rank test indicated that groups A and F had significantly different survival outcomes (p-value = 0.038, hazard ratio = 1.88) with median survival times of 78.1 months for group A and 26.7 months for group F. Similarly, we found patient subgroups having significantly different survival outcomes in lung cancer, ovarian cancer, and colorectal cancer (Figure 5B-D and Supplementary Figure 3). Secondly, we verified the reproducibility of NTriPath for the identification of the cancer-type-specific altered pathways. We collected the

level 2 somatic mutation data from 19 human cancer types those were updated after we collected initial dataset used in the original experiments from the TCGA data portal. We found that there are 1891 newly updated patients' mutation data from 15 cancer types (see Supplementary Table 5). We re-ran NTriPath to identify cancer-type-specific pathways across 19 cancers using 6681 patients' somatic mutation data including those of newly updated patients' mutation data. Interestingly, we found that many top ranked pathways identified by NTriPath in the original experiments were consistently highly ranked in the new experiments (see Supplementary Table 6).

These results reassure that NTriPath is a robust tool to detect the altered pathways across cancers, and the altered pathways identified by NTriPath can serve as robust prognostic signatures for identifying patient subgroups with different survival outcomes across multiple cancer types.

#### **NTriPath identified potential therapeutic targets in poor prognosis patient subgroups**

We further investigated whether we could identify potential targets or the therapy for the identified poor prognosis patient subgroups. Interestingly, we found that many known drug targets in the cancer-type specific altered pathways are often up-regulated in poor prognosis patient subgroups across cancers (see Method section and Supplementary 7). For example, in TCGA KIRC cohort, *LRPI* and *MMP9*, targets of FDA-approved drugs Tenecteplase and Captopril, were significantly up-regulated in poor prognosis group compared to good prognosis group (FDR-adjusted p-value < 0.05 with t-test). Tenecteplase binds to *LRPI* and induces both *LRPI* and *MMP9* expression, and Captopril inhibits *MMP9* expression. Thus, combinatorial therapy of these drugs can be beneficial for the KIRC patients with high *LRPI* and *MMP9* expression<sup>27-44</sup>. Another notable example includes DNA Topoisomerase I (*TOP1*), a target of well-known FDA-approved anticancer drugs such as Irinotecan and

Topotecan, identified by NTriPath as a new member gene into the cancer-type-specific altered pathways across many cancers including HNSC. Interestingly, we found that *TOP1* was up-regulated in poor prognosis subgroups in HNSC from both TCGA and UTSW cohorts (Group E and F in Figure 3C and 5A, respectively). In addition, we found that some patients with overexpression of *TOP1* in TCGA HNSC poor prognosis subgroup have developed therapy resistance against single chemotherapeutic agent such as Cisplatin. Interestingly, there is an ongoing trial in advanced HNSC showing efficacy of *TOP1* inhibitor Irinotecan with Cisplatin in a poor prognosis patient subgroup<sup>45</sup>. These observations may suggest that *TOP1* inhibitors-based combinations might offer an effective treatment option for HNSC patients with overexpression of *TOP1*. Taken together, these findings suggested that the use of NTriPath-derived altered pathways containing available drug targets may allow for the development of more tailored therapeutics.

## Discussion

Systematic understanding of how somatic mutations influence clinical outcomes is essential for the development and application of personalized therapies. Especially organizing alterations at the individual gene level and in the molecular pathways can correlate altered pathways and vulnerabilities with specific genetic lesions, and provide novel insights into cancer biology, biomarkers for patient stratification in clinical trials, and potential targeted drug development<sup>46</sup>. Here, we systematically identified biological and clinical relevant cancer-type-specific altered cross multiple cancer types. In particular, the integration of somatic mutation with biological prior knowledge led to the identification of altered pathways that contain recurrently mutated genes as a hallmark of specific cancer types. Interestingly, we found that several genes, while not frequently mutated or not mutated at all in patients, were part

of cancer-type-specific altered pathways that have been causally implicated in the development of corresponding cancer types, and expressions of those genes are significantly associated with clinical outcomes (Supplementary Figure 4). For example, no mutation of *MMP7* has been reported, but high expression of *MMP7* ( $p = 0.00191$ , HR=1.7 (95%CI,1.21–2.38)) is significantly associated with poor survival in TCGA KIRC patients. Other examples include *CABLES1* ( $p = 0.00272$ , HR=0.486 (95%CI,0.301–0.787)) in TCGA HNSC and LUAD, and *GCHI* ( $p = 0.0000528$ , HR=0.52 (95%CI,0.367–0.763)) in TCGA SKCM are not frequently or not mutated but low or high expression of those genes are significantly associated with poor survival. In addition, we found that known drug targets are not frequently mutated but often up-regulated in poor prognosis patient subgroups across many cancers. These results further corroborate that the integrative analysis of somatic mutations with additional biological prior knowledge may elucidate potential candidate genes associated with clinical outcomes and could be potentially used to design targeted therapy, which cannot be readily identified by somatic mutation analysis alone.

In our analysis, we did not remove synonymous mutations or further select a shorter list of recurrent mutated genes in cohorts with stringent criteria either <sup>3, 47</sup>. However, Hopfree et al<sup>7</sup> showed that filtering synonymous mutations resulted in a decreased ability to detect patient subgroups with different survival outcomes. Another recent study also showed that synonymous mutations could affect functions of oncogene and tumor suppressors <sup>48</sup>. In addition, our experimental results for patient stratification in comparison with recurrent mutated gene signatures identified by the TCGA Pan-Cancer<sup>1</sup> indicated that the use of NTriPath-derived pathways showed a comparable or better performance in discovering patient subgroups with different survival outcomes across cancers. To evaluate the impact of different network resources, we used

networks from the HPRD <sup>49</sup> and Rossin, E.J. *et al*<sup>50</sup> and repeated experiments. We summarize the results of altered pathways and patient stratification using different network resources and provide on our supplement website.

Lastly, NTriPath is a general computational algorithm and can be applied to other data types such as gene expression, copy number alteration, and methylation to identify altered pathways by different types of genomic aberrations. NTriPath can also be used to find altered pathways across associated with other cancer-related phenotypes (e.g., patient groups having therapy resistance vs. sensitivity, metastatic vs. non-metastatic).

## Conclusions

We have described an integrative somatic mutation analysis for discovering altered pathways in human cancers. NTriPath integrates somatic mutation data and prior biological knowledge from the pathway database and molecular networks to identify significantly altered pathways and their associations with specific cancer types. Specifically, NTriPath effectively utilizes mutation patterns that exist in only a subset of samples (or specific cancer types), thus revealing pathways altered by complex mutation patterns across cancer types. Furthermore, the use of gene-gene interaction networks and the pathway database provides the potential to identify altered pathways enriched with genes harboring mutations at high/intermediate frequencies, as well as those not mutated per se but nevertheless playing critical roles in tumorigenesis in network and pathway contexts. Thus, NTriPath is uniquely suited to provide a global analysis of altered pathways by somatic mutation across cancer types.

We applied NTriPath to somatic mutation data from 19 types of cancers, and discovered cancer-type-specific altered pathways based on these mutations in human cancers. Functional enrichment analysis of cancer-type-specific pathways demonstrated that the identified cancer-type-specific altered pathways are biologically

meaningful to each cancer type. It also provided unique pathway views of key biological processes underlying each cancer type. Of particular significance, we identified a patient subgroup with poor survival by cancer-type-specific altered pathway signatures from TCGA cohorts, which in independent cohorts. These results implied the potential utility of cancer-type-specific altered pathway signatures to serve as a guide to tailored treatment in a patient subgroup.

## Materials and Methods

### Somatic mutation, human gene-gene interaction networks, and pathway data

The level 2 somatic mutation data from 19 human cancer types were collected from the TCGA data portal on May 19<sup>th</sup> 2013<sup>51</sup>. We constructed a gene-gene interaction network by combining networks from Zhang, S. *et al*<sup>52</sup>, the Human Protein Reference Database (Dec. 2013)<sup>53</sup> and Rossin, E.J. *et al*<sup>50</sup>. Four sets of pathways were used in the analysis: 1) 4,620 conserved subnetworks from the human gene-gene interaction network<sup>12</sup>, 2) KEGG, 3) Biocarta, and 4) Reactome gene sets from MsigDB (Sept. 2010)<sup>54</sup>.

### Algorithm

The algorithm identifies pathways disrupted by mutated genes. Disrupted pathways are found based on the factorization results from the network regularized nonnegative tri-matrix factorization.

#### 1. Notations

We construct a binary data matrix  $\mathbf{X} \in R^{N \times M}$  from the mutation data, where  $N$  is the number of patients,  $M$  is the number of genes and the  $(i, j)$ <sup>th</sup> element of the matrix  $\mathbf{X}$ ,  $[\mathbf{X}]_{ij}$ , is '1' if the  $i$ th patient has a mutation on the  $j$ th gene, '0' otherwise. We derive the adjacency matrix from the human gene-gene interaction networks and denote it as

$\mathbf{A}$ , where  $[\mathbf{A}]_{ij}='1'$  if the  $i$ th gene is interacting with the  $j$ th genes in the networks and '0' otherwise. We define the graph Laplacian matrix by  $\mathbf{L} = \mathbf{D} - \mathbf{A}$ , where each diagonal element in the diagonal matrix  $\mathbf{D}$  is given by  $[\mathbf{D}]_{ii} = \sum_j [\mathbf{A}]_{ij}$ . We construct a binary matrix  $\mathbf{U} \in R^{N \times K_1}$  denoting patient cluster, where  $K_1$  indicates the number of cancer types and  $[\mathbf{U}]_{ij}=1$  indicates the  $i$ th patient has  $j$ th cancer type. We construct a binary matrix  $\mathbf{V}_0 \in R^{M \times K_2}$  from the specific pathway database denoting pathway information, where  $K_2$  is the number of pathways and  $[\mathbf{V}_0]_{ij}=1$  if the  $i$ th gene is annotated in  $j$ th pathway as a member in the pathway database, otherwise 0. Since current pathway database annotation is still incomplete, we define a matrix  $\mathbf{V} \in R^{M \times K_2}$  denoting newly updated pathway information including newly added member genes by NTriPath. We define a matrix  $\mathbf{S} \in R^{K_1 \times K_2}$  denoting cancer type and pathway associations, where each element of  $[\mathbf{S}]_{ij}$  represents associations between  $i$ th cancer type with  $j$ th pathway. Higher values of elements indicate stronger associations between cancer types and pathways. Since  $\mathbf{V}$  and  $\mathbf{S}$  are unknown, we need to learn about those matrices during optimization (see below section for details)

## 2. Network regularized non-negative tri-matrix factorization

The network regularized nonnegative tri-matrix factorization is an extension of Nonnegative Tri Matrix Factorization (NTMF); in this work, somatic mutation data  $\mathbf{X}$  is factorized as the products of three element-wise non-negative matrices  $\mathbf{U}$ ,  $\mathbf{S}$ , and  $\mathbf{V}$  denoting patient's cancer type, cancer-type and pathway associations, and cancer-related pathways, respectively. We here consider a weighted loss function to deal with the sparseness of the data matrix  $\mathbf{X}$  (More than 98% of entries are zero). It enables us to focus on an approximation error at nonzero entries, which correspond to mutated genes. In addition, to incorporate the prior knowledge from human gene-gene interaction networks and pathway datasets into factorizations, we enforce constraints



on parameters, which involve the graph Laplaican  $\mathbf{L}$  and the pathway information  $\mathbf{V}_0$ .

All these ideas are accomplished by minimizing the following objective function

$$\min_{\mathbf{S}, \mathbf{V} \geq 0} \|\mathbf{W} \circ (\mathbf{X} - \mathbf{USV}^T)\|_F^2 + \lambda_S \|\mathbf{S}\|_1^2 + \lambda_V \|\mathbf{V}\|_1^2 + \lambda_0 \|\mathbf{V} - \mathbf{V}_0\|_F^2 + \lambda_L \text{tr}\{\mathbf{V}^T \mathbf{L} \mathbf{V}\},$$

where  $\mathbf{W} \in R^{N \times M}$  is a weight matrix where  $[\mathbf{W}]_{ij} = 1$  if  $[\mathbf{X}]_{ij} > 0$  otherwise 0, and the operator  $\circ$  represents the element-wise multiplication. Here, we are only interested in learning of  $\mathbf{S}$  and  $\mathbf{V}$  among three factor matrices, since factor  $\mathbf{U}$  can be obtained from the patient's clinical information.

To solve our minimization problem, we adapt the multiplicative update method for NTMF proposed in a recent study<sup>55</sup>, which contains a routine for avoiding 'inadmissible zeros problem' where the solution of multiplicative update rules is stuck at zero when an entry in the factor becomes.

### Step 1: Initialization

Initialize the factor matrices  $\mathbf{S} = \mathbf{1}$  and  $\mathbf{V} = \mathbf{V}_0$ , where  $\mathbf{1} \in R^{K_1 \times K_2}$  is a matrix whose elements are all one. Set the regularization parameters  $\lambda_S = \lambda_V = \lambda_L = 1$  and  $\lambda_0 = 0.1$ . Set the user specified parameters for avoiding the inadmissible zeros problem,  $\kappa_{tol} = 10^{-10}$ ,  $\kappa = 10^{-6}$  and  $\epsilon = 10^{-10}$ .

### Step 2: Iteration

Iterate until it converges or reaches the maximum number of iterations:

$$[\mathbf{S}]_{ij} \leftarrow ([\mathbf{S}]_{ij} + \kappa_{ij}^S) \tau_{ij}^S$$

$$[\mathbf{V}]_{ij} \leftarrow ([\mathbf{V}]_{ij} + \kappa_{ij}^V) \tau_{ij}^V$$

where  $\kappa_{ij}^M$  is set to  $\kappa$  if  $[\mathbf{M}]_{ij} \geq \kappa_{tol}$  and  $\tau_{ij}^M > 1$ , otherwise 0, and

$$\tau_{ij}^S = \frac{[\mathbf{U}^T \mathbf{X} \mathbf{V}]_{ij}}{[\mathbf{U}^T (\mathbf{W} \circ (\mathbf{USV}^T)) \mathbf{V}]_{ij} + \lambda_S \|\mathbf{S}\|_1 + \epsilon},$$

$$\tau_{ij}^V = \frac{[\mathbf{X}^T \mathbf{U} \mathbf{S} + \lambda_L \mathbf{A} \mathbf{V} + \lambda_0 \mathbf{V}_0]_{ij}}{[(\mathbf{W} \circ (\mathbf{USV}^T))^T \mathbf{U} \mathbf{S} + \lambda_0 \mathbf{V} + \lambda_L \mathbf{D} \mathbf{V}]_{ij} + \lambda_V \|\mathbf{V}\|_1 + \epsilon}.$$

Empirically, the algorithm converges fast within 50 iterations in the experiments.

### **3. Identification of cancer-type-specific altered pathways**

Once the above optimization problem is solved, we used  $\mathbf{S}$  matrix to identify cancer-type-specific altered pathways across cancers. Specifically, we ranked pathways based on values of elements of the  $\mathbf{S}$  matrix for each cancer type (e.g., rank pathways based on all values of  $i$ th row which indicate association scores between  $i$ th cancer type and all pathways). In addition, to measure statistical significance of cancer-type and pathway associations, we performed a permutation test (e.g., we randomly permuted somatic mutation data and repeated experiments 5,000 times to calculate empirical p-values) and defined cancer-type-specific altered pathways based on the following strict criteria: 1) Pathways must be ranked within the top  $K$ th compared to other pathways in each cancer type based on their association scores in matrix  $\mathbf{S}$ . 2) Pathways must have significant BH-adjusted p-values (Benjamini-Hochberg adjusted p-values using a false discovery rate cutoff of 0.1) (See Supplementary X). In this work, we selected the top 3 ranked pathways having significant BH-adjusted p-values per each cancer type. Top ranked pathways for KICH, KIRP, and THCA were excluded for further analysis, due to the insignificant BH-adjusted p-values.

#### **Gene expression data and clustering**

We collected RNA-seq data for TCGA BLCA, BRCA, HNSC, KIRC, LAML, LUAD, SKCM, STAD, UCEC from cBioPortal<sup>56-58</sup> using CGDS MATLAB toolbox with RNA Seq V2 RSEM option. We collected microarray gene expression profiles for TCGA GBM from the TCGA dataportal<sup>51</sup> and TCGA OV and two others from Zhang, W. *et al.*<sup>59</sup>. We collected colon cancer data from GSE39582 and lung cancer data from Shedden, K. *et al.*<sup>60</sup>. RNA-seq data were z-transformed while other expression data were quantile normalized, log transformed, and expression values

were median centered. To perform consensus clustering, we used Matlab K-means clustering and used two-way hierarchical clustering.

## **Authors' contributions**

SP and THH designed the study. SP and THH performed the analysis. THH supervised the project. SP, JN, SL, and THH wrote the paper. All authors read and approved the final manuscript.

## **Acknowledgements**

We would like to thank to Quantitative Biomedical Research Center to allow us to use their computational resources.

## Figures

### Figure 1. Overview of NTriPath

This figure describes steps to discover altered pathways across multiple cancer types. The aim of the approach is to integrate the somatic mutation data with gene-gene interaction networks and a pathway database for discovering altered pathways across cancers.

### Figure 2. KIRC-specific altered pathways

(A) Diagrams of the top three ranked altered pathways in patients with KIRC. Red color indicates genes that are frequently mutated. A circular-shaped node represents the original member genes annotated in the pathway database, and a diamond-shaped node represents newly identified members genes of the pathways by NTriPath (B) Protein and mRNA expression and mutation status for all genes identified in the top three KIRC altered pathways. Each row represents a member gene in the TCGA

KIRC-specific altered pathway, and each column represents a patient sample in TCGA KIRC cohort.

Figure 3. Cancer-type-specific altered pathways across cancers correlate with survival outcomes

Kaplan-Meier survival plots based on patient subgroups defined by consensus clustering using genes from the top 3 altered pathways for (a) kidney renal cell carcinoma (KIRC), (b) bladder urothelial carcinoma (BLCA), (c) head and neck squamous carcinoma (HNSC), and (d) skin cutaneous melanoma (SKCM).

Figure 4. Comparing NTriPath-derived signatures with mutation-frequency-based signatures

This figure describes comparisons of patient stratification using signatures derived from NTriPath and mutation frequency reported in Kandoth, C. *et al.* <sup>1</sup>.

Figure 5. Validation in independent cohort

This figure describes Kaplan-Meier survival plots for patient subgroups from (a) UTSW HNSC (b) Lung adenocarcinoma (c) Colon cancer, (d) Ovarian cancer.

	Cancer Type	Number of Patients
1	Acute Myeloid Leukemia (LAML)	75
2	Bladder Urothelial Carcinoma (BLCA)	136
3	Brain Lower Grade Glioma (LGG),	217
4	Breast Invasive Carcinoma (BRCA)	772
5	Cervical Squamous Cell Carcinoma and Endocervical Adenocarcinoma (CESC)	41
6	Colon Adenocarcinoma (COAD)	270
7	Glioblastoma Multiforme (GBM)	290
8	Head and Neck Squamous Cell Carcinoma (HNSC)	323

9	Kidney Chromophobe Renal Cell Carcinoma (KICH)	65
10	Kidney Renal Clear Cell Carcinoma (KIRC)	210
11	Kidney Renal Papillary Cell Carcinoma (KIRP)	111
12	Lung Adenocarcinoma (LUAD)	380
13	Ovarian Serous Cystadenocarcinoma (OV)	463
14	Prostate Adenocarcinoma (PRAD)	171
15	Rectum Adenocarcinoma (READ)	116
16	Skin Cutaneous Melanoma (SKCM)	269
17	Stomach Adenocarcinoma (STAD)	264
18	Thyroid Carcinoma (THCA)	369
19	Uterine Corpus Endometrioid Carcinoma	248

Table 1. A list of cancer types used in the analysis.

## References

1. Kandoth, C. *et al.* Mutational landscape and significance across 12 major cancer types. *Nature* 502, 333-339 (2013).
2. Tamborero, D. *et al.* Comprehensive identification of mutational cancer driver genes across 12 tumor types. *Sci. Rep.* 3 (2013).
3. Lawrence, M.S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature* 499, 214-218 (2013).
4. Lawrence, M.S. *et al.* Discovery and saturation analysis of cancer genes across 21 tumour types. *Nature* 505, 495-501 (2014).
5. Osmanbeyoglu, H.U., Pelossof, R., Bromberg, J.F. & Leslie, C.S. Linking signaling pathways to transcriptional programs in breast cancer. *Genome Res* 24, 1869-1880 (2014).
6. Baselga, J. Targeting the phosphoinositide-3 (PI3) kinase pathway in breast cancer. *The oncologist* 16 Suppl 1, 12-19 (2011).
7. Hofree, M., Shen, J.P., Carter, H., Gross, A. & Ideker, T. Network-based stratification of tumor mutations. *Nat Meth* 10, 1108-1115 (2013).
8. Hwang, T. *et al.* Large-scale integrative network-based analysis identifies common pathways disrupted by copy number alterations across cancers. *BMC Genomics* 14, 1-13 (2013).
9. Vaske, C.J. *et al.* Inference of patient-specific pathway activities from multi-dimensional cancer genomics data using PARADIGM. *Bioinformatics* 26, i237-i245 (2010).
10. Cerami, E., Demir, E., Schultz, N., Taylor, B.S. & Sander, C. Automated Network Analysis Identifies Core Pathways in Glioblastoma. *PLoS ONE* 5, e8918 (2010).

11. Vandin, F., Upfal, E. & Raphael, B. Algorithms for detecting significantly mutated pathways in cancer. *Journal of computational biology : a journal of computational molecular cell biology* 18, 507-522 (2011).
12. Suthram, S. *et al.* Network-Based Elucidation of Human Disease Similarities Reveals Common Functional Modules Enriched for Pluripotent Drug Targets. *PLoS Comput Biol* 6, e1000662 (2010).
13. Hofree, M., Shen, J.P., Carter, H., Gross, A. & Ideker, T. Network-based stratification of tumor mutations. *Nature methods* 10, 1108-1115 (2013).
14. Forbes, S.A. *et al.* COSMIC: exploring the world's knowledge of somatic mutations in human cancer. *Nucleic acids research* 43, D805-811 (2015).
15. Pena-Llopis, S. *et al.* BAP1 loss defines a new class of renal cell carcinoma. *Nat Genet* 44, 751-759 (2012).
16. Sato, Y. *et al.* Integrated molecular analysis of clear-cell renal cell carcinoma. *Nat Genet* 45, 860-867 (2013).
17. Peña-Llopis, S. & Brugarolas, J. Simultaneous isolation of high-quality DNA, RNA, miRNA and proteins from tissues for genomic applications. *Nat. Protocols* 8, 2240-2255 (2013).
18. Nickerson, M.L. *et al.* Improved Identification of von Hippel-Lindau Gene Alterations in Clear Cell Renal Tumors. *Clinical Cancer Research* 14, 4726-4734 (2008).
19. The Cancer Genome Atlas Research, N. Comprehensive molecular characterization of clear cell renal cell carcinoma. *Nature* 499, 43-49 (2013).
20. Gialeli, C., Theocharis, A.D. & Karamanos, N.K. Roles of matrix metalloproteinases in cancer progression and their pharmacological targeting. *The FEBS journal* 278, 16-27 (2011).
21. Staudt, N.D. *et al.* Myeloid Cell Receptor LRP1/CD91 Regulates Monocyte Recruitment and Angiogenesis in Tumors. *Cancer Research* 73, 3902-3912 (2013).
22. Langlois, B. *et al.* LRP-1 Promotes Cancer Cell Invasion by Supporting ERK and Inhibiting JNK Signaling Pathways. *PLoS ONE* 5, e11584 (2010).
23. Song, H., Li, Y., Lee, J., Schwartz, A.L. & Bu, G. Low-density lipoprotein receptor-related protein 1 promotes cancer cell migration and invasion by inducing the expression of matrix metalloproteinases 2 and 9. *Cancer Res* 69, 879-886 (2009).
24. Duan, D. *et al.* Inhibition of transcription elongation by the VHL tumor suppressor protein. *Science* 269, 1402-1406 (1995).
25. Ohh, M. *et al.* Ubiquitination of hypoxia-inducible factor requires direct binding to the [bgr]-domain of the von Hippel-Lindau protein. *Nat Cell Biol* 2, 423-427 (2000).
26. Tanimoto, K., Makino, Y., Pereira, T. & Poellinger, L. Mechanism of regulation of the hypoxia-inducible factor-1 $\alpha$  by the von Hippel-Lindau tumor suppressor protein. *The EMBO Journal* 19, 4298-4309 (2000).
27. Yamamoto, D., Takai, S. & Miyazaki, M. Inhibitory profiles of captopril on matrix metalloproteinase-9 activity. *European journal of pharmacology* 588, 277-279 (2008).
28. Williams, R.N., Parsons, S.L., Morris, T.M., Rowlands, B.J. & Watson, S.A. Inhibition of matrix metalloproteinase activity and growth of gastric adenocarcinoma cells by an angiotensin converting enzyme inhibitor in in vitro and murine models. *European journal of surgical oncology : the journal*

- of the European Society of Surgical Oncology and the British Association of Surgical Oncology 31, 1042-1050 (2005).
29. Underwood, C.K., Min, D., Lyons, J.G. & Hambley, T.W. The interaction of metal ions and Marimastat with matrix metalloproteinase 9. *Journal of inorganic biochemistry* 95, 165-170 (2003).
  30. Reinhardt, D. *et al.* Cardiac remodelling in end stage heart failure: upregulation of matrix metalloproteinase (MMP) irrespective of the underlying disease, and evidence for a direct inhibitory effect of ACE inhibitors on MMP. *Heart* 88, 525-530 (2002).
  31. Rask-Andersen, M., Almen, M.S. & Schioth, H.B. Trends in the exploitation of novel drug targets. *Nature reviews. Drug discovery* 10, 579-590 (2011).
  32. Rajapakse, N., Mendis, E., Kim, M.M. & Kim, S.K. Sulfated glucosamine inhibits MMP-2 and MMP-9 expressions in human fibrosarcoma cells. *Bioorganic & medicinal chemistry* 15, 4891-4896 (2007).
  33. Overington, J.P., Al-Lazikani, B. & Hopkins, A.L. How many drug targets are there? *Nature reviews. Drug discovery* 5, 993-996 (2006).
  34. Okada, M. *et al.* Captopril attenuates matrix metalloproteinase-2 and -9 in monocrotaline-induced right ventricular hypertrophy in rats. *Journal of pharmacological sciences* 108, 487-494 (2008).
  35. Nenan, S. *et al.* Metalloelastase (MMP-12) induced inflammatory response in mice airways: effects of dexamethasone, rolipram and marimastat. *European journal of pharmacology* 559, 75-81 (2007).
  36. Mendis, E., Kim, M.M., Rajapakse, N. & Kim, S.K. Carboxy derivatized glucosamine is a potent inhibitor of matrix metalloproteinase-9 in HT1080 cells. *Bioorganic & medicinal chemistry letters* 16, 3105-3110 (2006).
  37. Imming, P., Sinning, C. & Meyer, A. Drugs, their targets and the nature and number of drug targets. *Nature reviews. Drug discovery* 5, 821-834 (2006).
  38. Fenton, J.I., Chlebek-Brown, K.A., Caron, J.P. & Orth, M.W. Effect of glucosamine on interleukin-1-conditioned articular cartilage. *Equine veterinary journal. Supplement*, 219-223 (2002).
  39. Dodge, G.R. & Jimenez, S.A. Glucosamine sulfate modulates the levels of aggrecan and matrix metalloproteinase-3 synthesized by cultured human osteoarthritis articular chondrocytes. *Osteoarthritis and cartilage / OARS, Osteoarthritis Research Society* 11, 424-432 (2003).
  40. Chu, S.C. *et al.* Glucosamine sulfate suppresses the expressions of urokinase plasminogen activator and inhibitor and gelatinases during the early stage of osteoarthritis. *Clinica chimica acta; international journal of clinical chemistry* 372, 167-172 (2006).
  41. Chen, X., Ji, Z.L. & Chen, Y.Z. TTD: Therapeutic Target Database. *Nucleic acids research* 30, 412-415 (2002).
  42. Berman, H.M. *et al.* The Protein Data Bank. *Nucleic acids research* 28, 235-242 (2000).
  43. Sato, A. *et al.* Inhibition of MMP-9 using a pyrrole-imidazole polyamide reduces cell invasion in renal cell carcinoma. *International journal of oncology* 43, 1441-1446 (2013).
  44. Hu, K. *et al.* Tissue-type plasminogen activator acts as a cytokine that triggers intracellular signal transduction and induces matrix metalloproteinase-9 gene expression. *The Journal of biological chemistry* 281, 2120-2127 (2006).

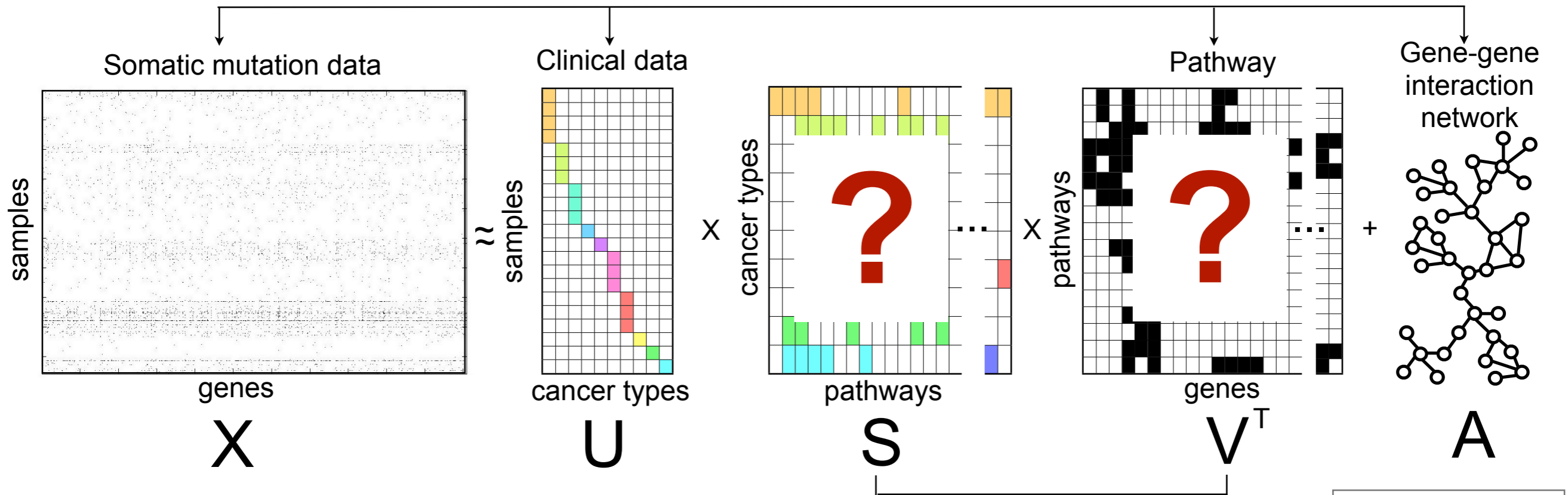
45. Gilbert, J. *et al.* Phase II trial of irinotecan plus cisplatin in patients with recurrent or metastatic squamous carcinoma of the head and neck. *Cancer* 113, 186-192 (2008).
46. Garraway, Levi A. & Lander, Eric S. Lessons from the Cancer Genome. *Cell* 153, 17-37 (2013).
47. Dees, N.D. *et al.* MuSiC: Identifying mutational significance in cancer genomes. *Genome Research* 22, 1589-1598 (2012).
48. Supek, F., Minana, B., Valcarcel, J., Gabaldon, T. & Lehner, B. Synonymous mutations frequently act as driver mutations in human cancers. *Cell* 156, 1324-1335 (2014).
49. Keshava Prasad, T.S. *et al.* Human Protein Reference Database--2009 update. *Nucleic acids research* 37, D767-772 (2009).
50. Rossin, E.J. *et al.* Proteins Encoded in Genomic Regions Associated with Immune-Mediated Disease Physically Interact and Suggest Underlying Biology. *PLoS Genet* 7, e1001273 (2011).
51. strel'tsov, S.A., Mikheikin, A.L. & Nechipurenko Iu, D. [Interaction of topotecan--a DNA topoisomerase I inhibitor--with dual-stranded polydeoxyribonucleotides. II. Formation of a complex containing several DNA molecules in the presence of topotecan]. *Molekuliarnaia biologii* 35, 442-450 (2001).
52. Zhang, S., Li, Q., Liu, J. & Zhou, X.J. A novel computational framework for simultaneous integration of multiple types of genomic data to identify microRNA-gene regulatory modules. *Bioinformatics* 27, i401-i409 (2011).
53. Keshava Prasad, T.S. *et al.* Human Protein Reference Database—2009 update. *Nucleic acids research* 37, D767-D772 (2009).
54. Subramanian, A. *et al.* Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences of the United States of America* 102, 15545-15550 (2005).
55. Seung-Jun, K., TaeHyun, H. & Giannakis, G.B. in Cognitive Information Processing (CIP), 2012 3rd International Workshop on 1-6 (2012).
56. Teicher, B.A. Next generation topoisomerase I inhibitors: Rationale and biomarker strategies. *Biochemical pharmacology* 75, 1262-1271 (2008).
57. Cerami, E. *et al.* The cBio Cancer Genomics Portal: An Open Platform for Exploring Multidimensional Cancer Genomics Data. *Cancer Discovery* 2, 401-404 (2012).
58. Gao, J. *et al.* Integrative Analysis of Complex Cancer Genomics and Clinical Profiles Using the cBioPortal. *Sci. Signal.* 6, p11- (2013).
59. Zhang, W. *et al.* Network-based Survival Analysis Reveals Subnetwork Signatures for Predicting Outcomes of Ovarian Cancer Treatment. *PLoS Comput Biol* 9, e1002975 (2013).
60. Shedden, K. *et al.* Gene expression-based survival prediction in lung adenocarcinoma: a multi-site, blinded validation study. *Nat Med* 14, 822-827 (2008).



Fig. 1

**Input:**

- 1) Somatic mutation
- 2) Clinical data (e.g., patient's cancer type)
- 3) Pathway information
- 4) Gene-gene interaction networks



$$\min_{S, V} \|X - USV^T\|_F^2 + \|S\|_1^2 + \|V - V_0\|_F^2 + \text{tr}(V^T(D-A)V)$$

**Output:**

- 1) Newly updated pathway information (V)
- 2) Cancer type-pathway association (S)

**Newly updated member genes in the pathways**

**X:** Somatic mutation data (#patients X #genes)  
**U:** Patient cluster (#patients X #cancer types)  
**S:** Cancer type-pathway association matrix (#cancer types X #pathways)  
 **$V_0$ :** Initial pathway information (#pathways X #genes)  
**V:** Newly updated pathway information (#pathways X #genes)  
**A:** Gene-gene interaction network (#genes X #genes)

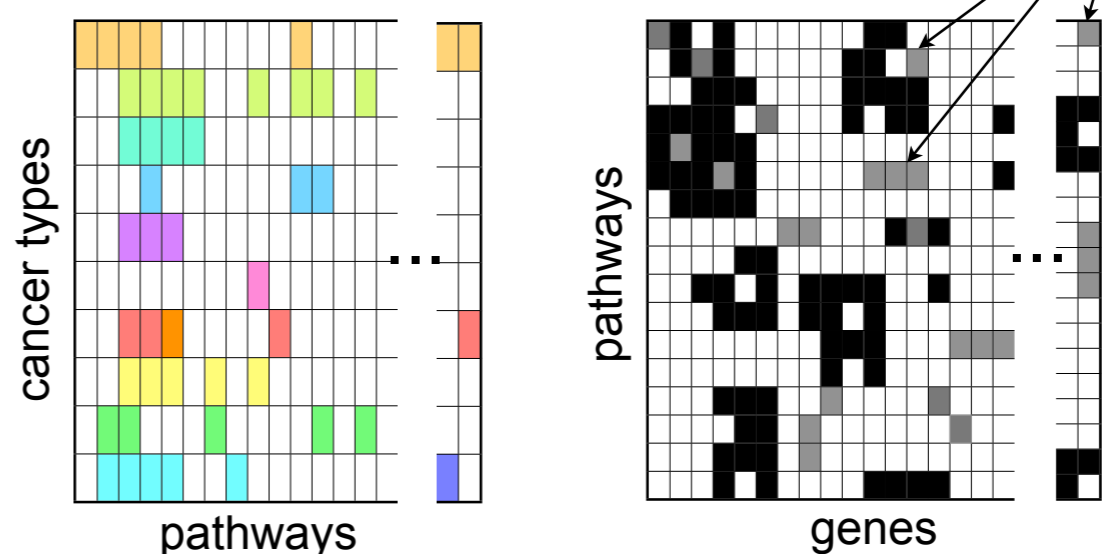


Fig. 2

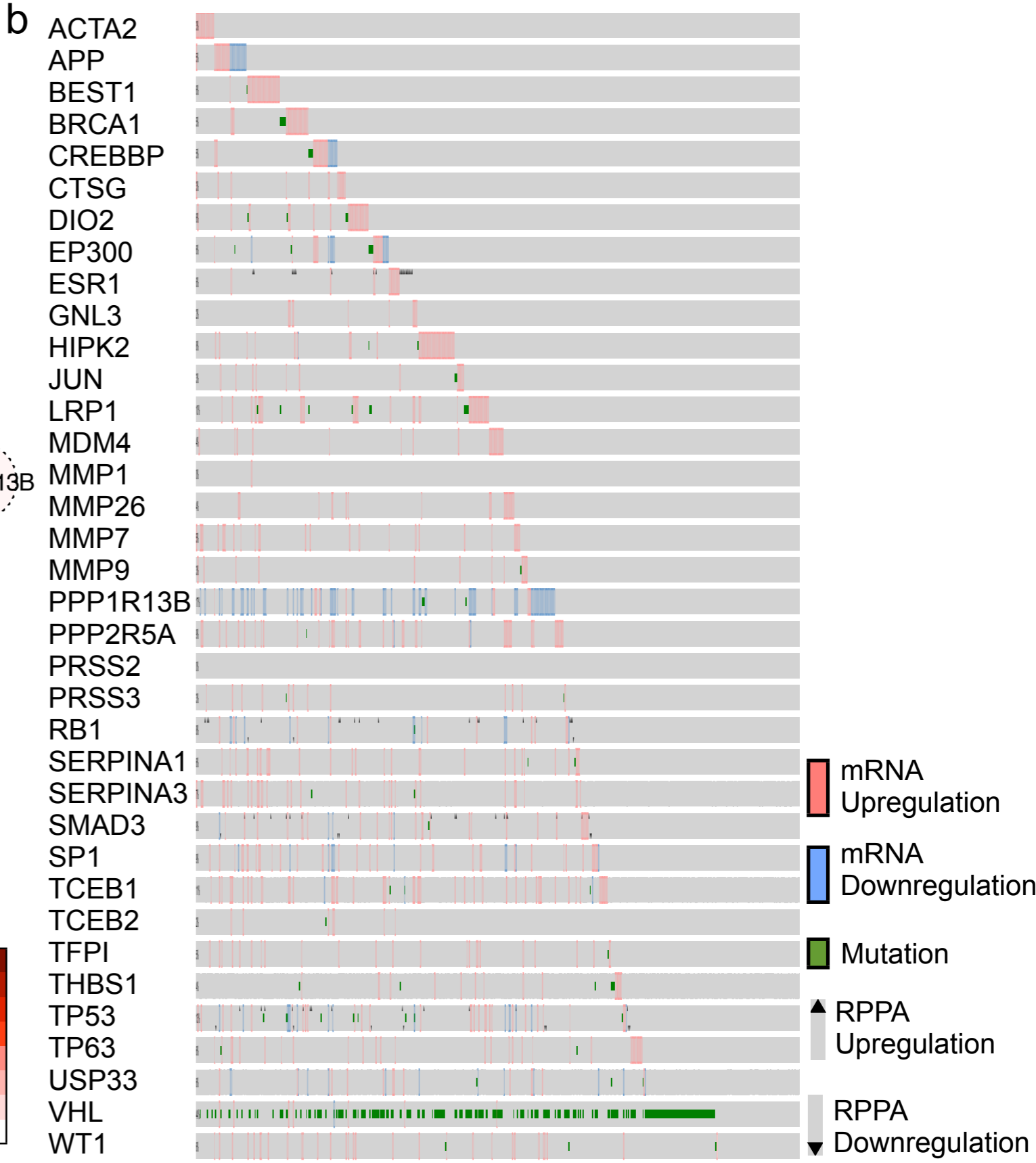
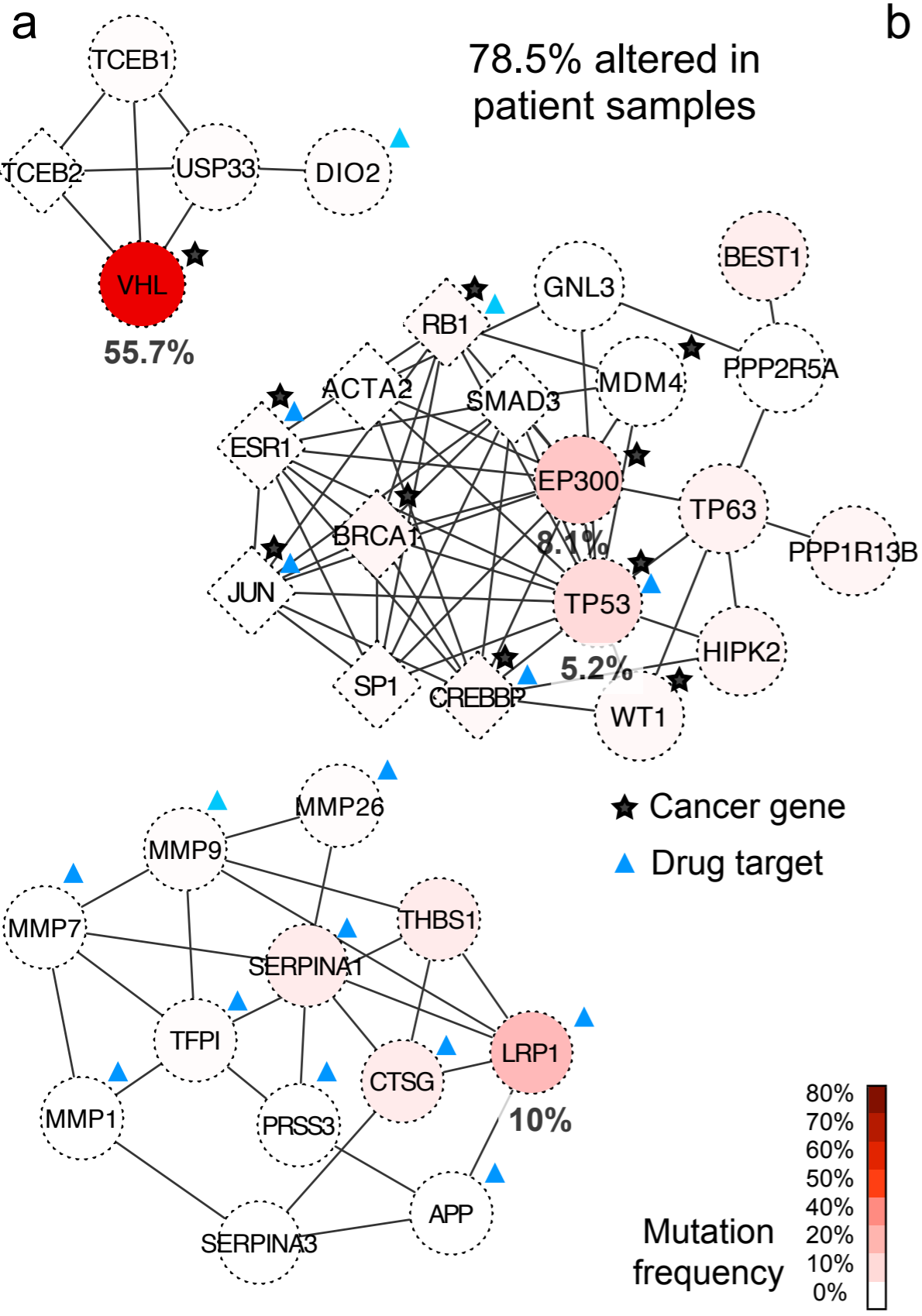


Fig. 3

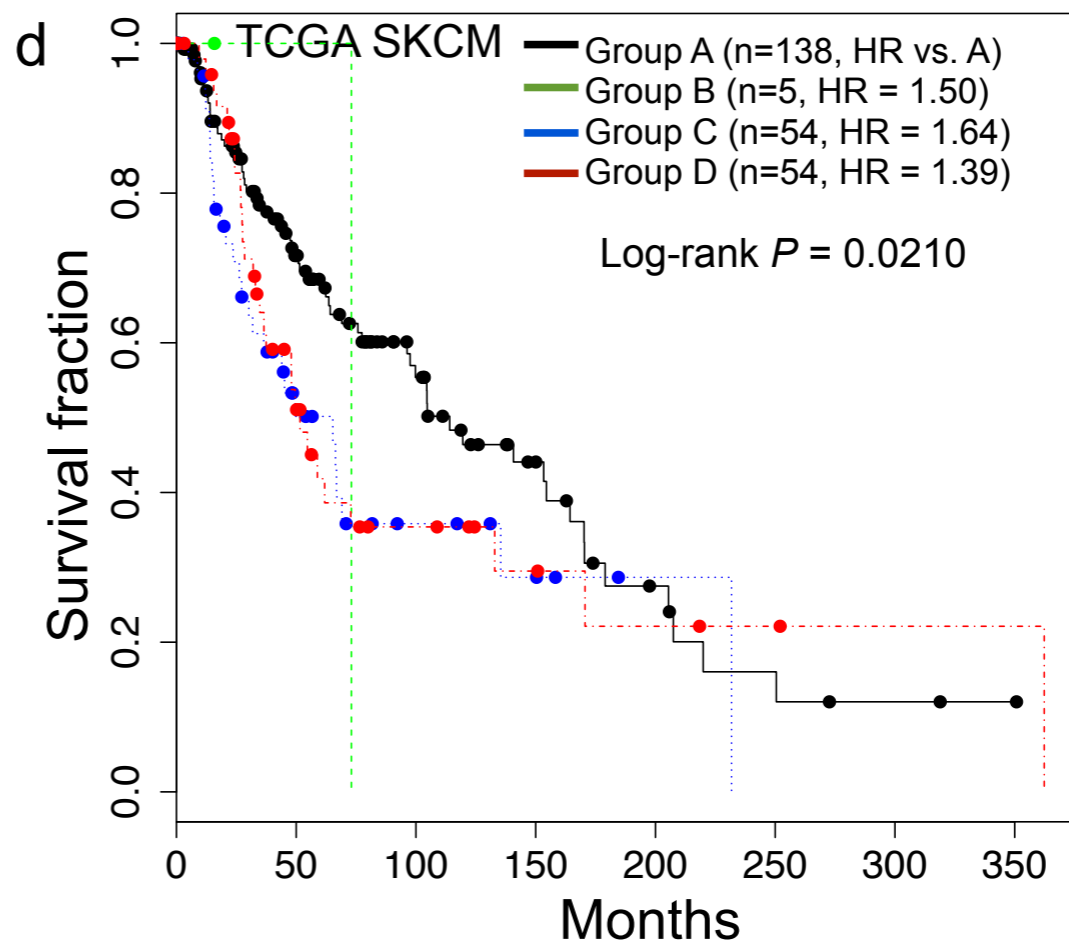
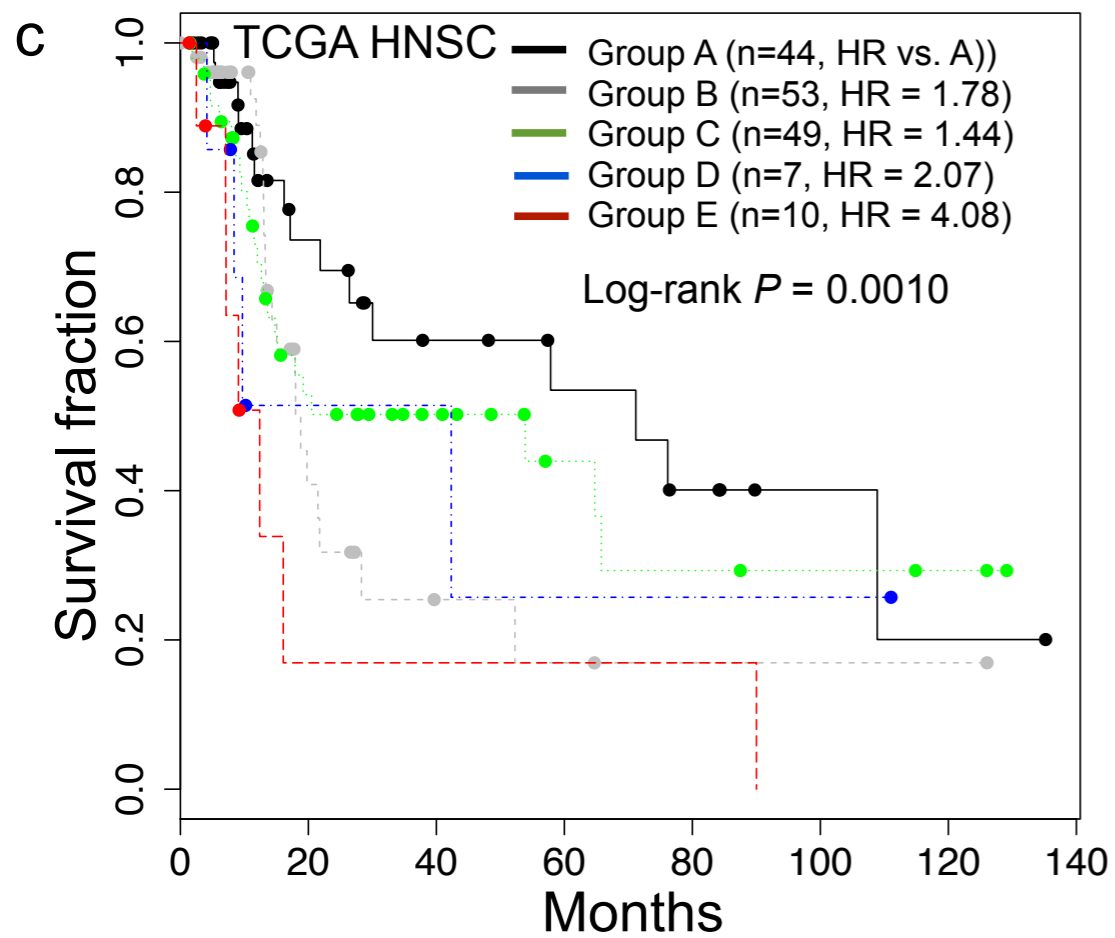
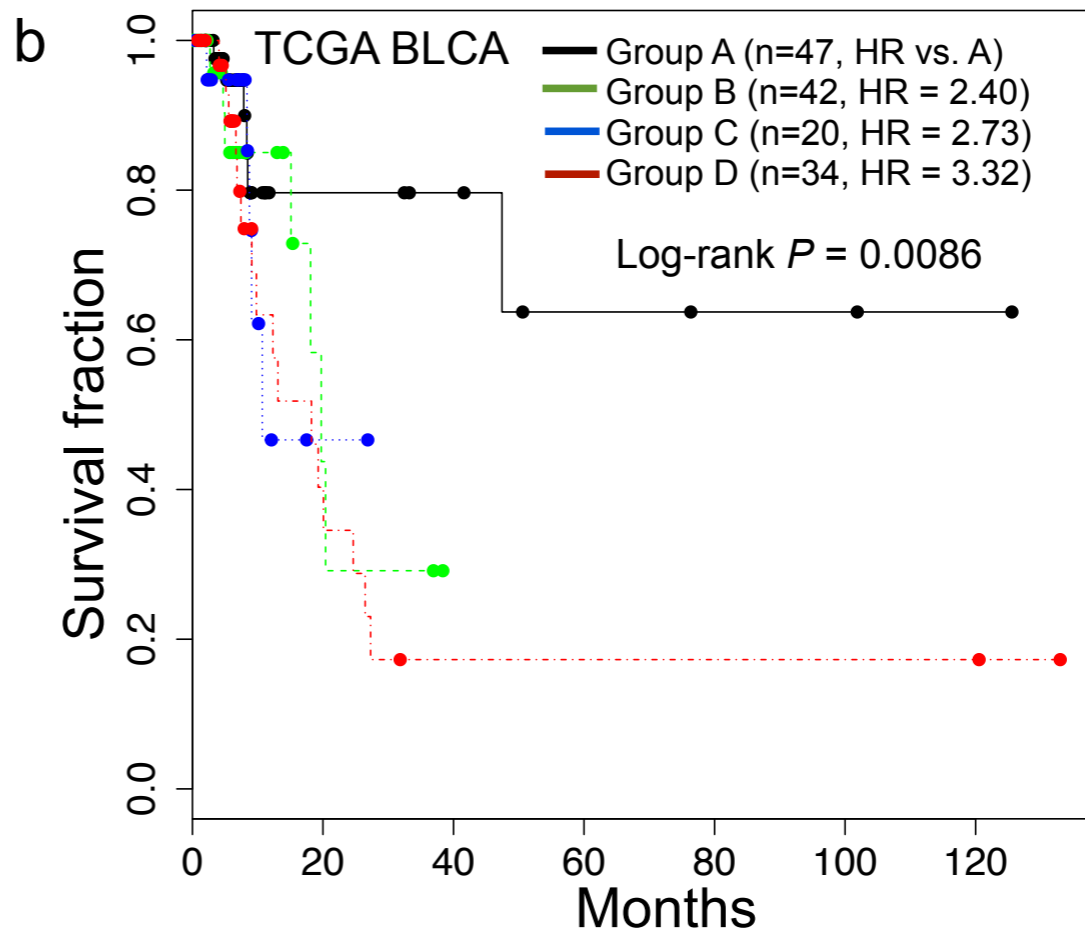
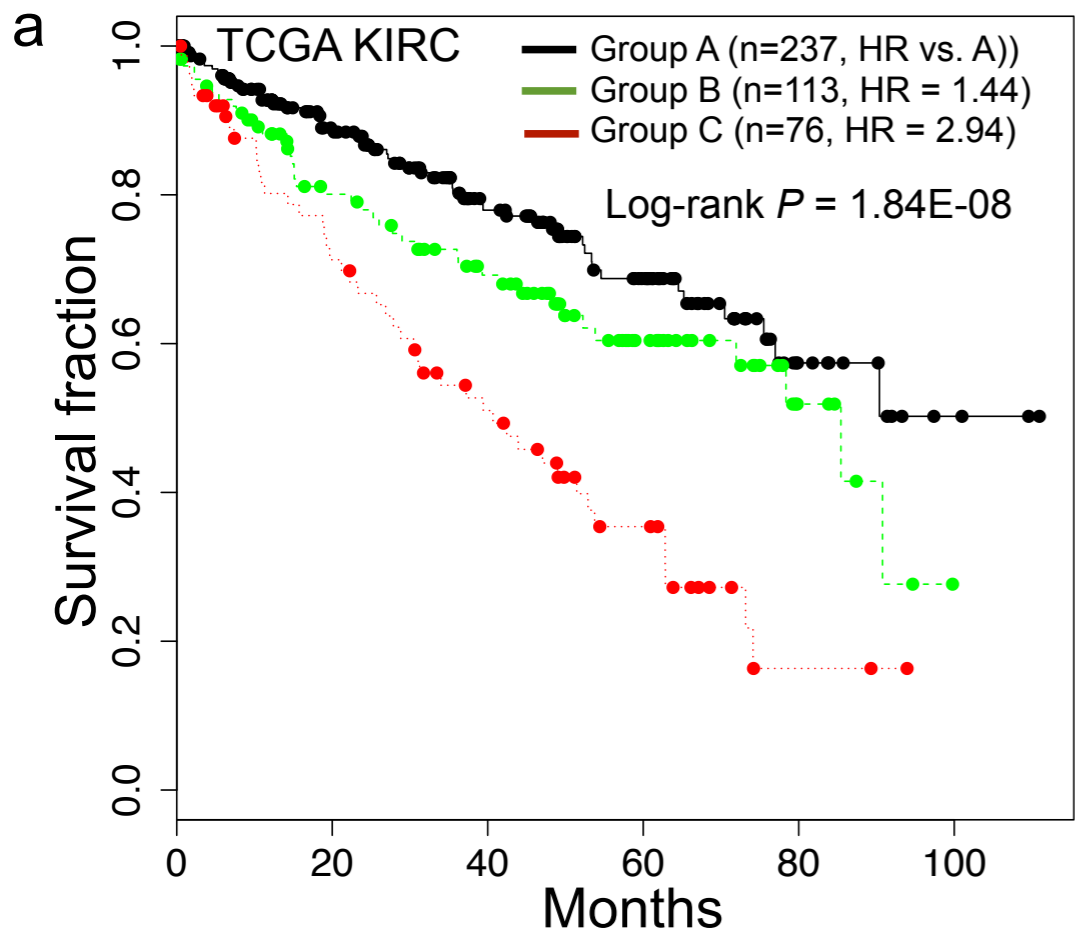


Fig. 4

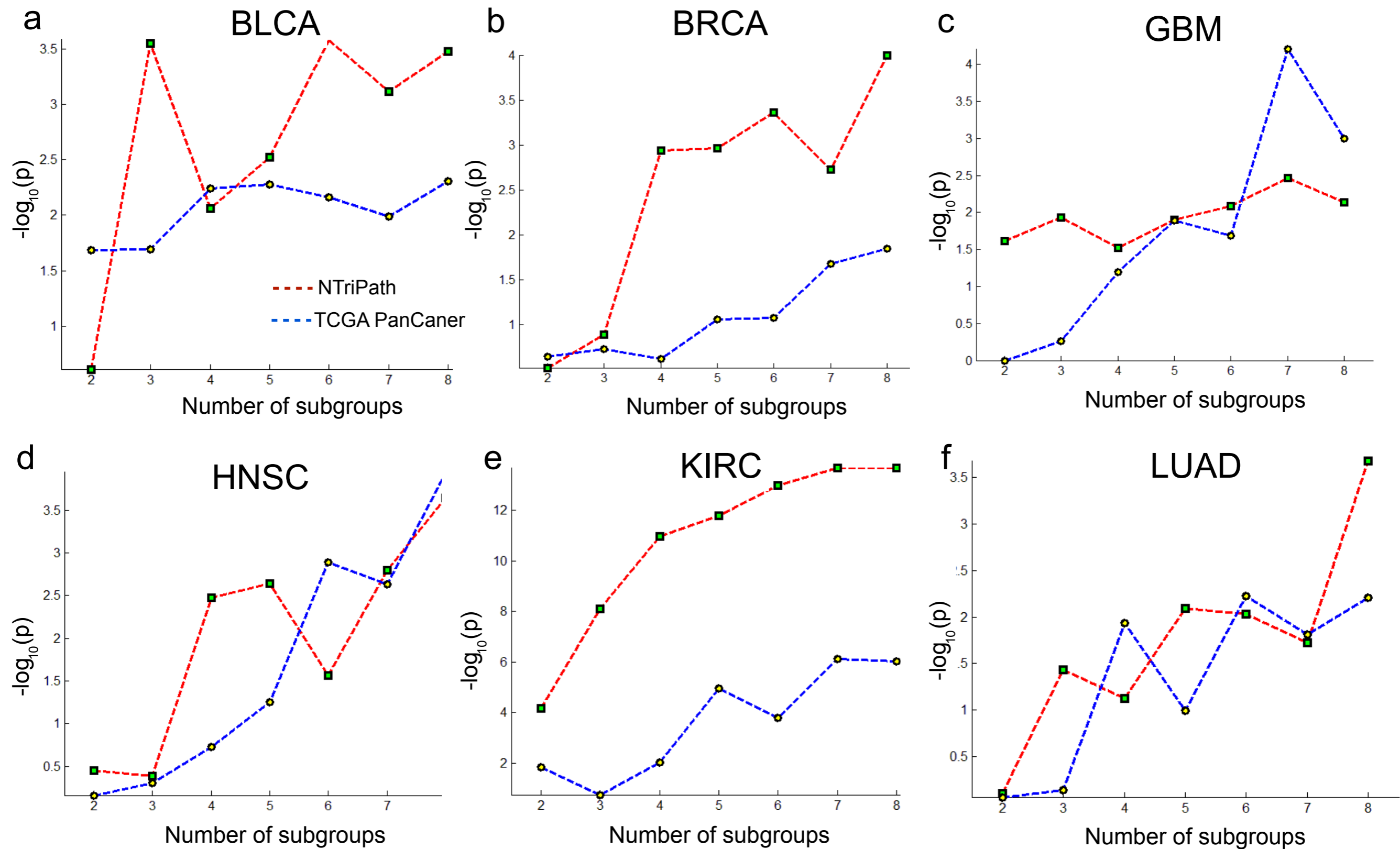
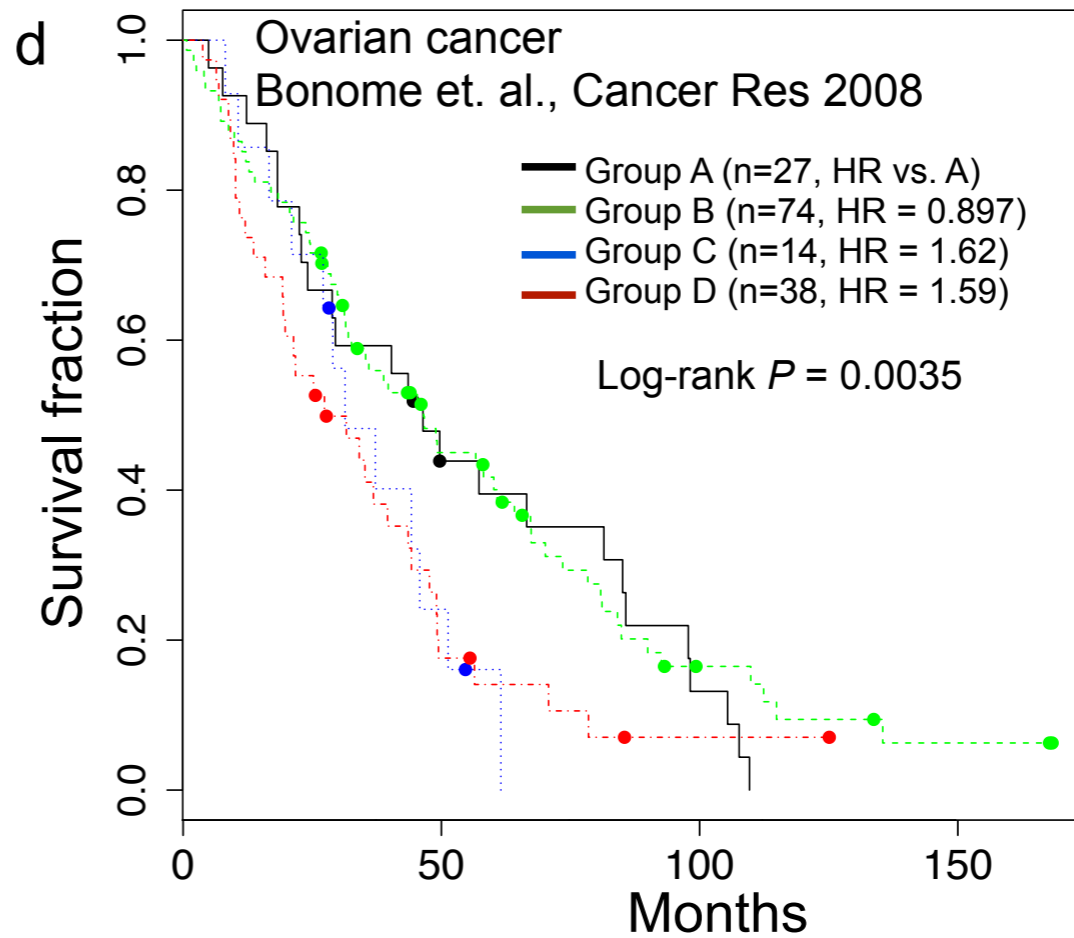
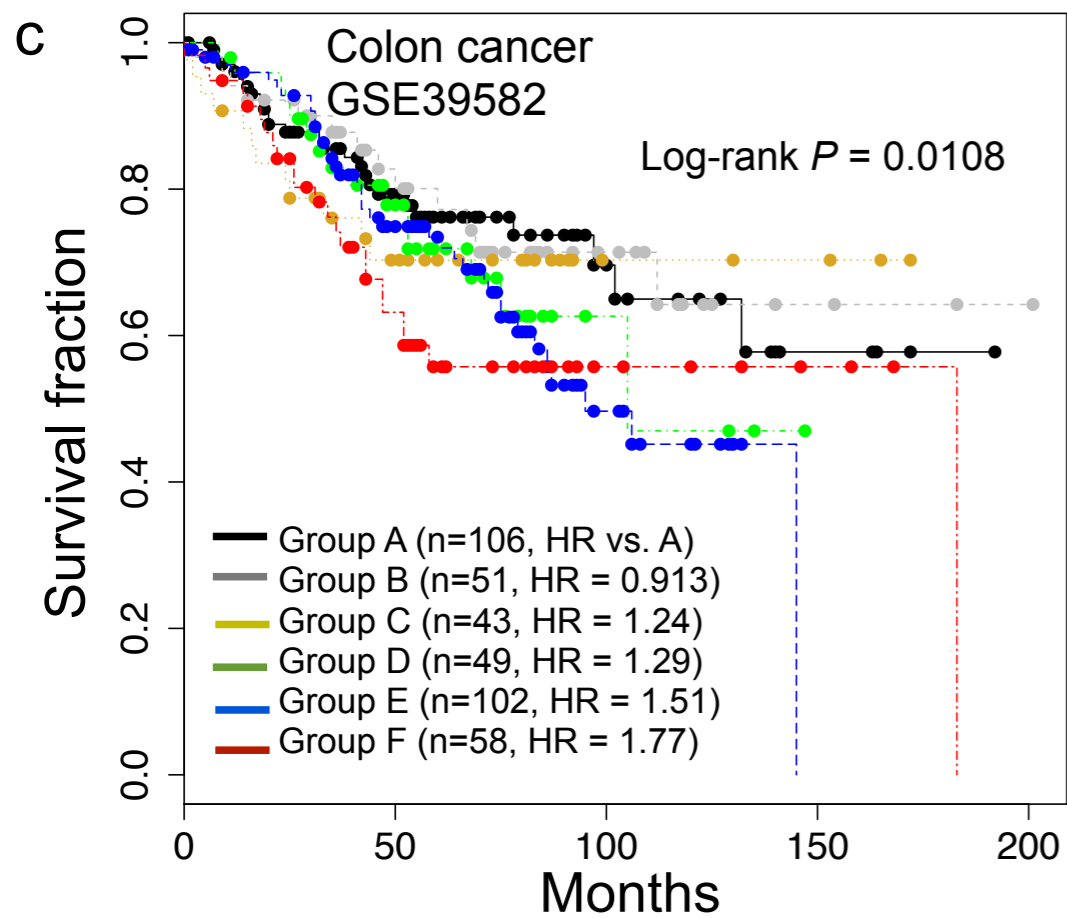
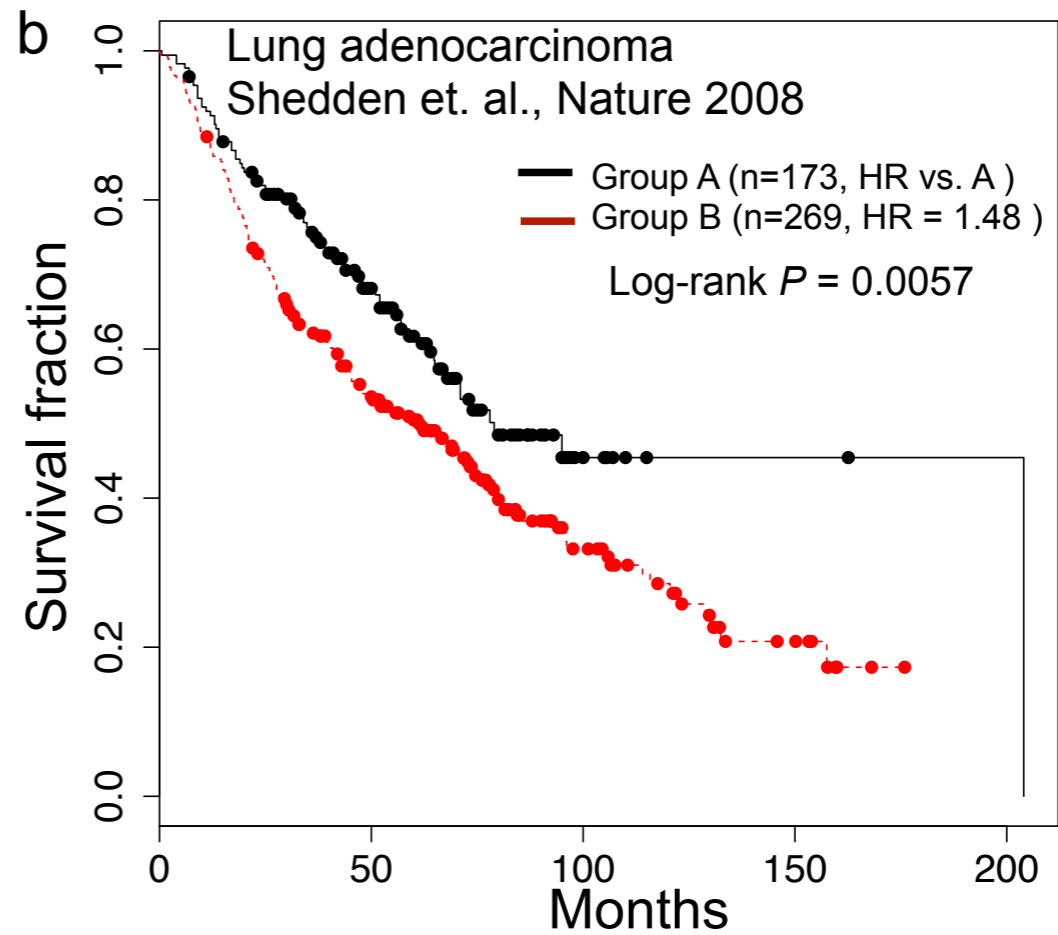
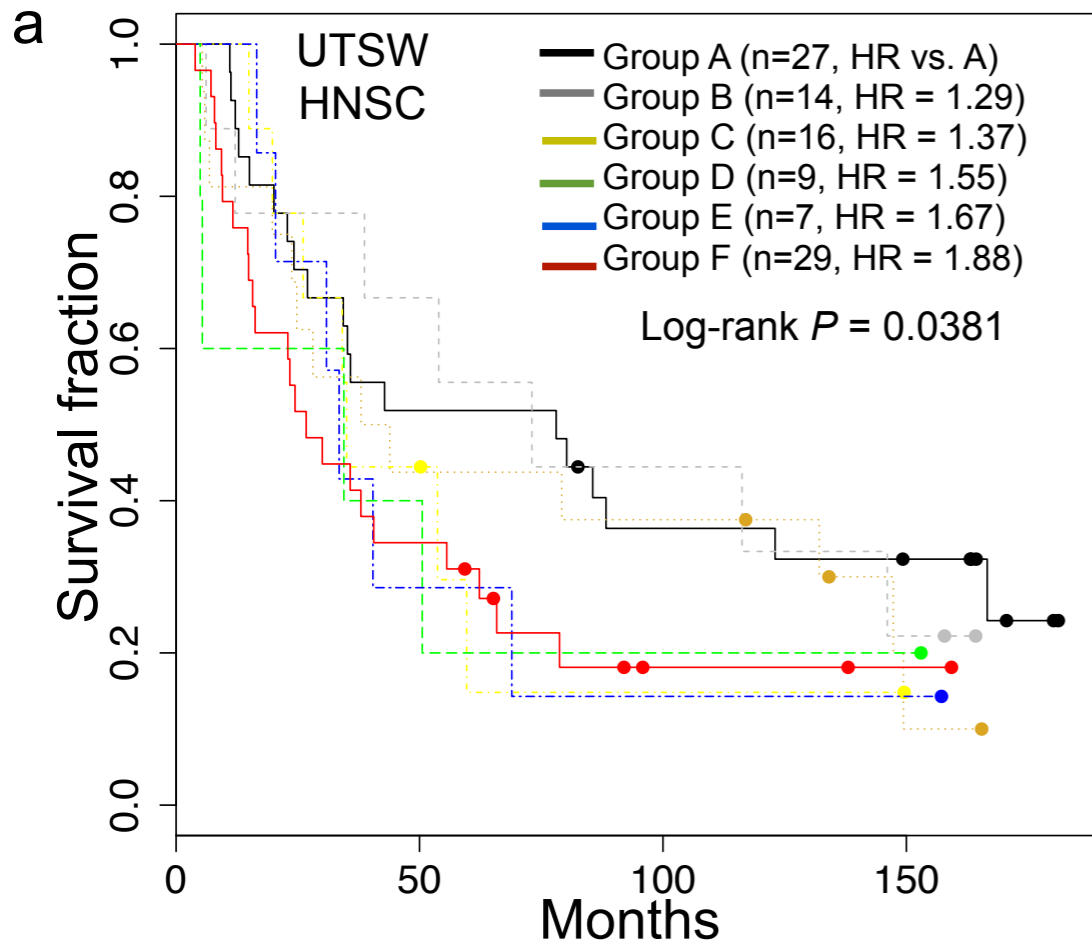
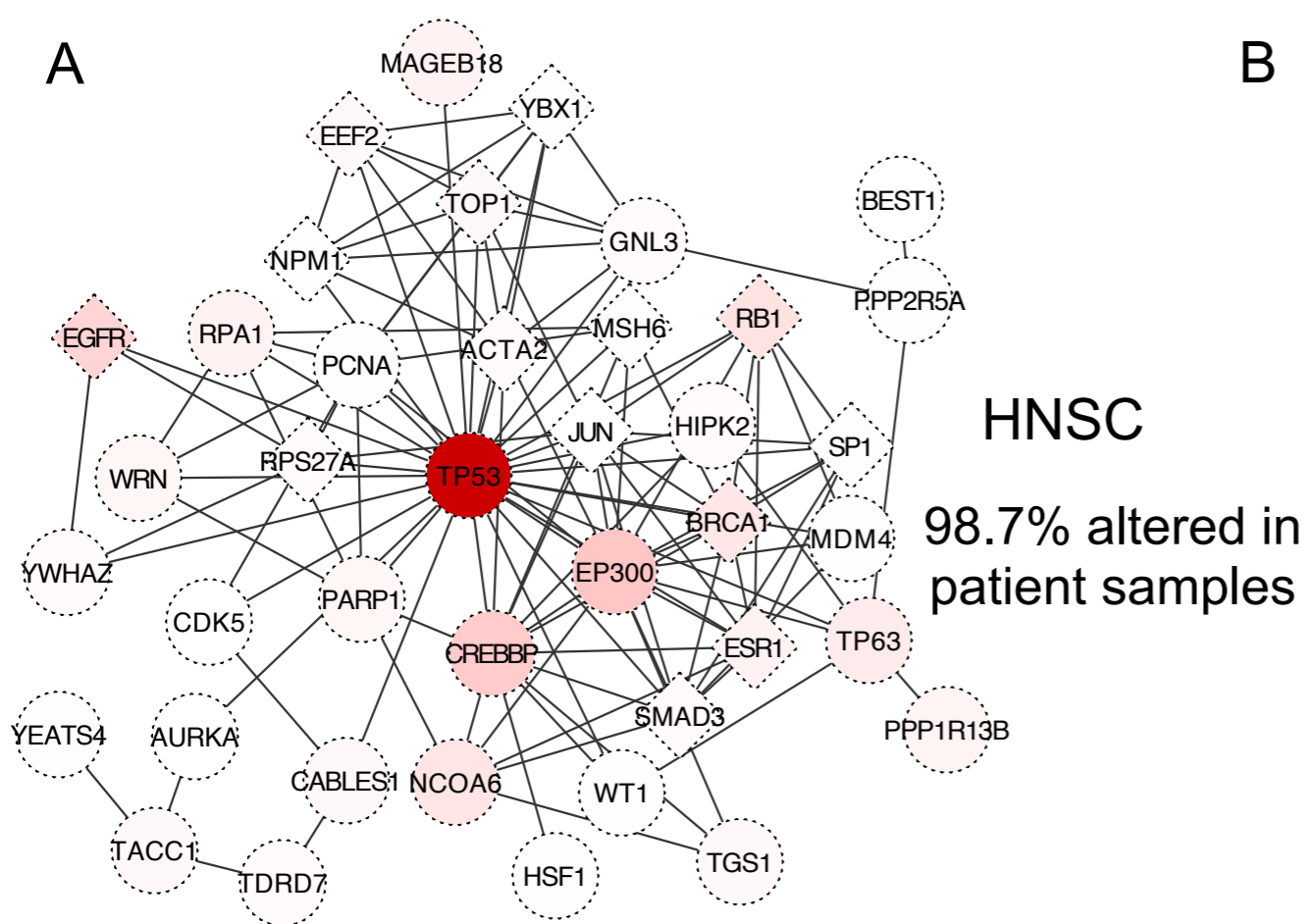
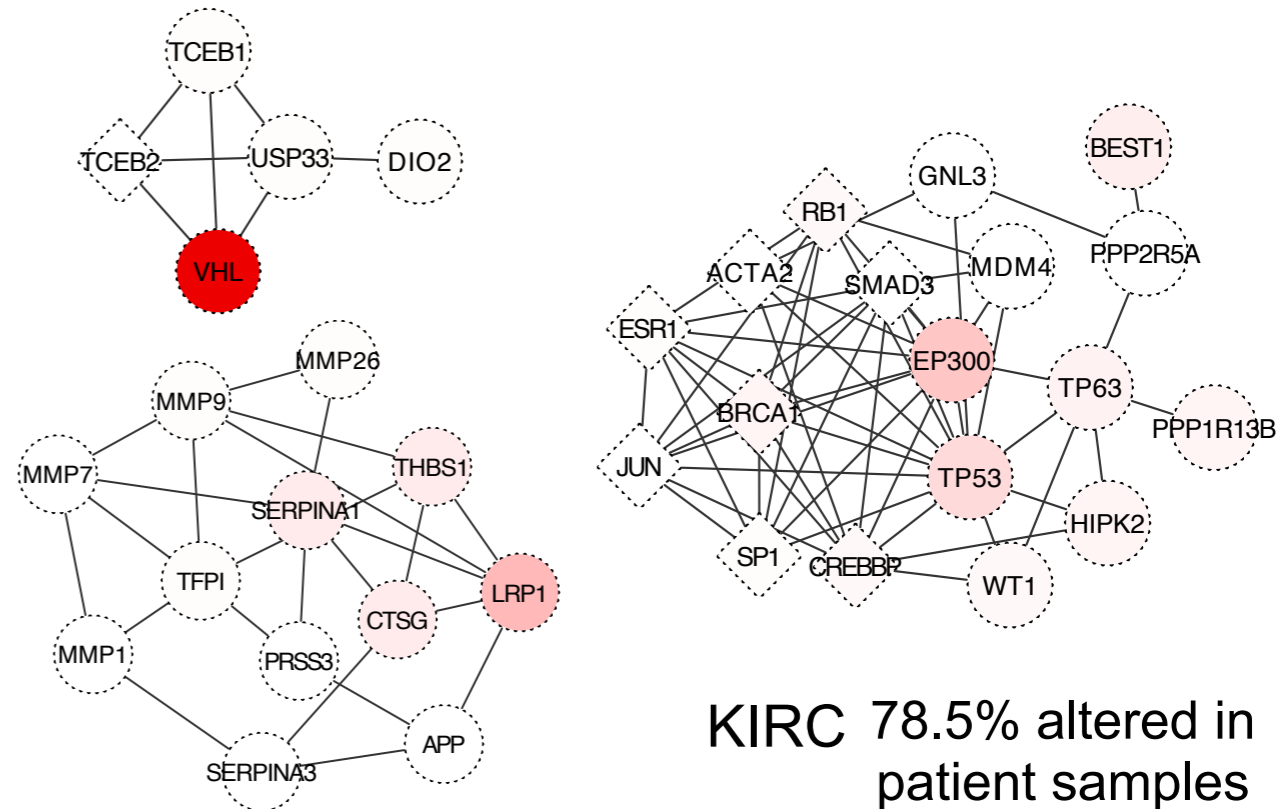
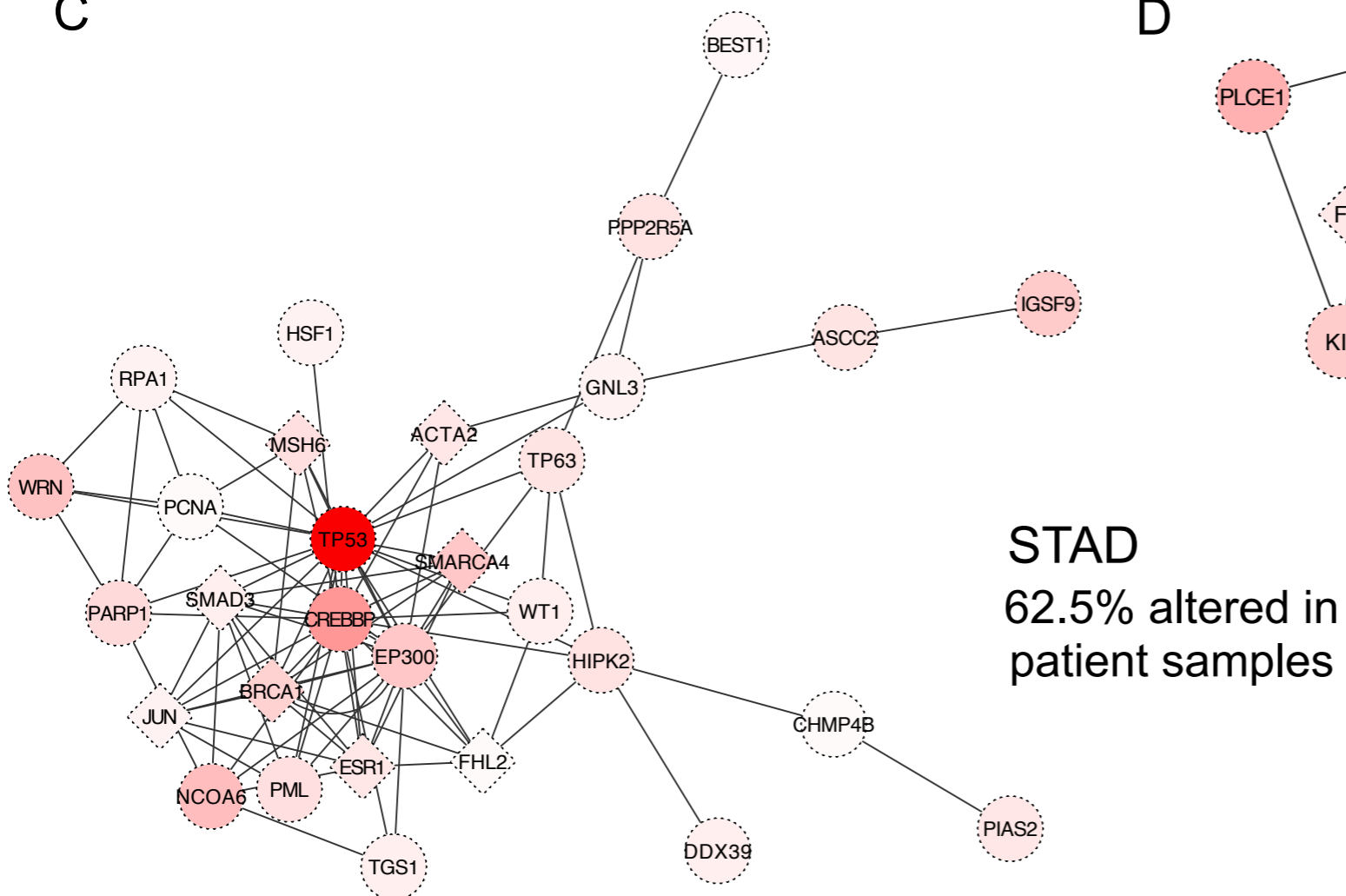
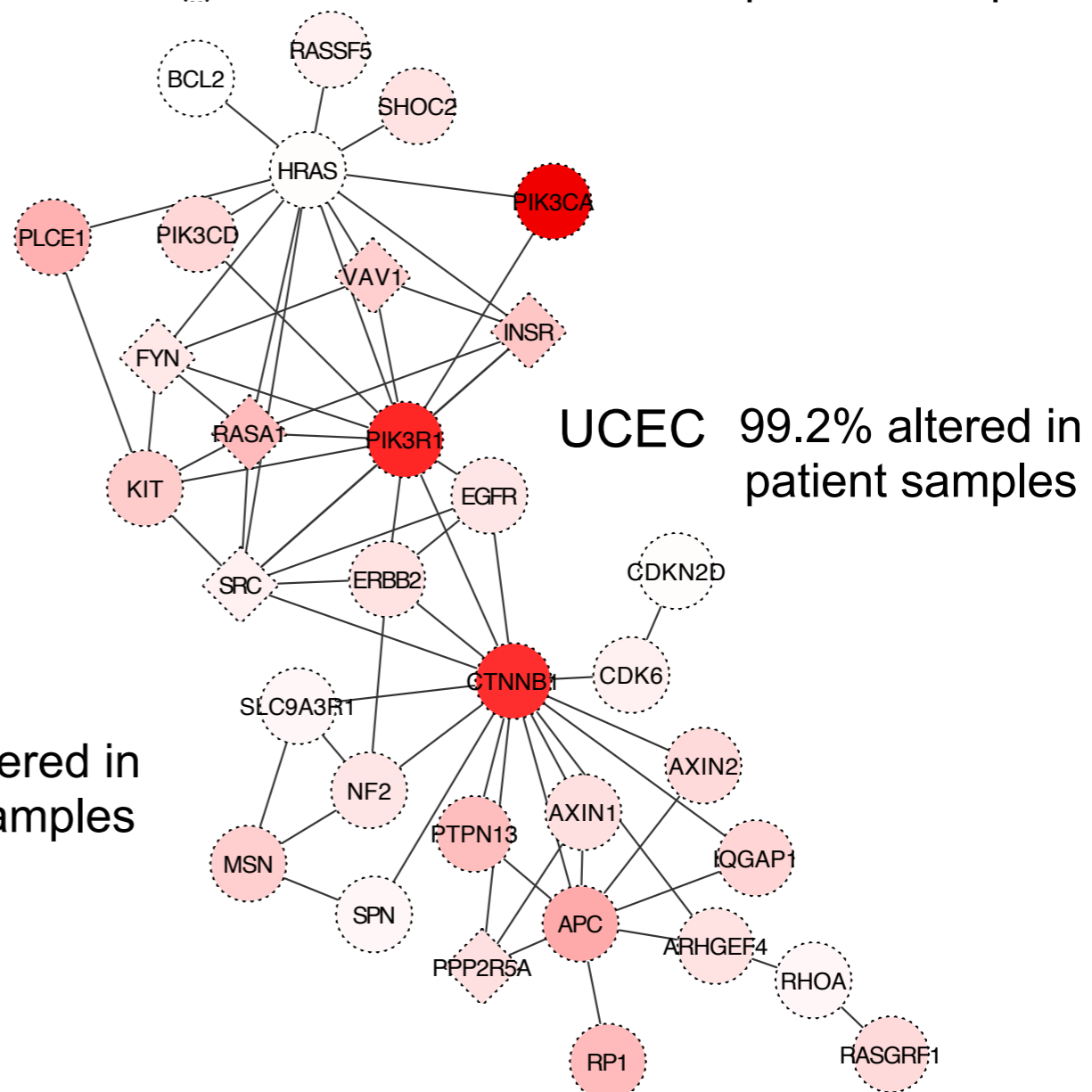
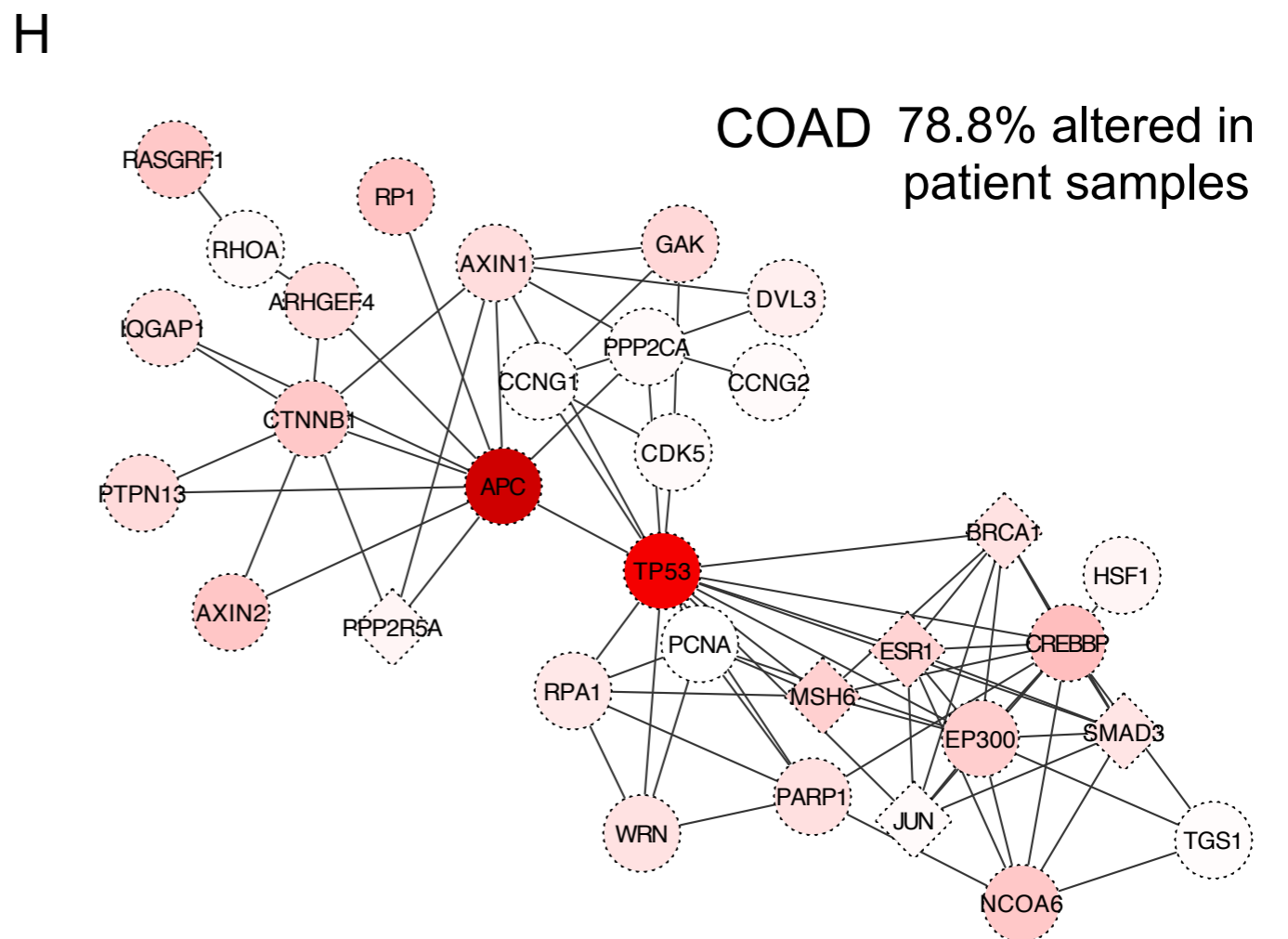
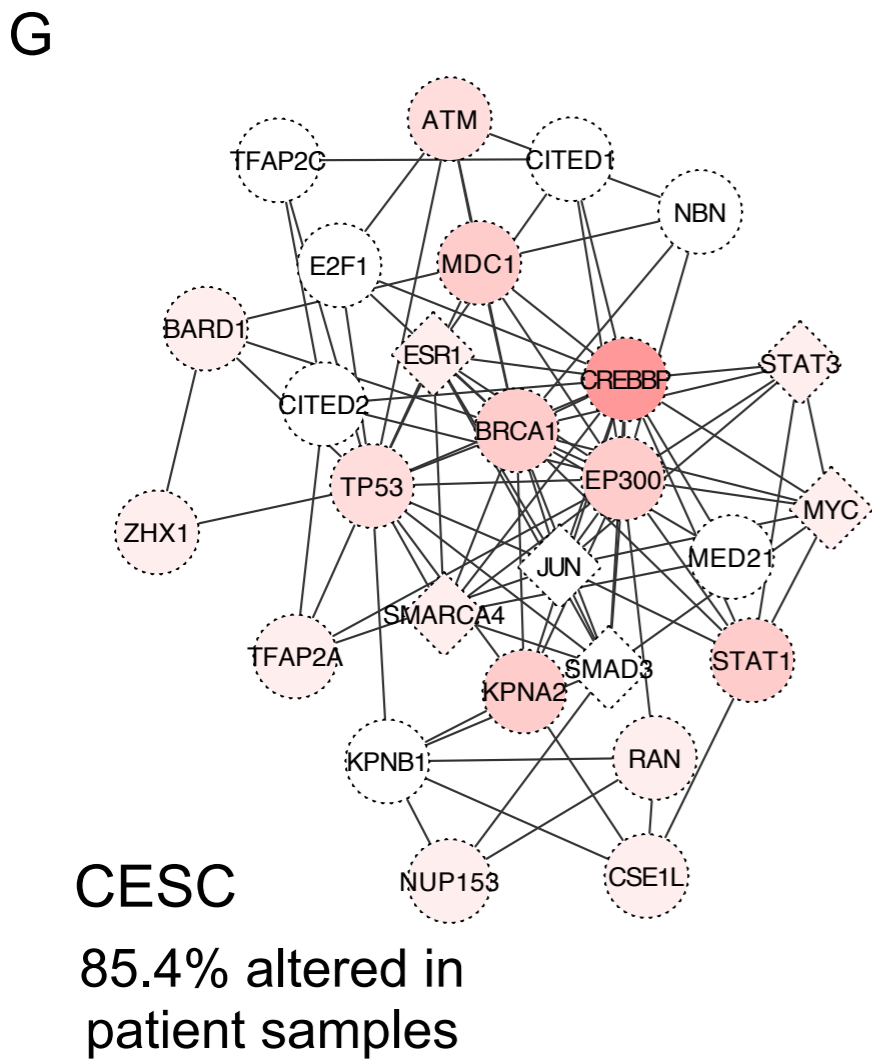
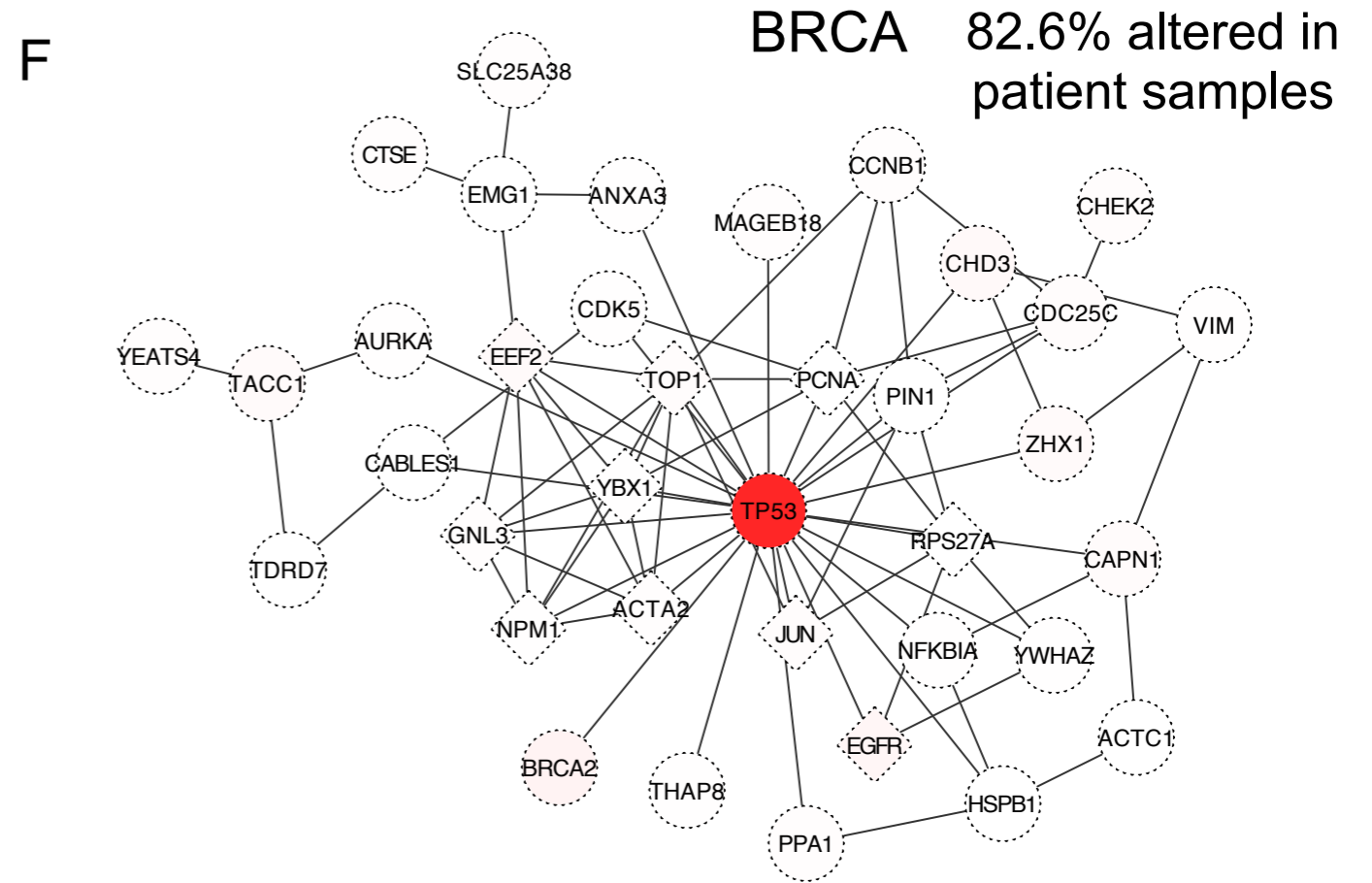
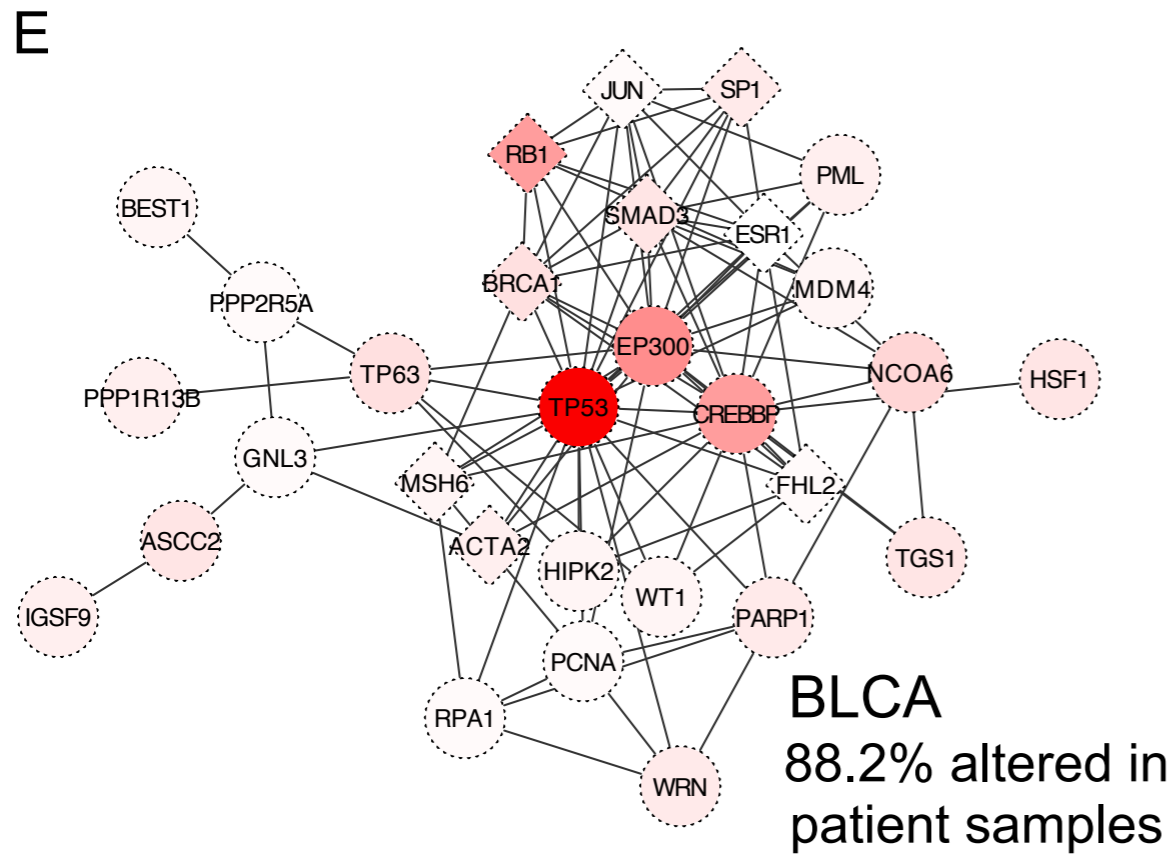


Fig. 5

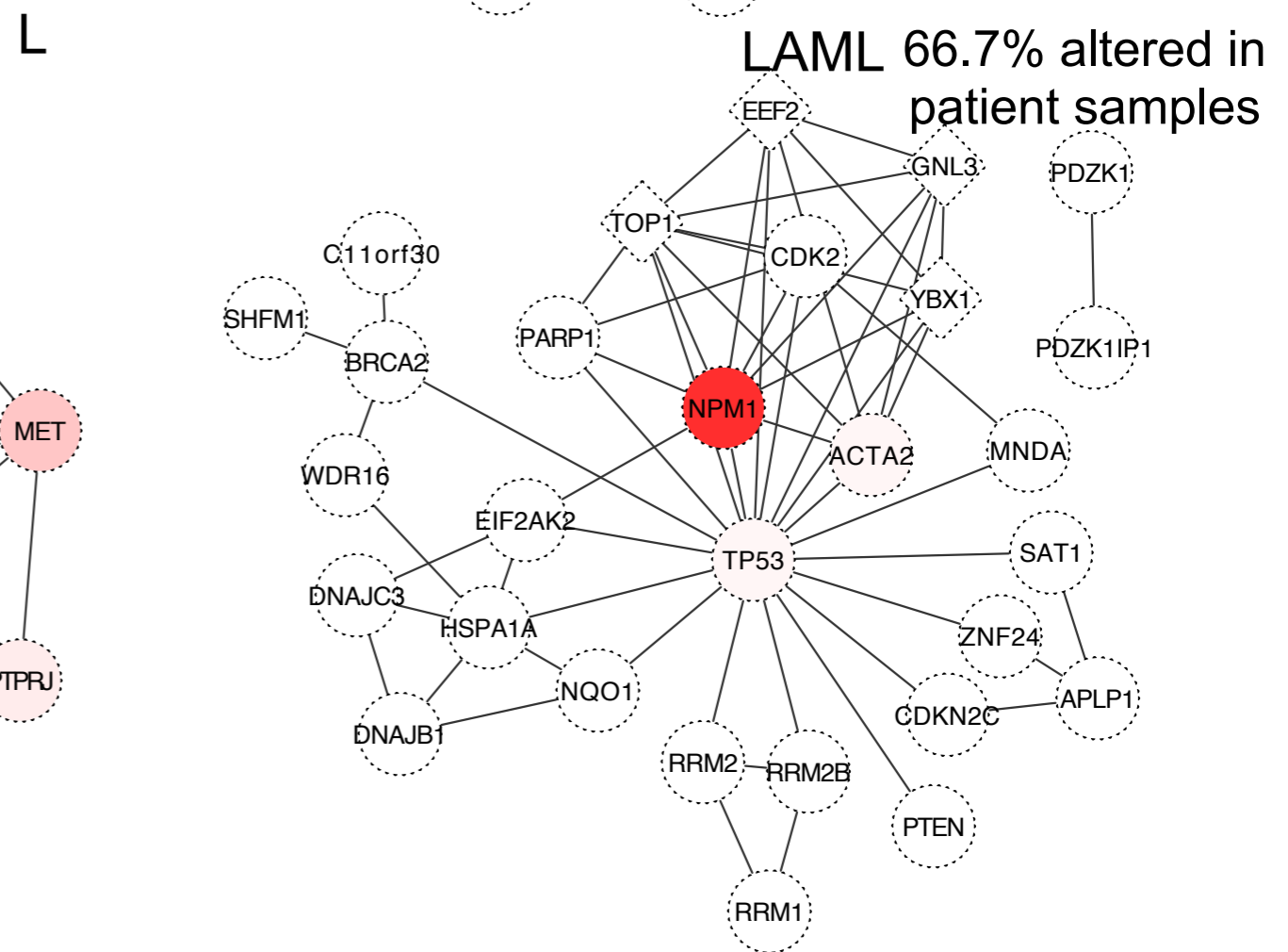
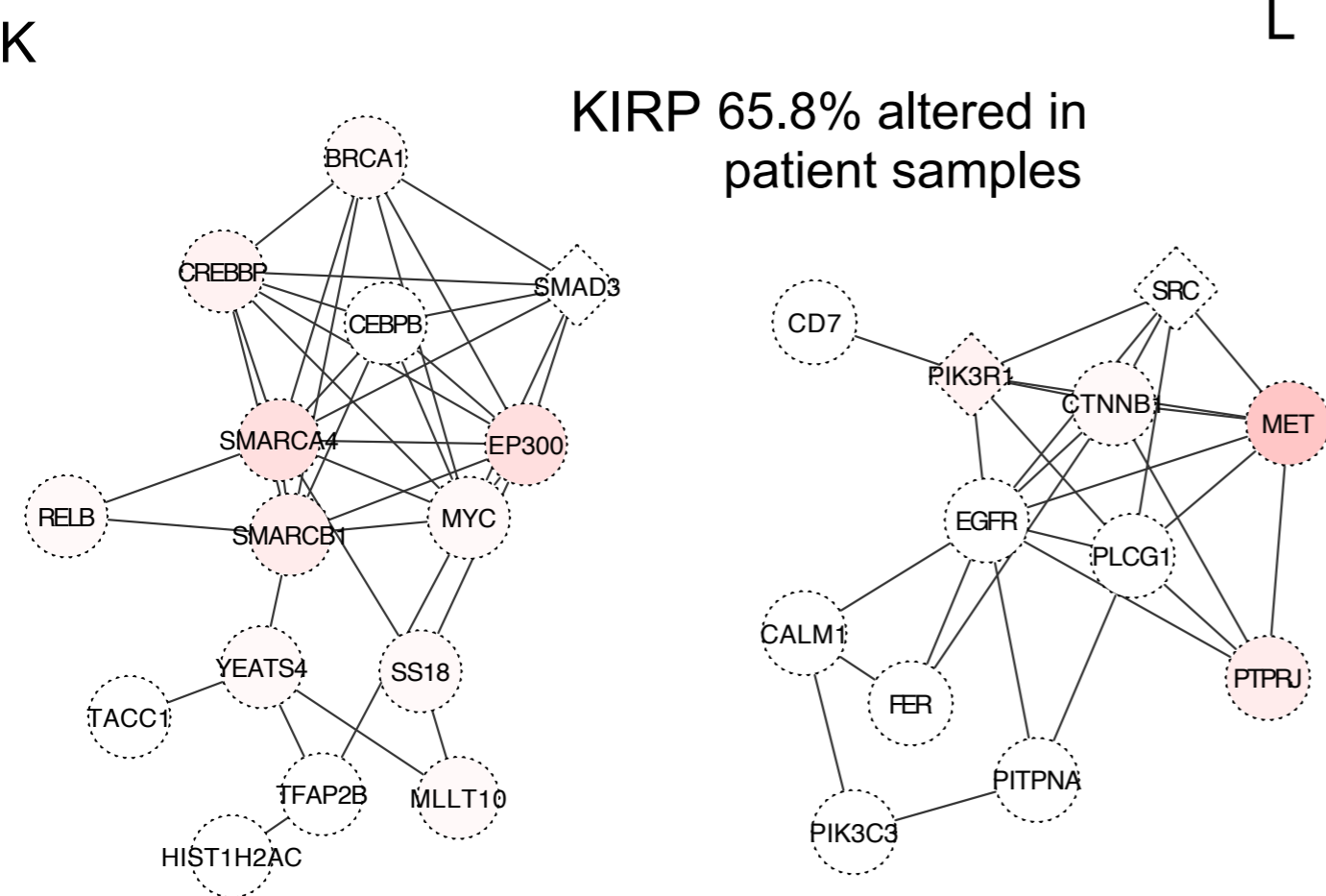
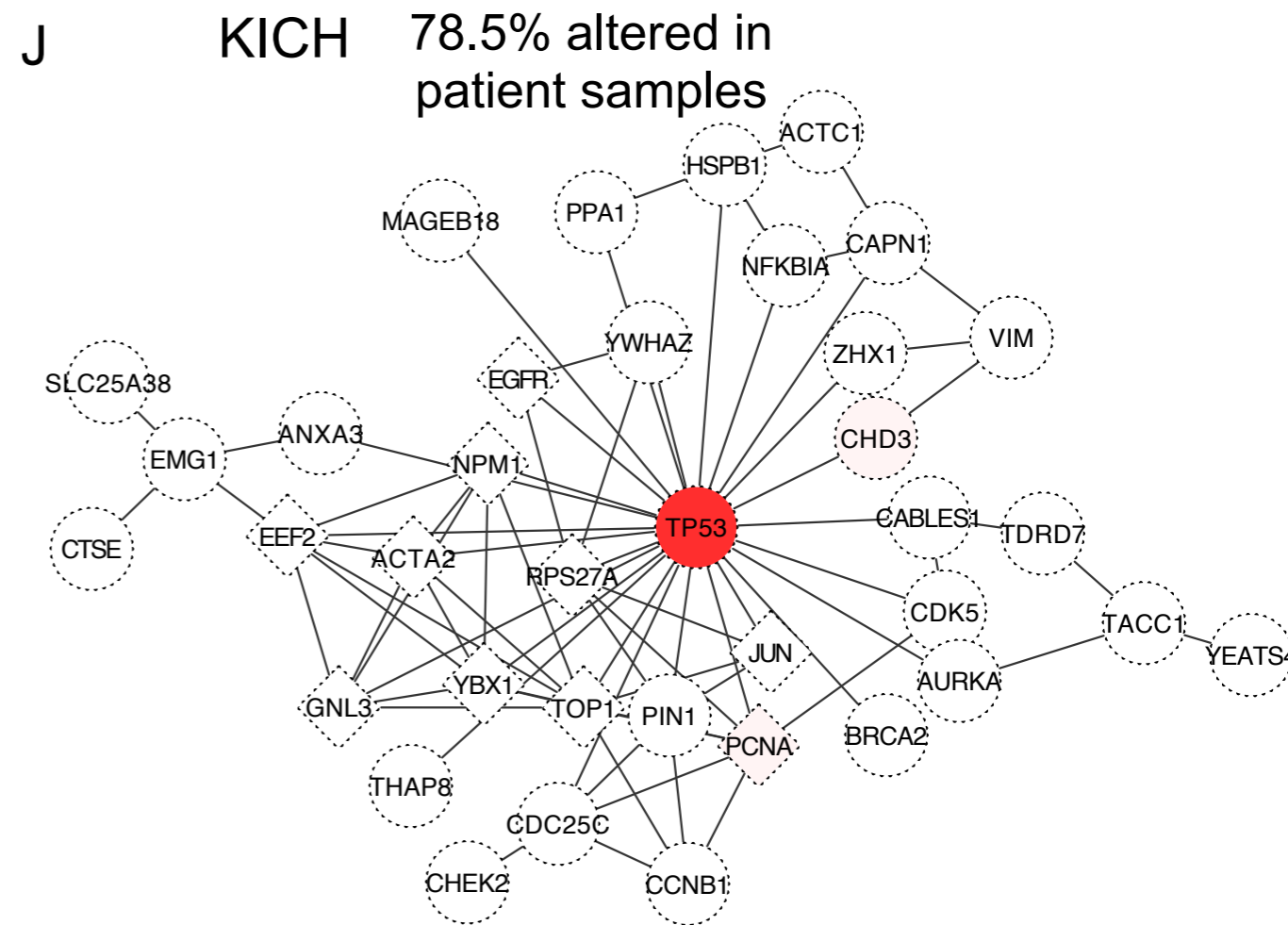
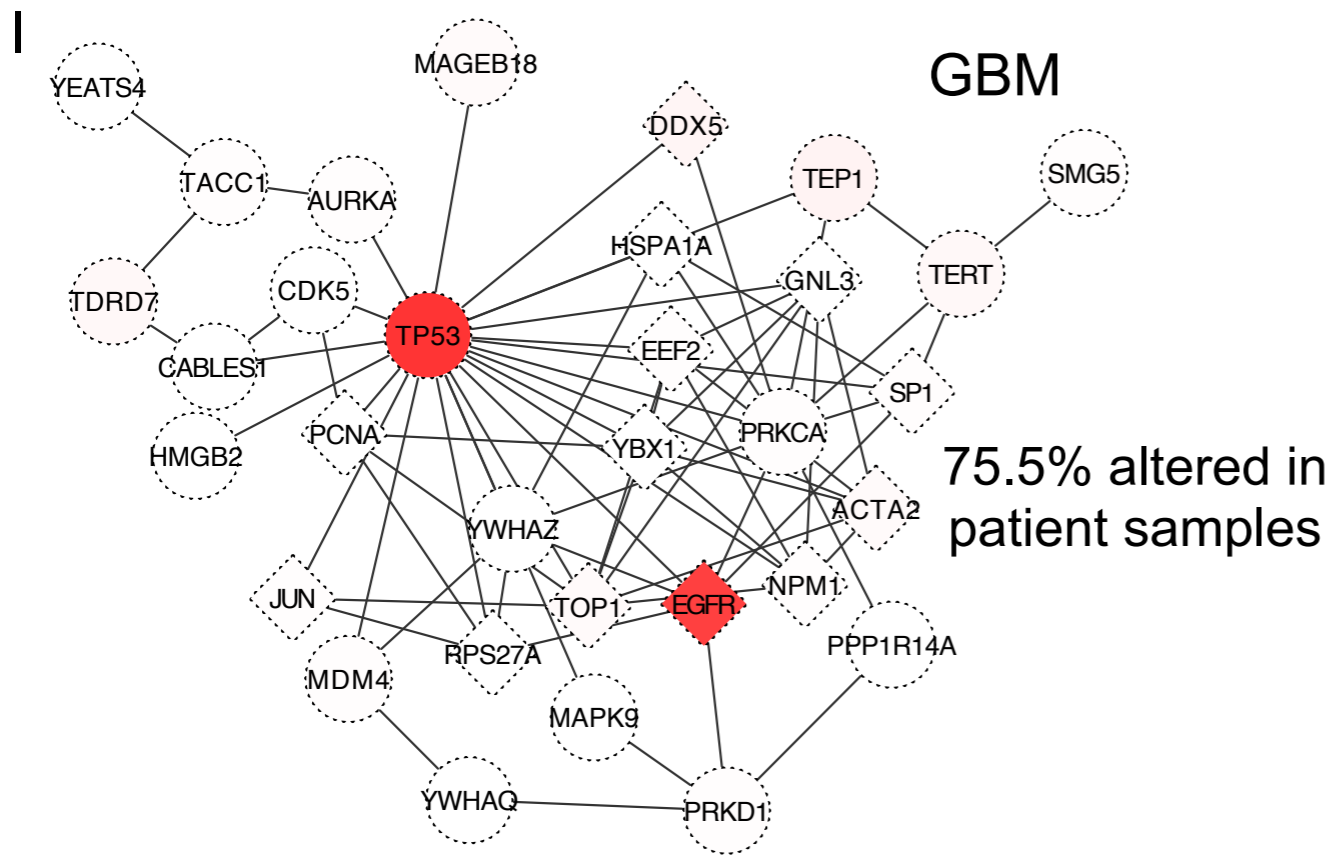


# Supplementary Figure 1

**A****B****C****D**

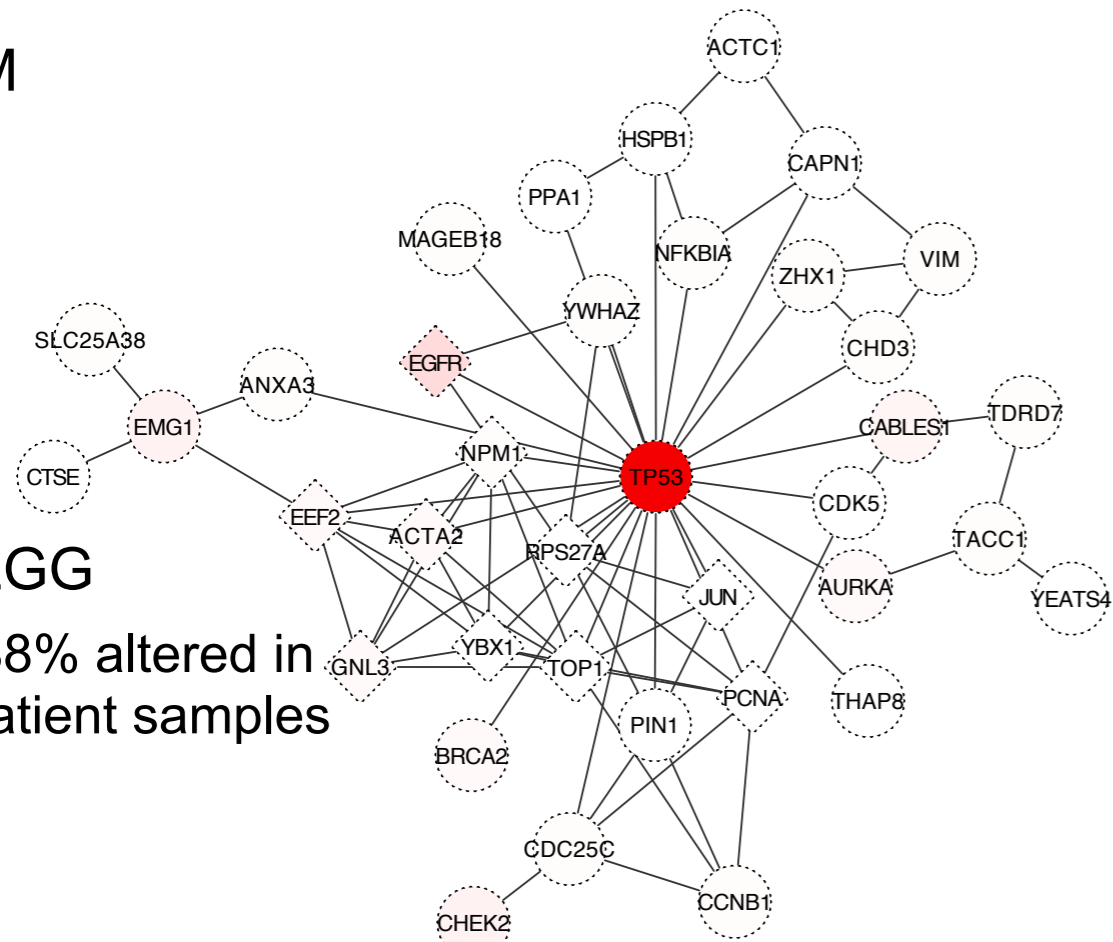




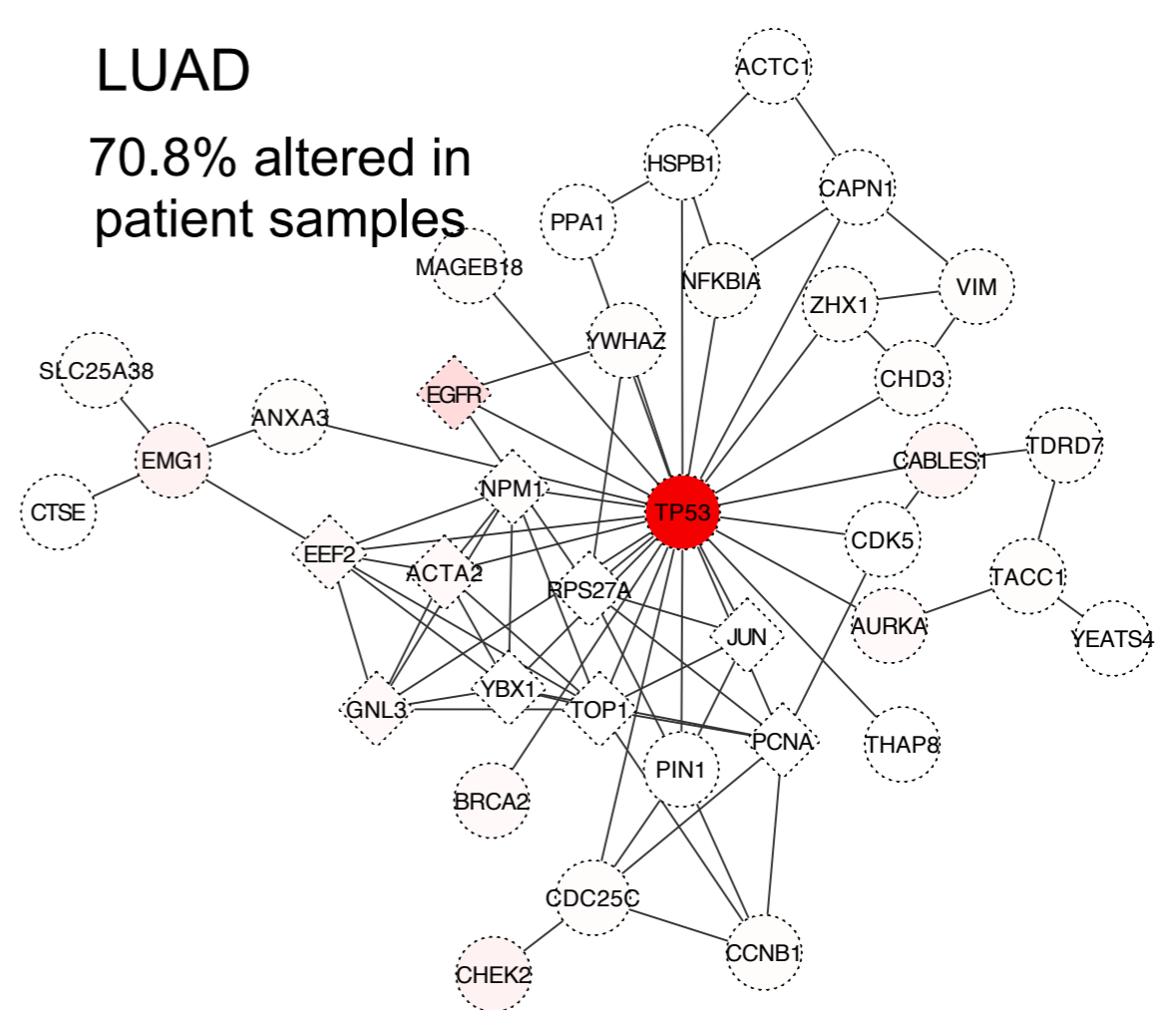


**M****LGG**

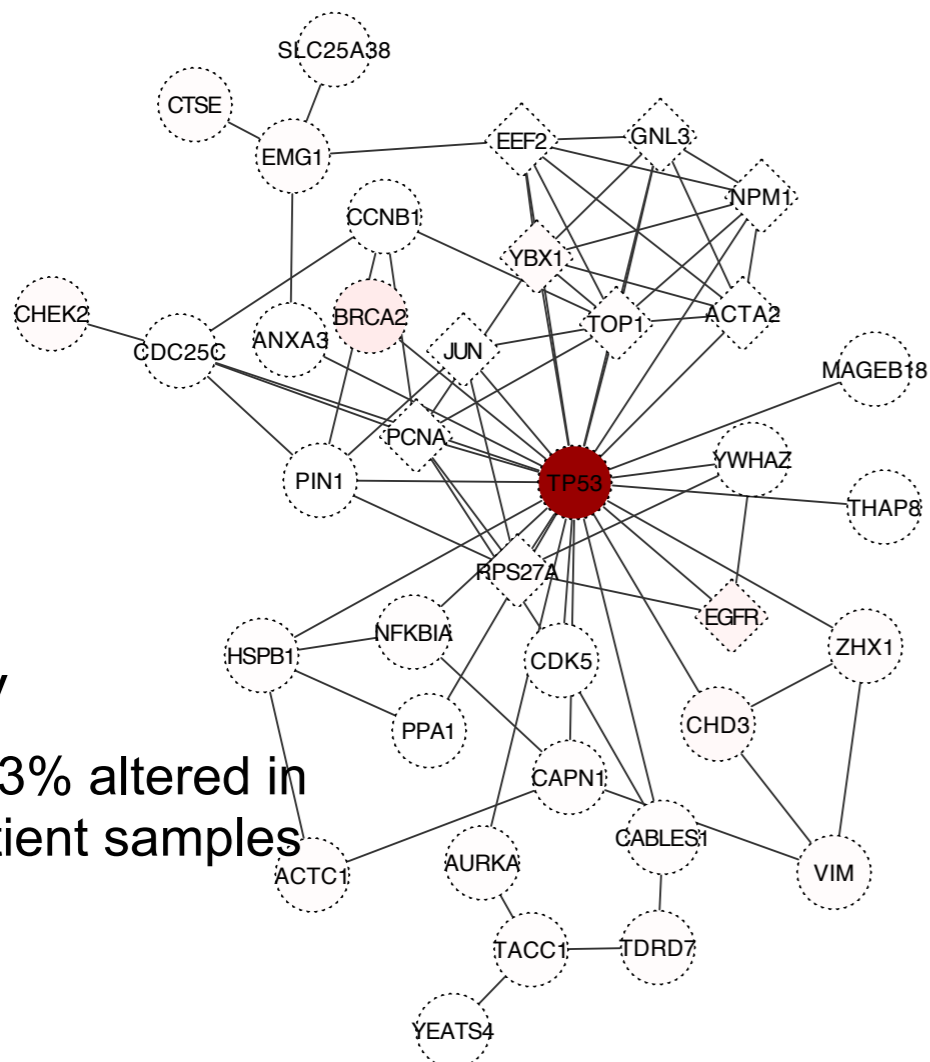
88% altered in patient samples

**N****LUAD**

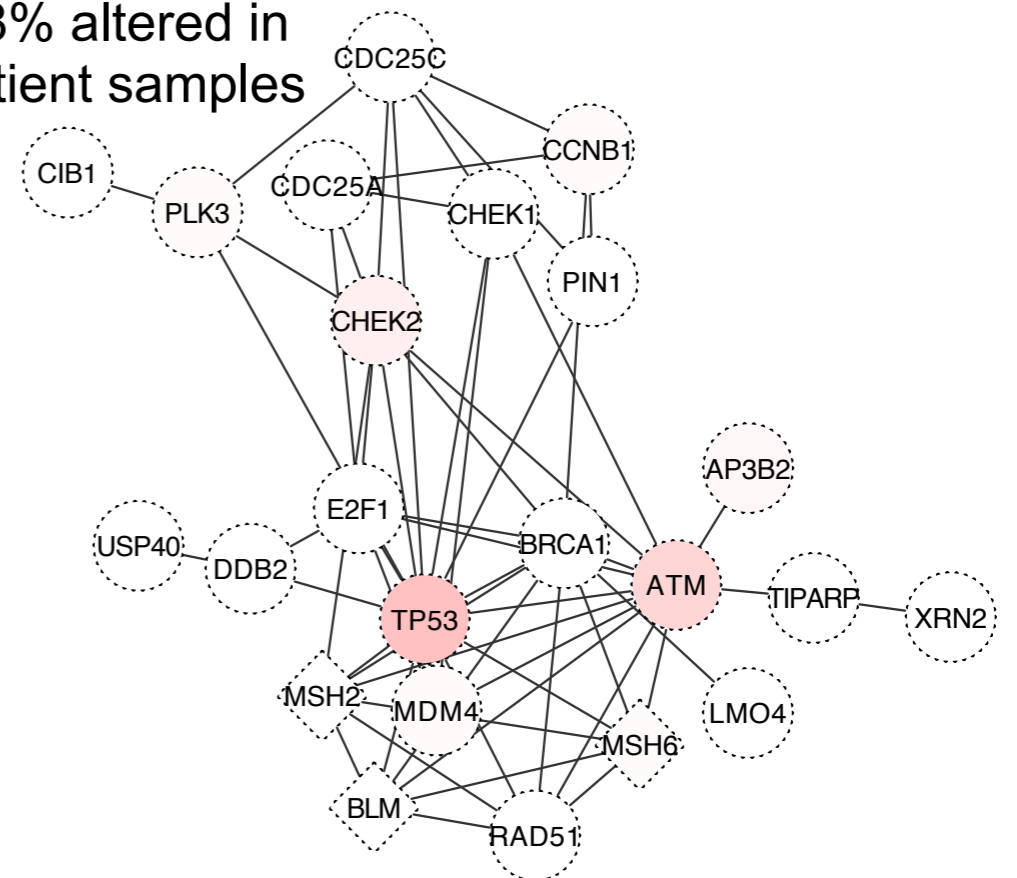
70.8% altered in patient samples

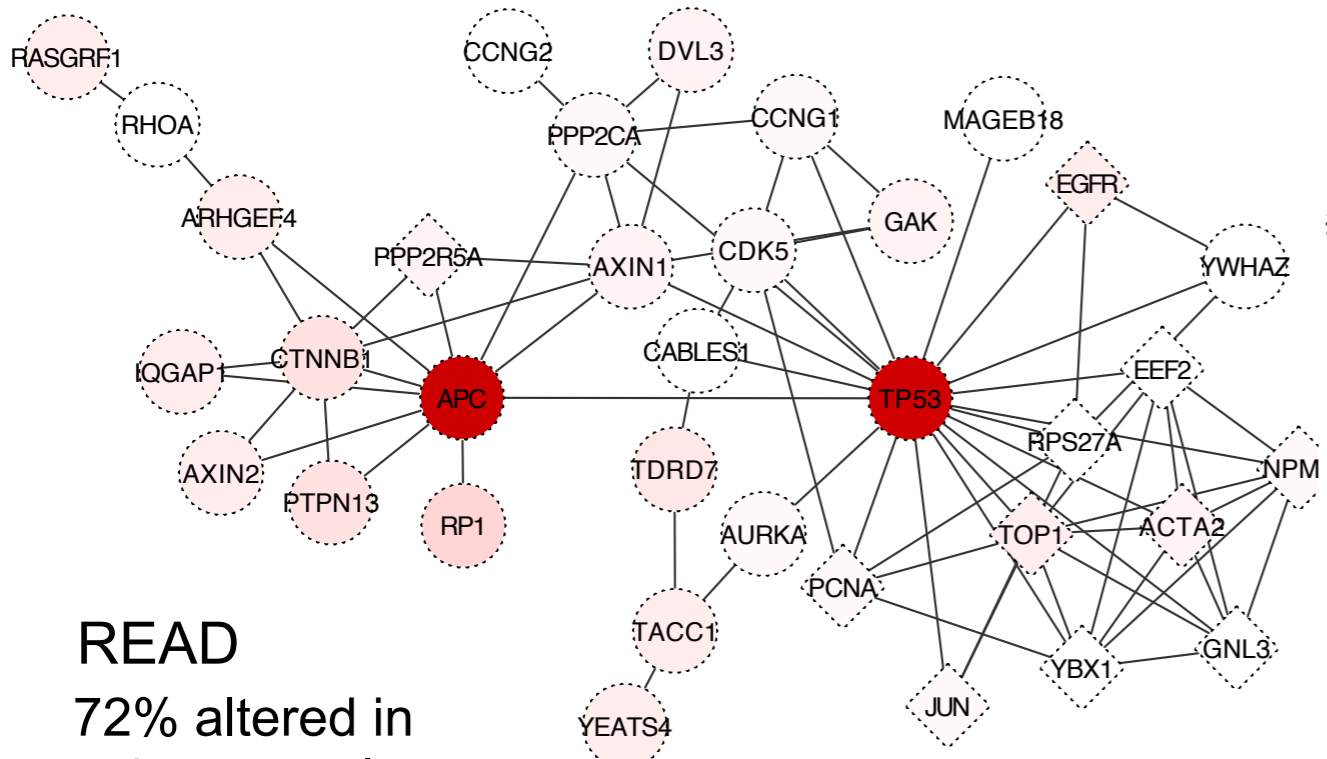
**O****OV**

80.3% altered in patient samples

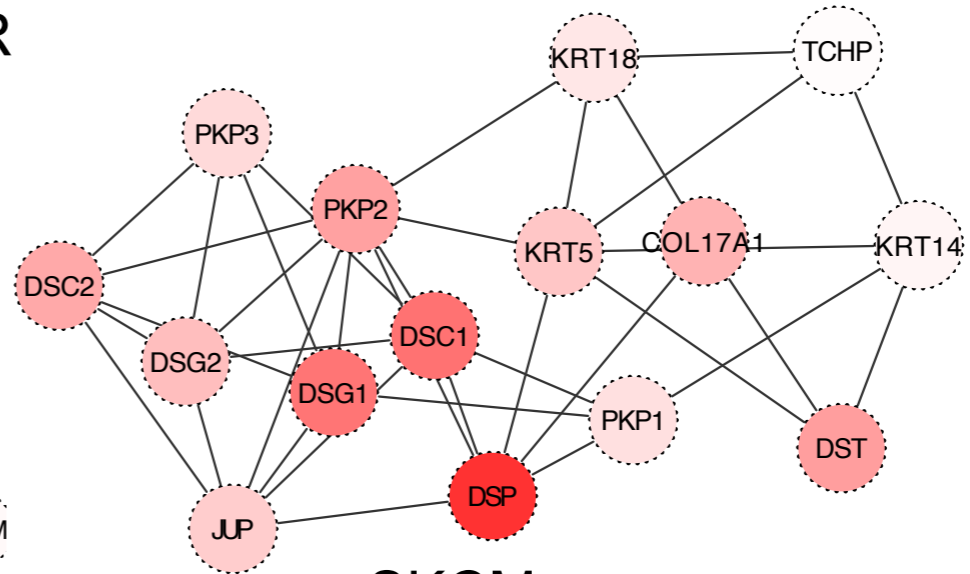
**P****PRAD**

48% altered in patient samples

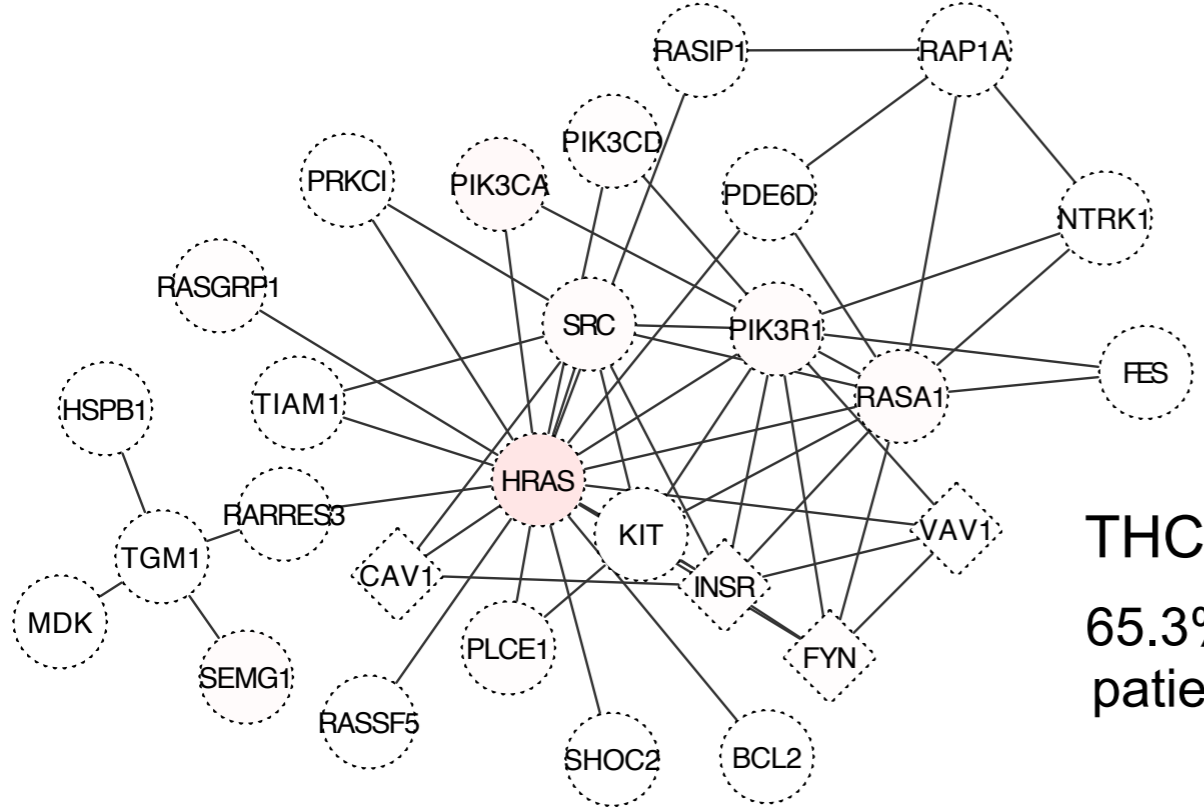


**Q**

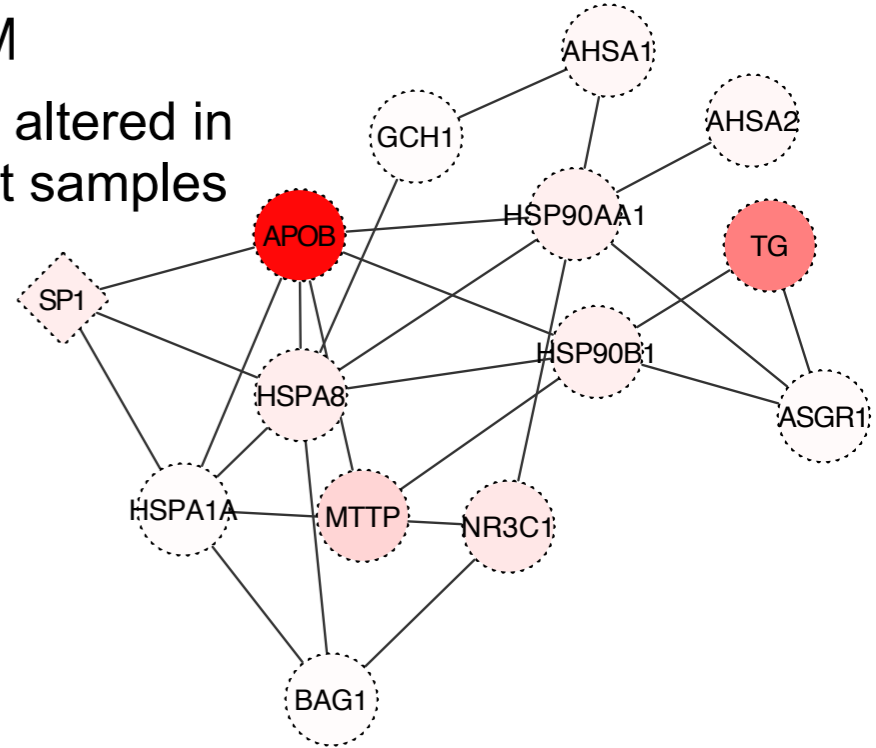
**READ**  
72% altered in patient samples

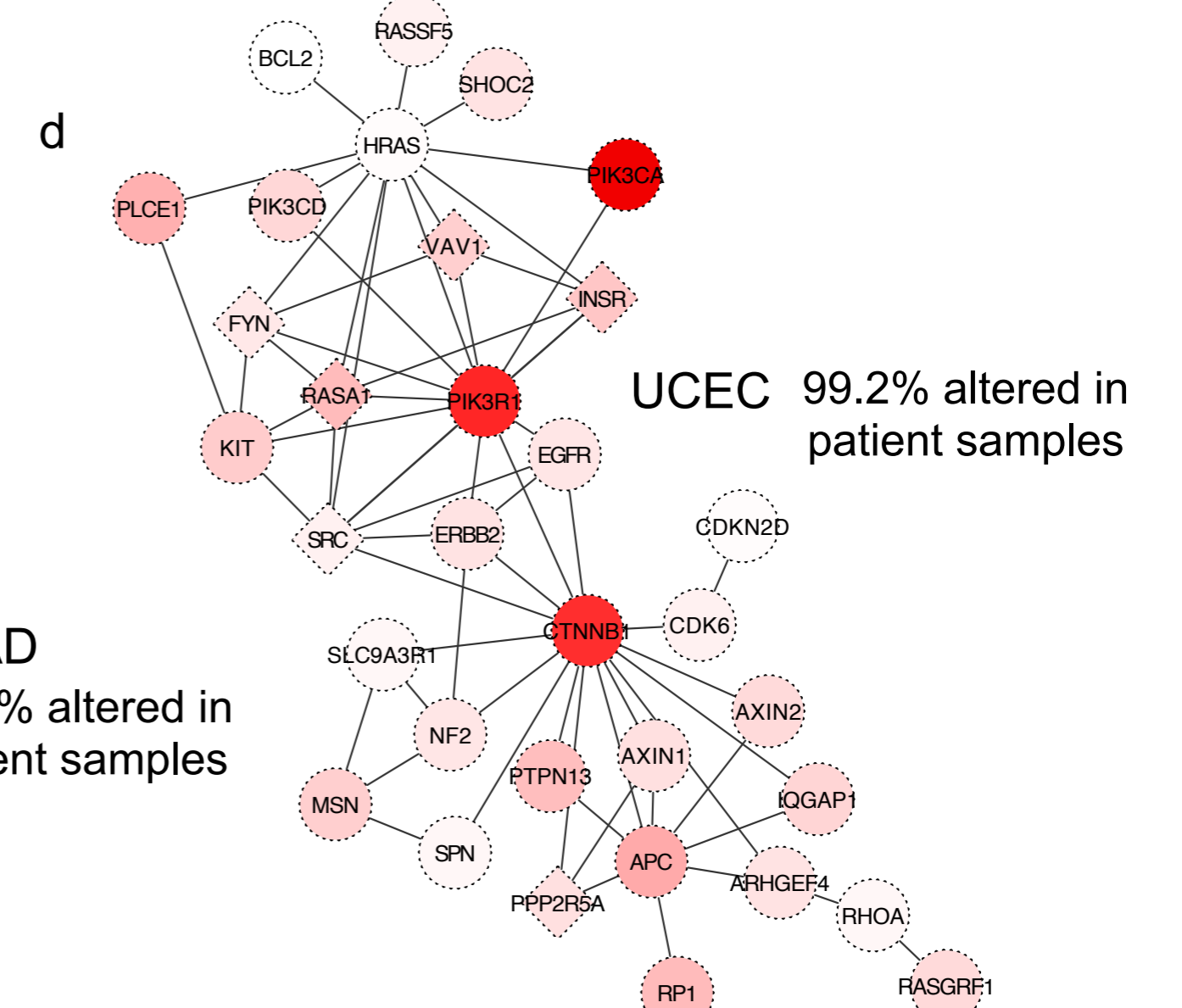
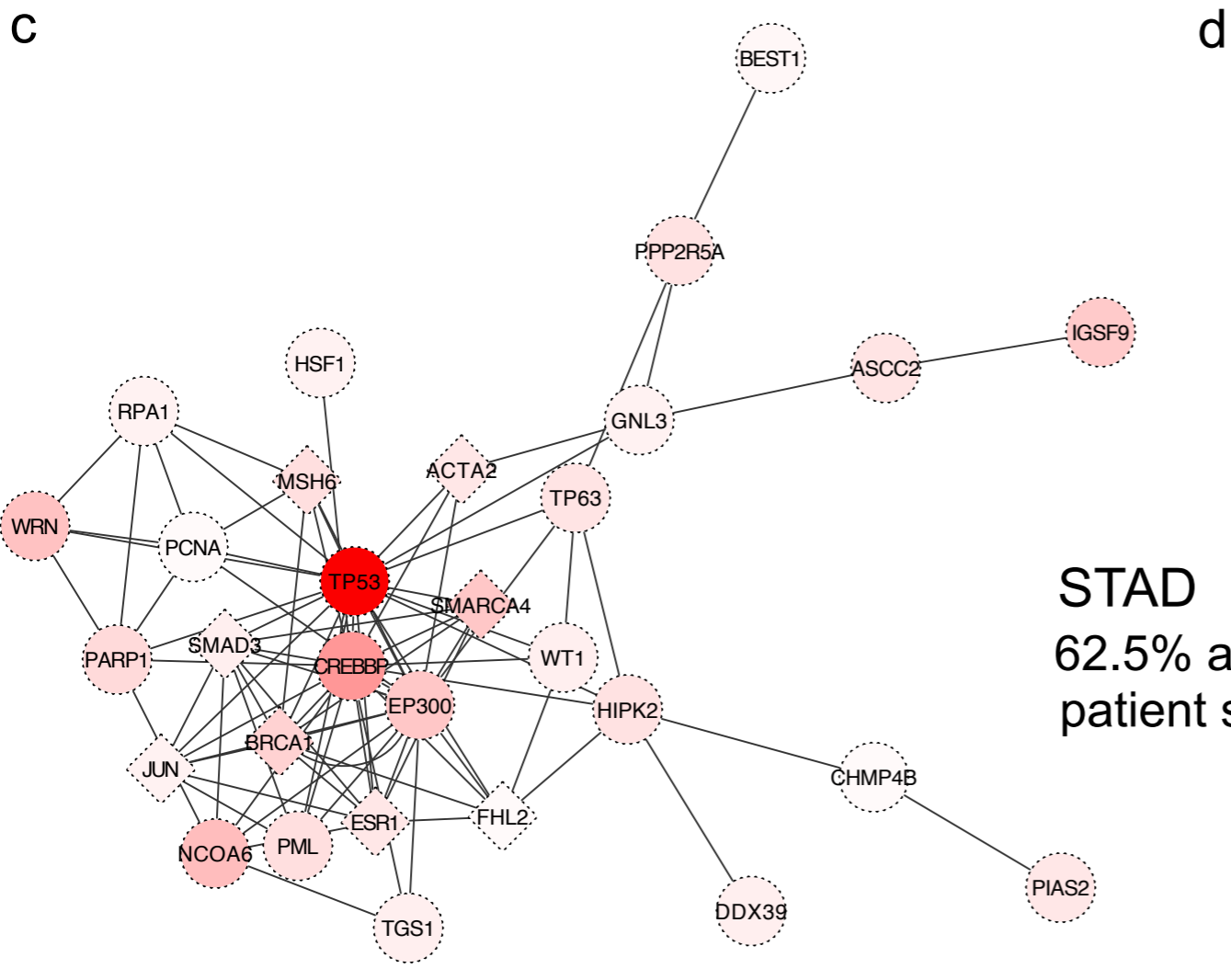
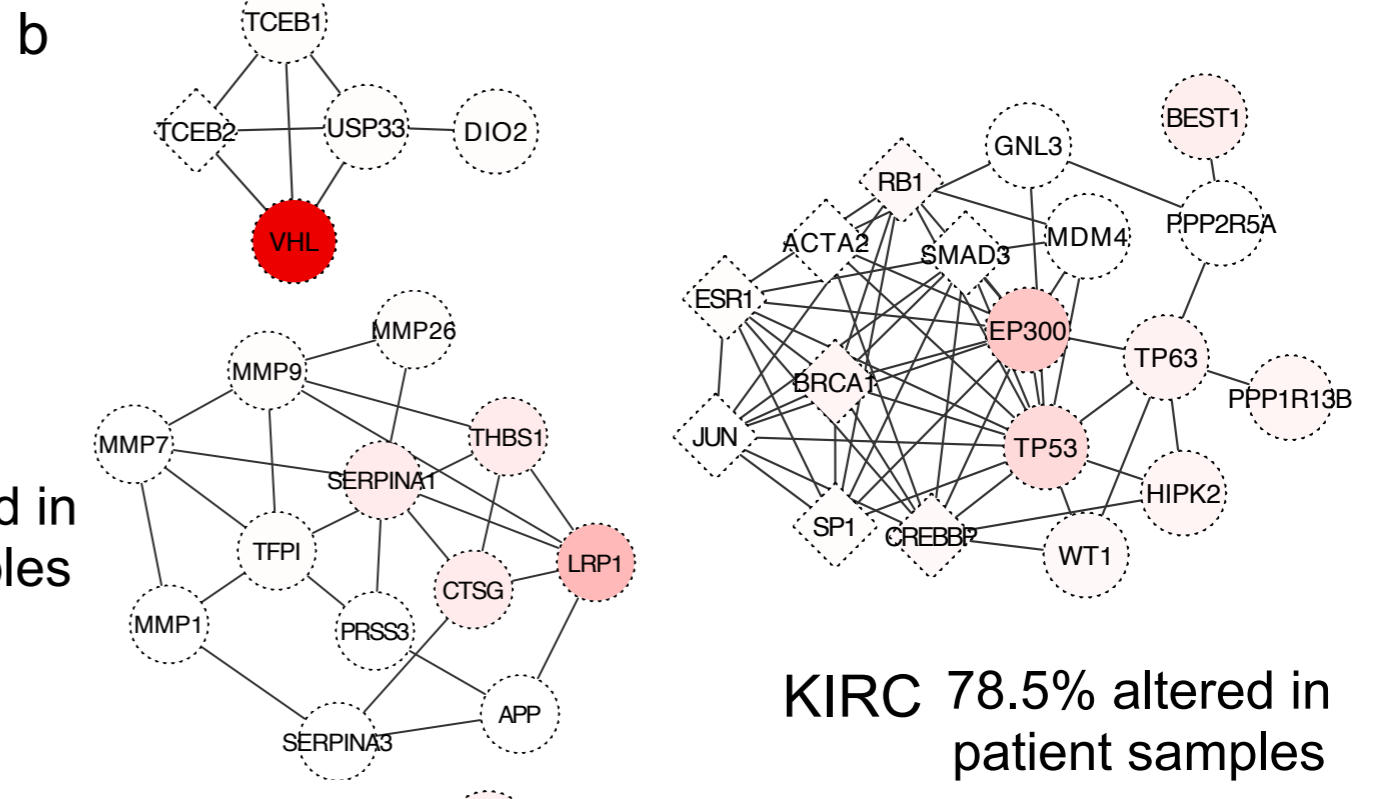
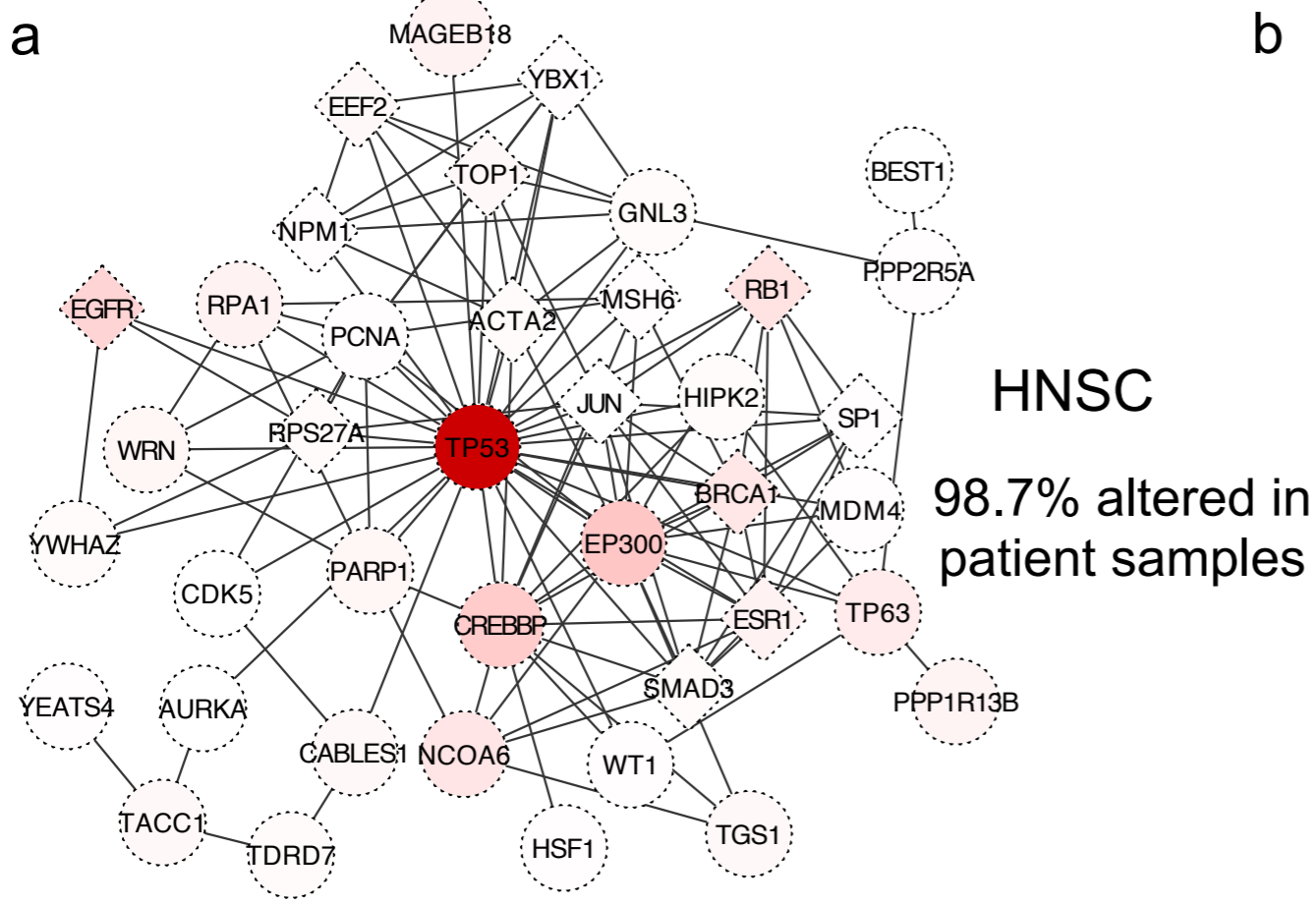
**R**

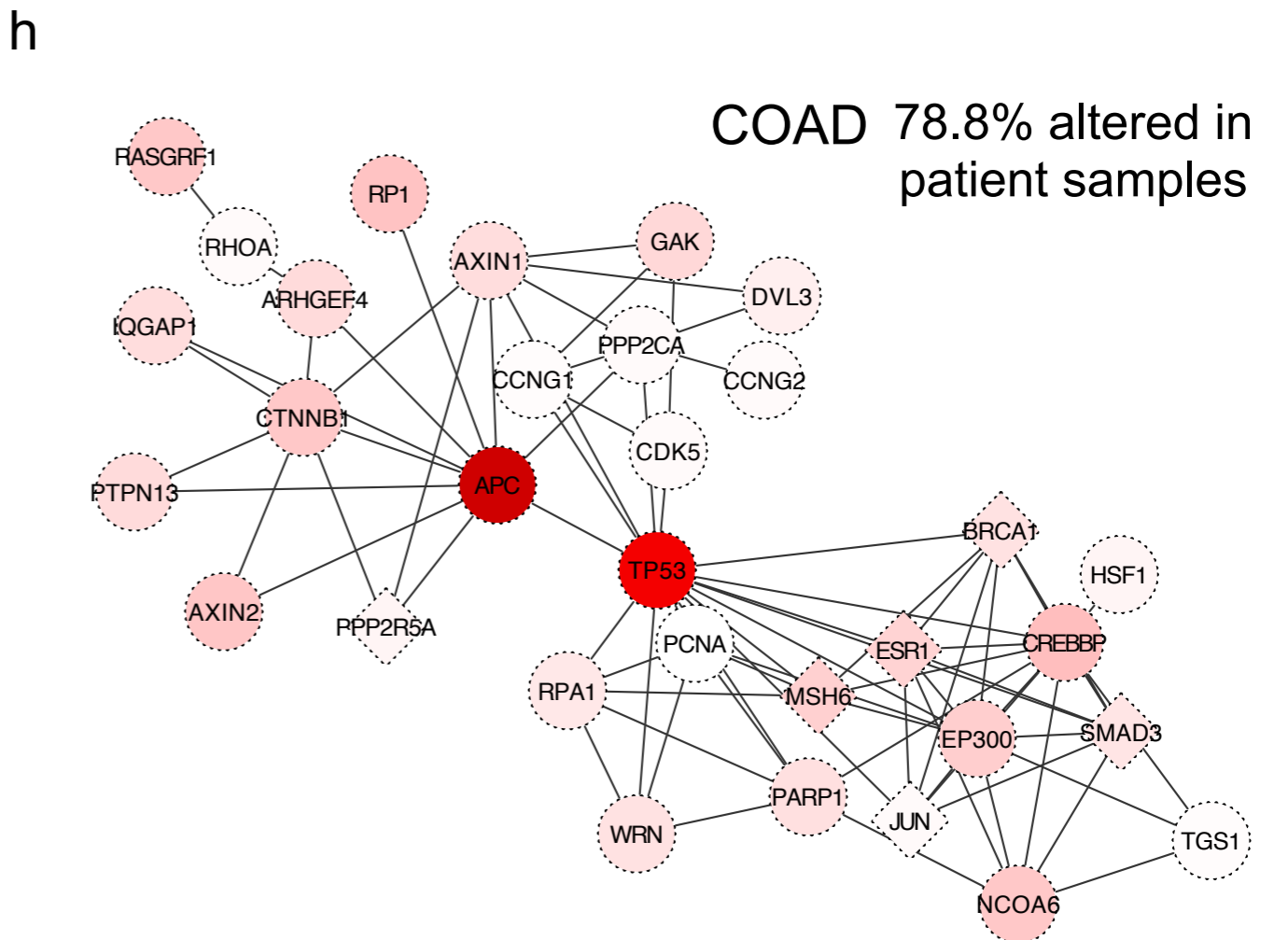
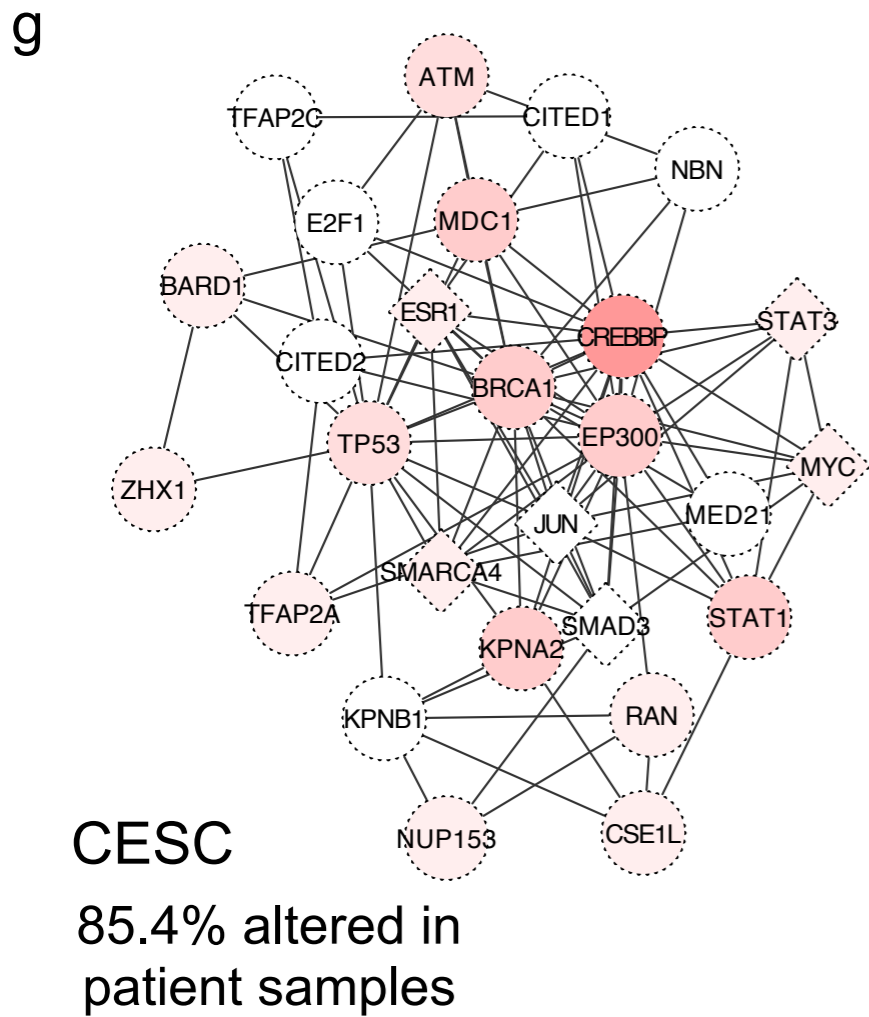
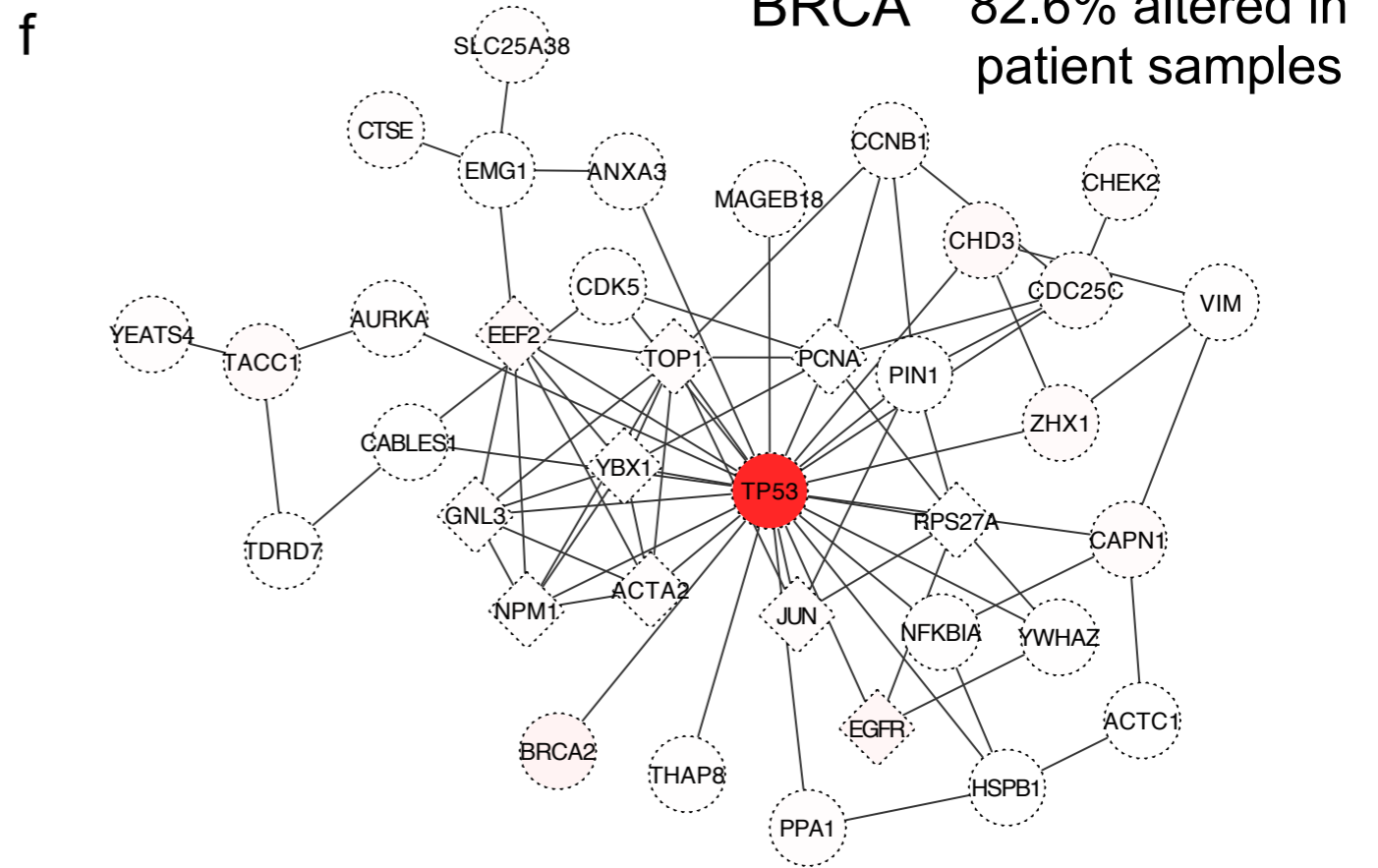
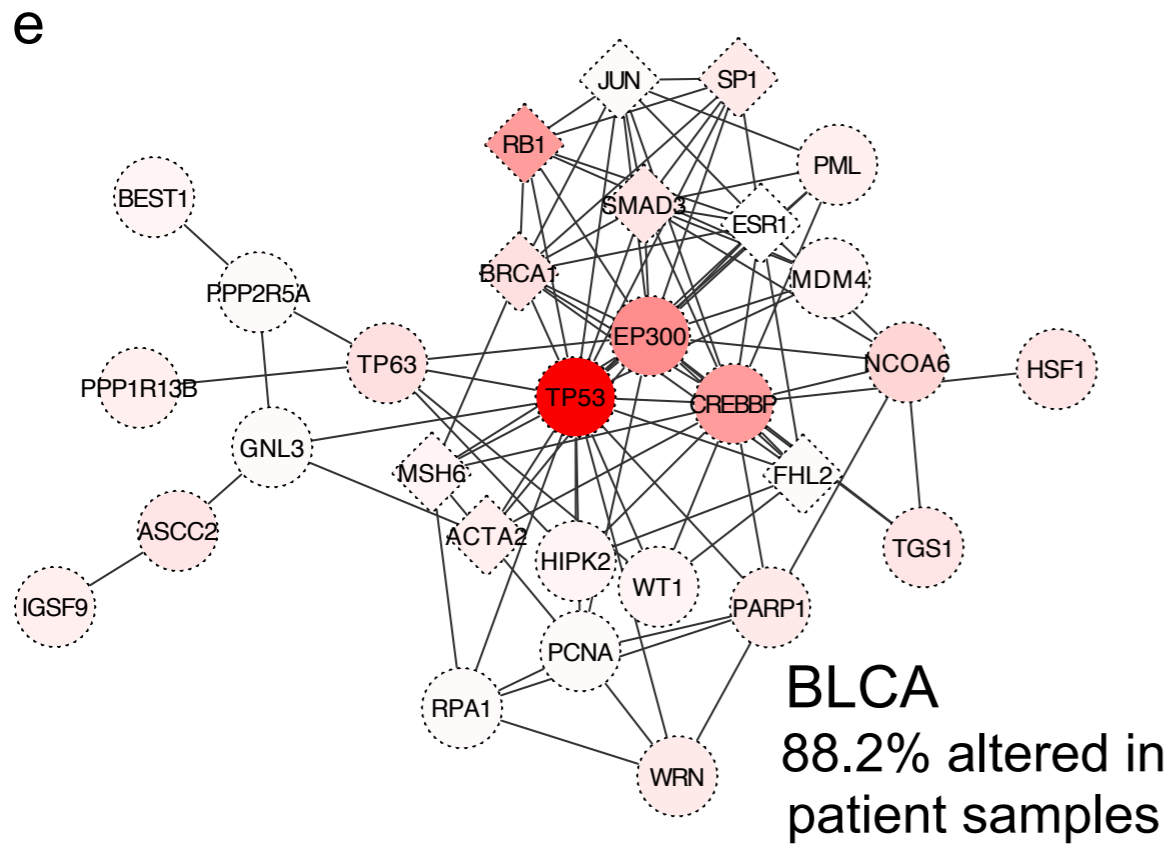
**SKCM**  
86.2% altered in patient samples

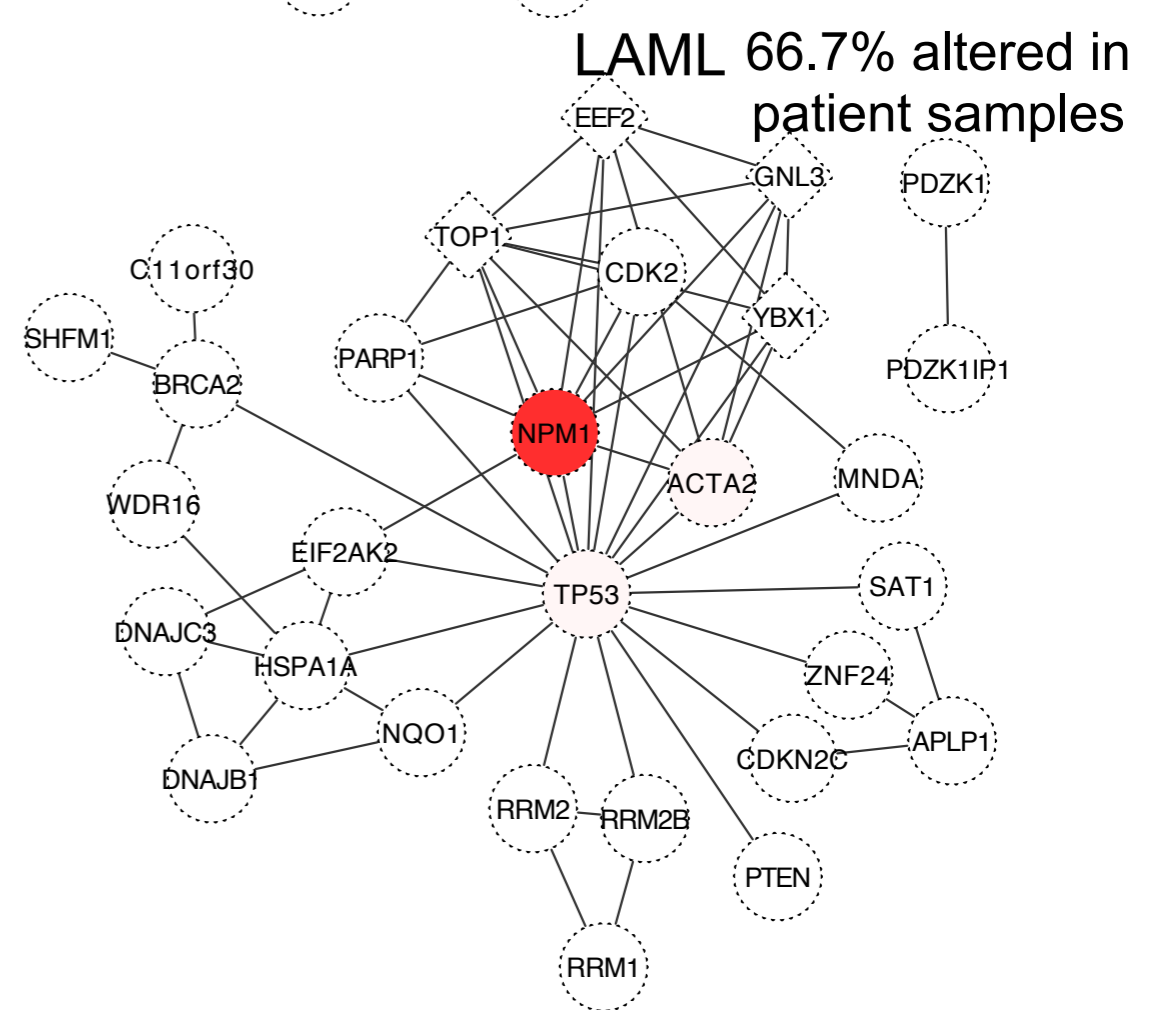
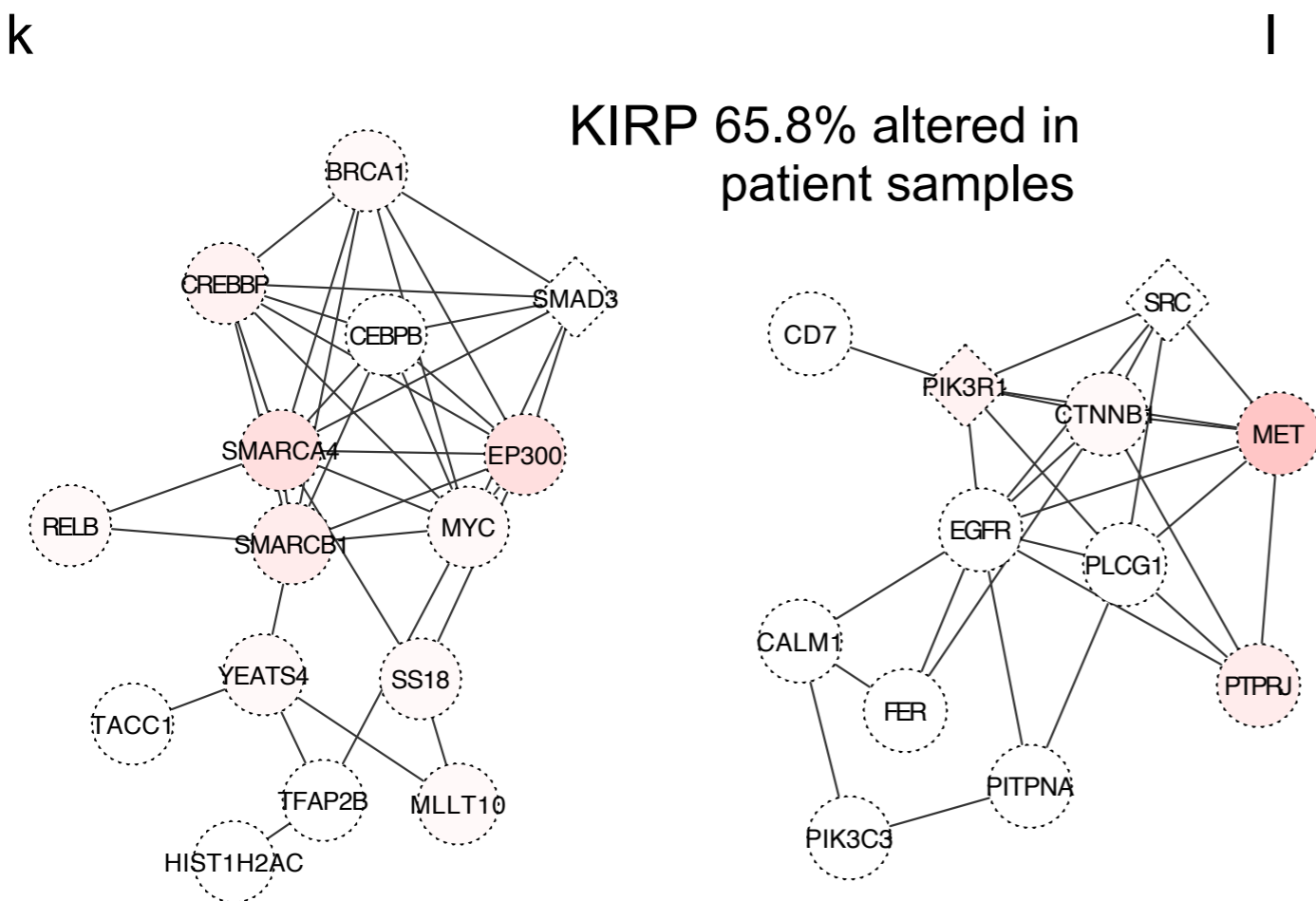
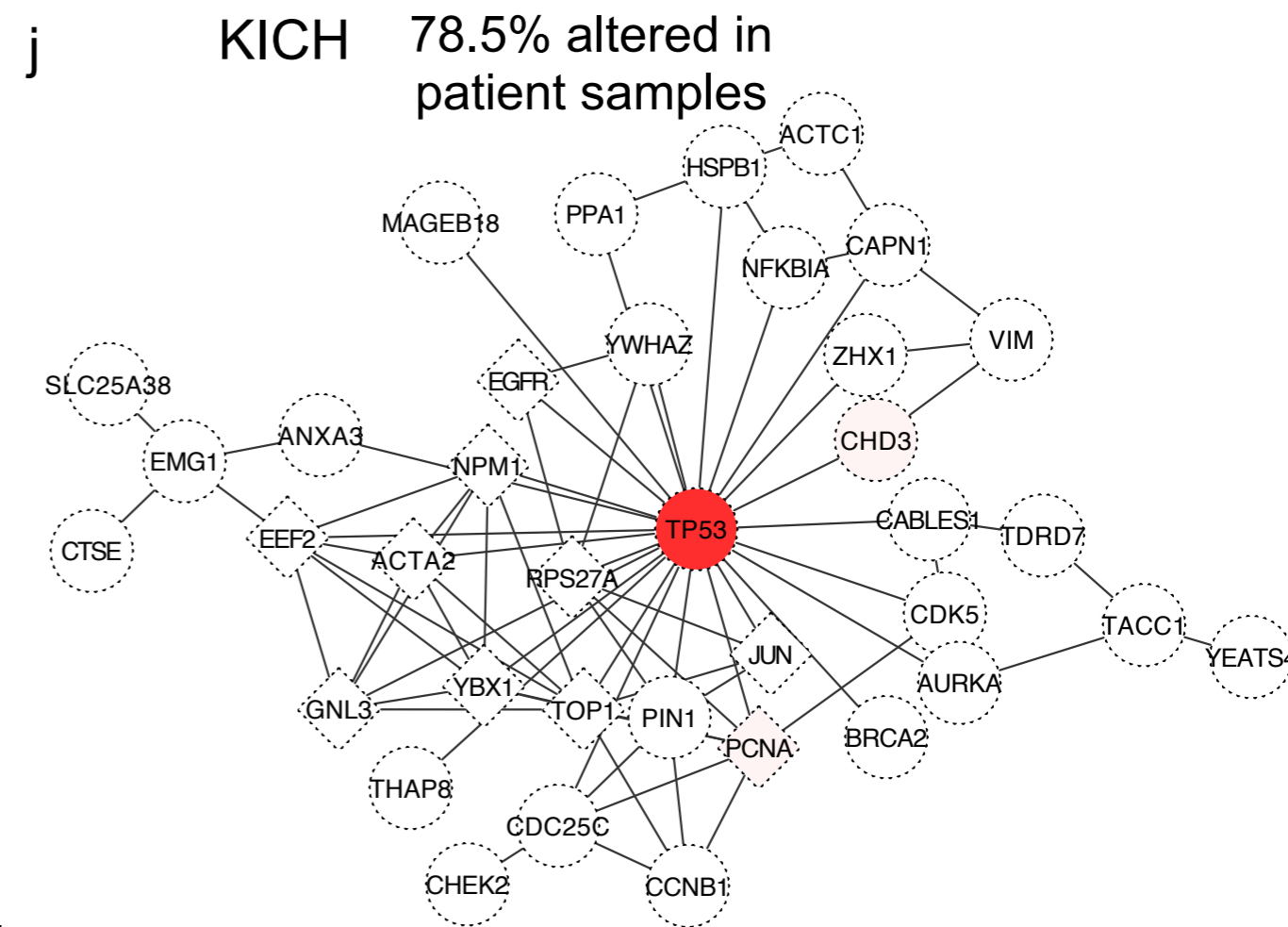
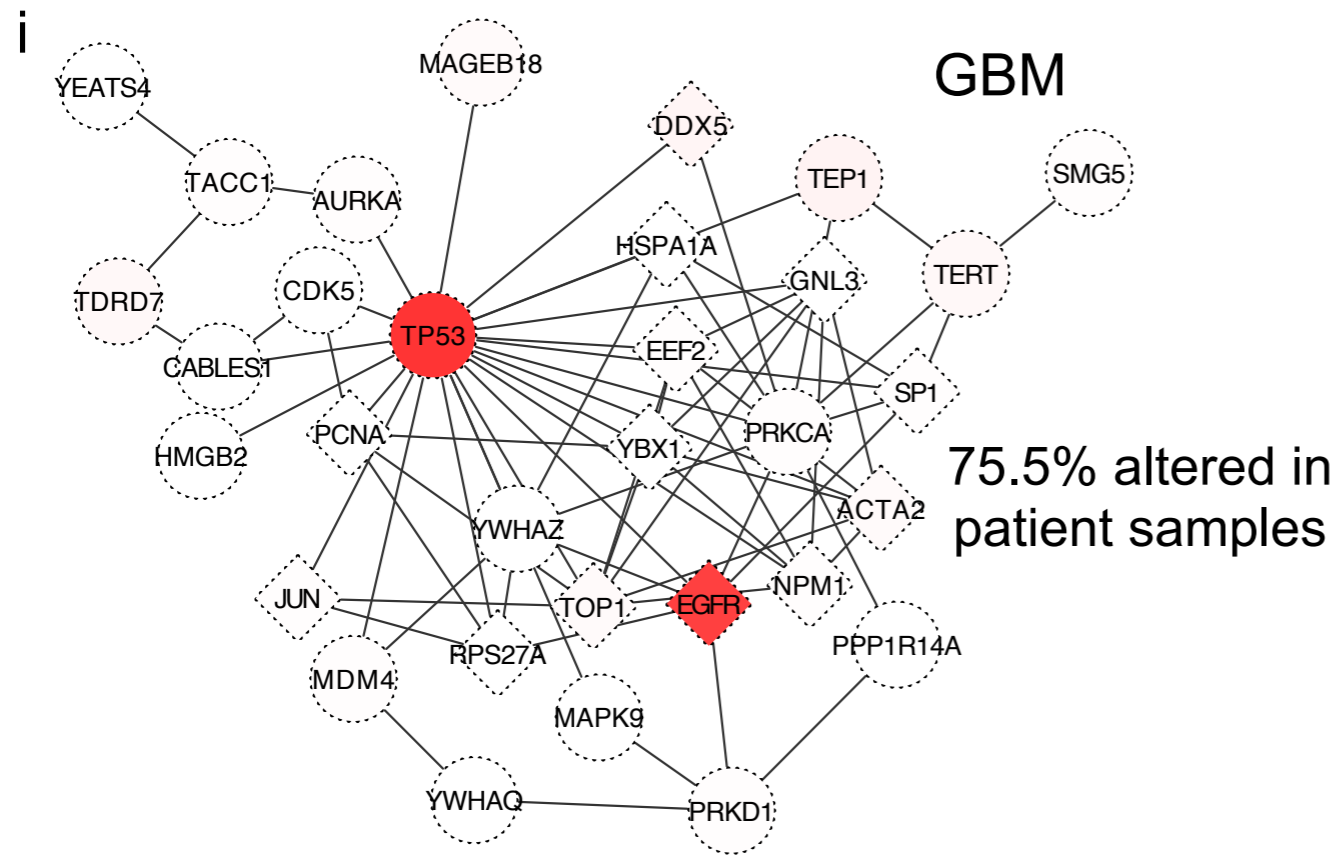
**S**

**THCA**  
65.3% altered in patient samples





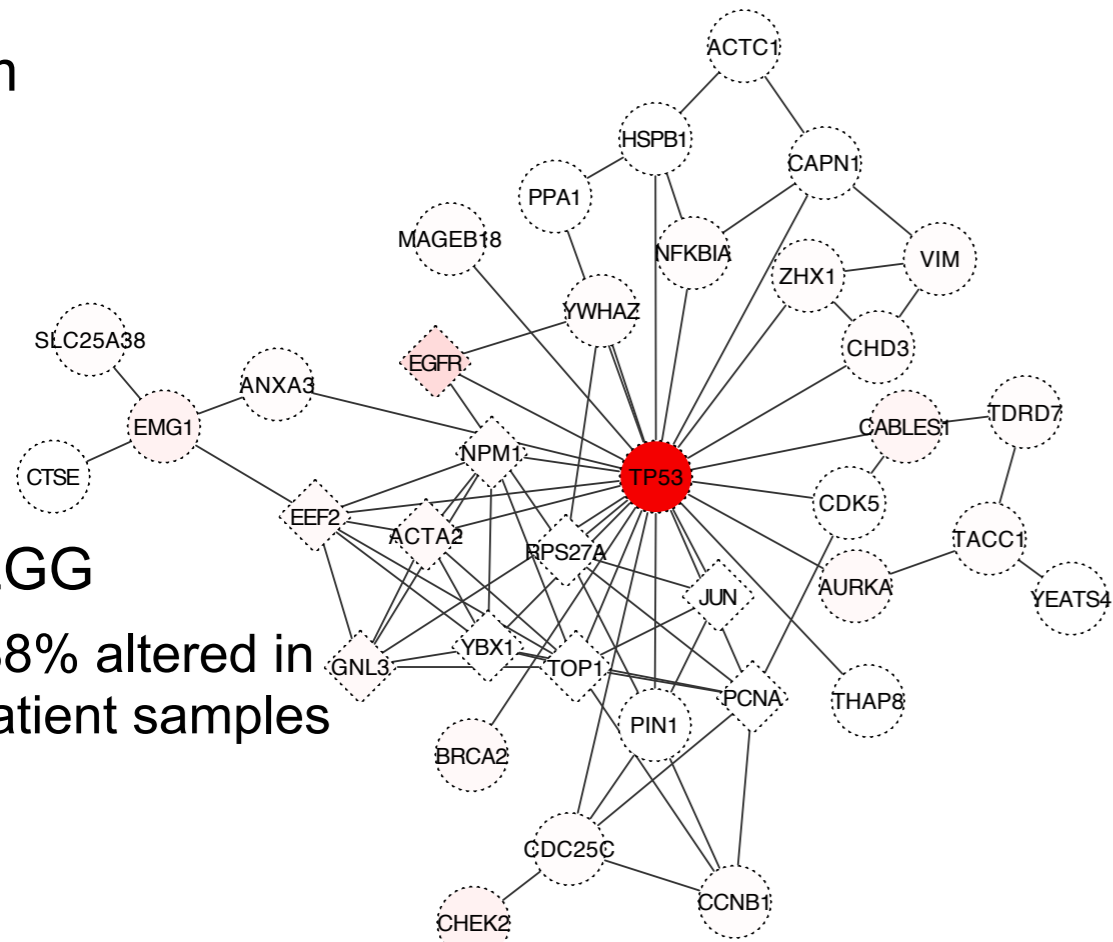




m

LGG

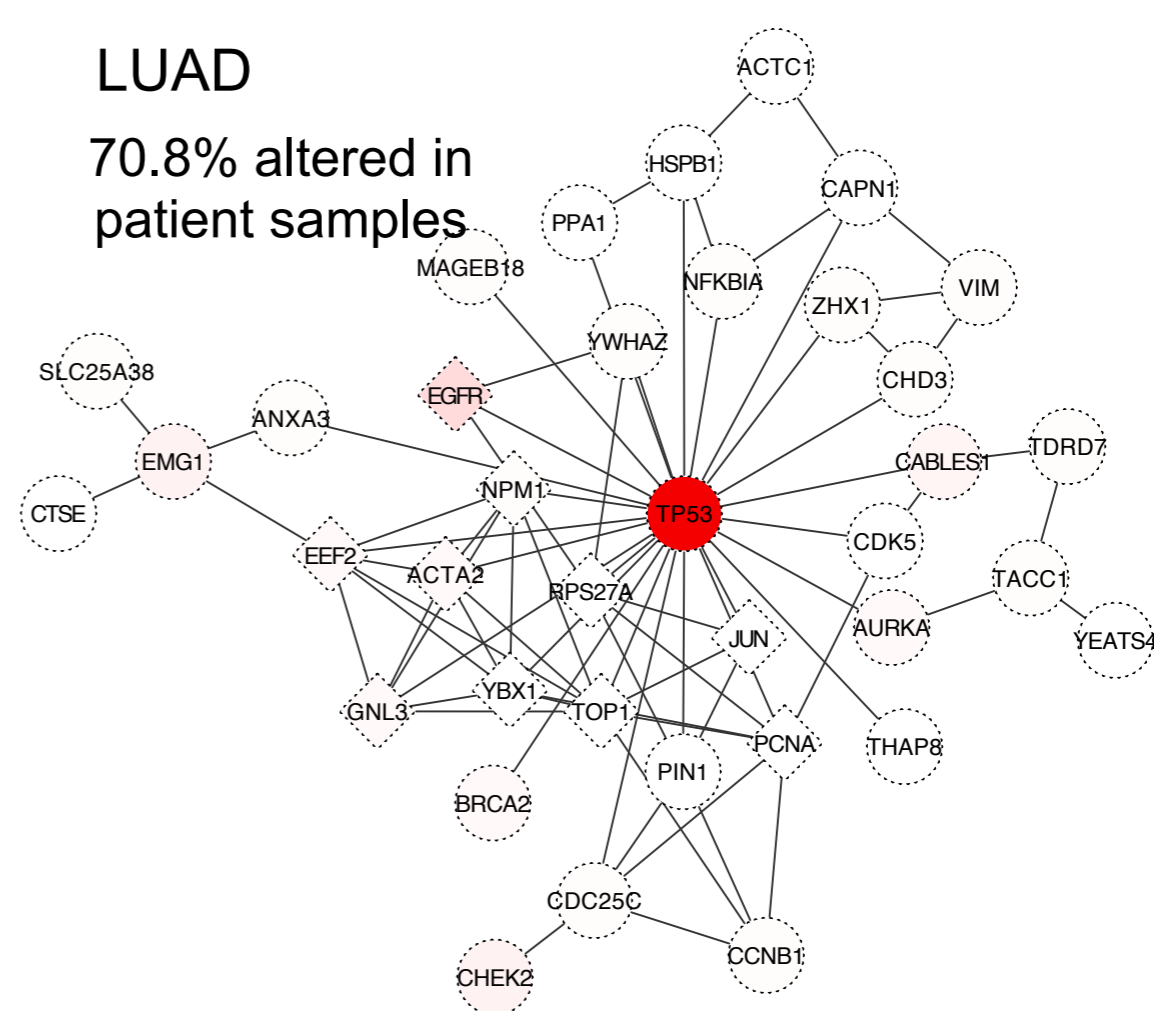
88% altered in patient samples



n

LUAD

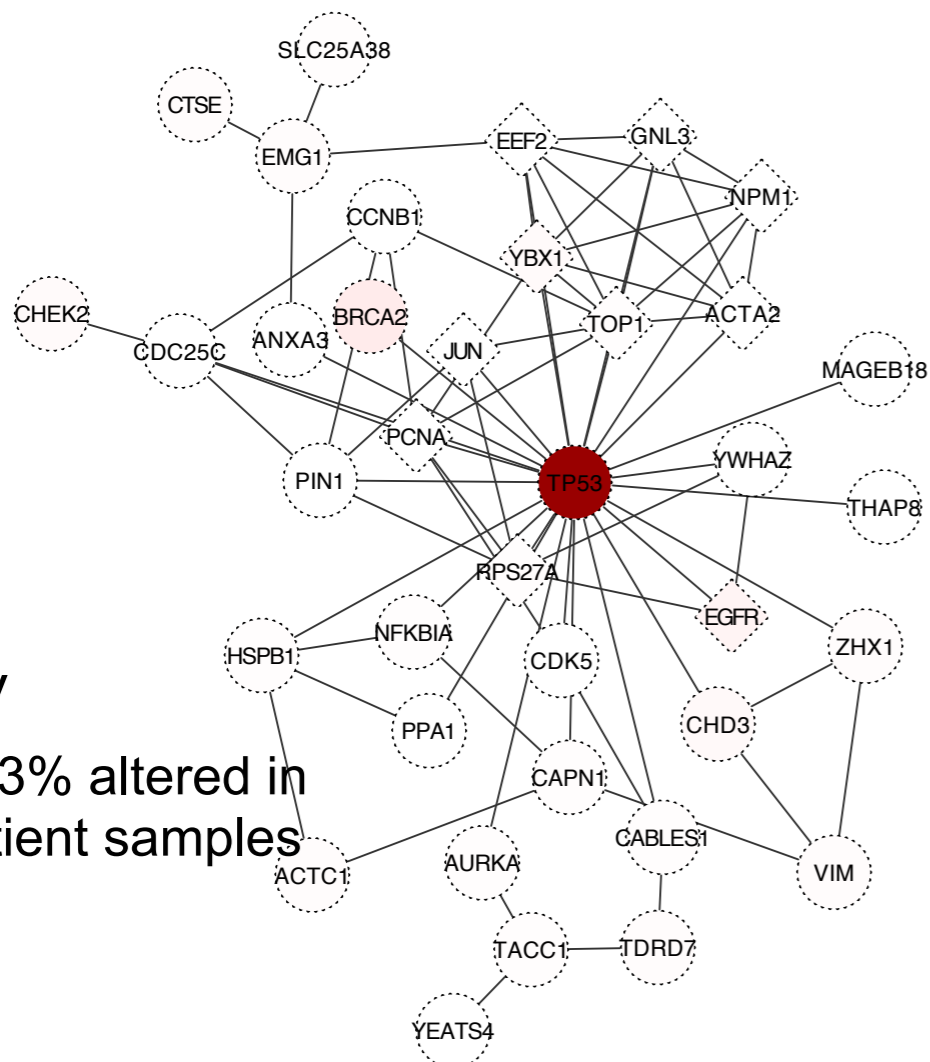
70.8% altered in patient samples



o

OV

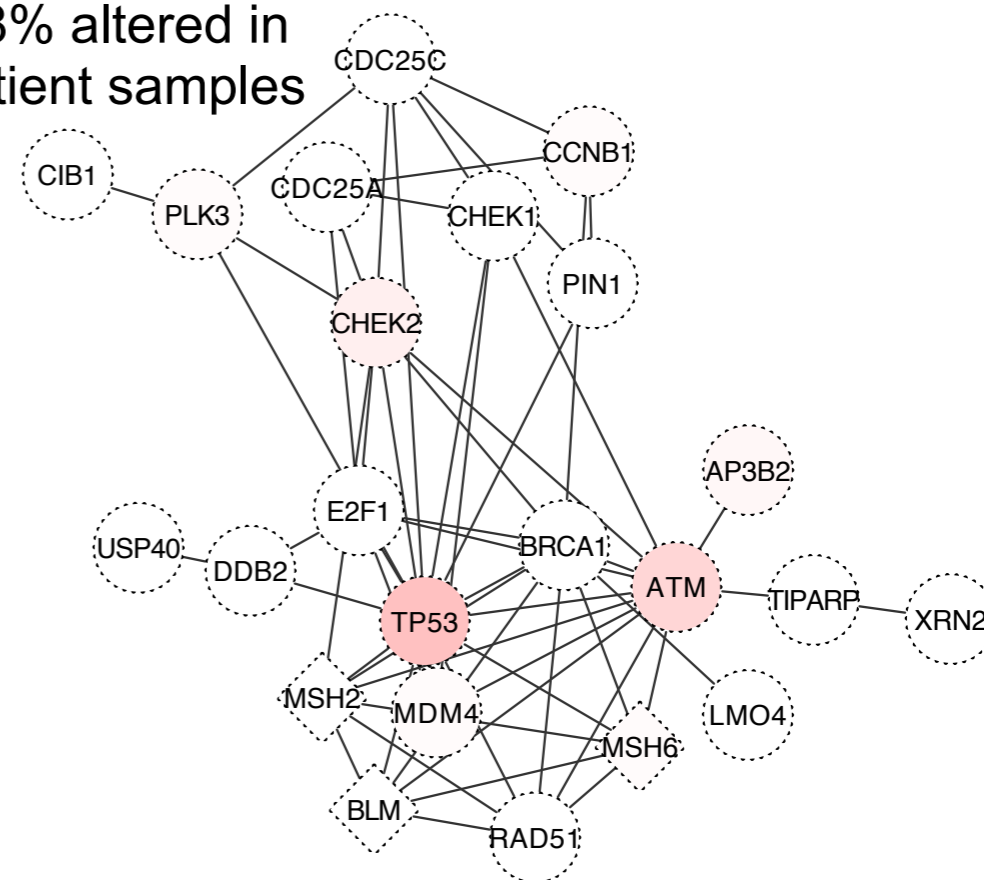
80.3% altered in patient samples

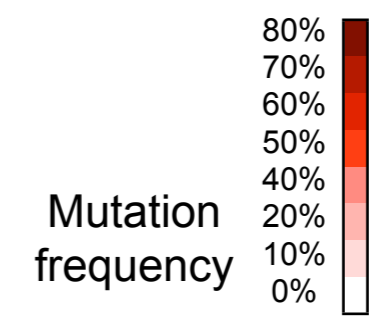
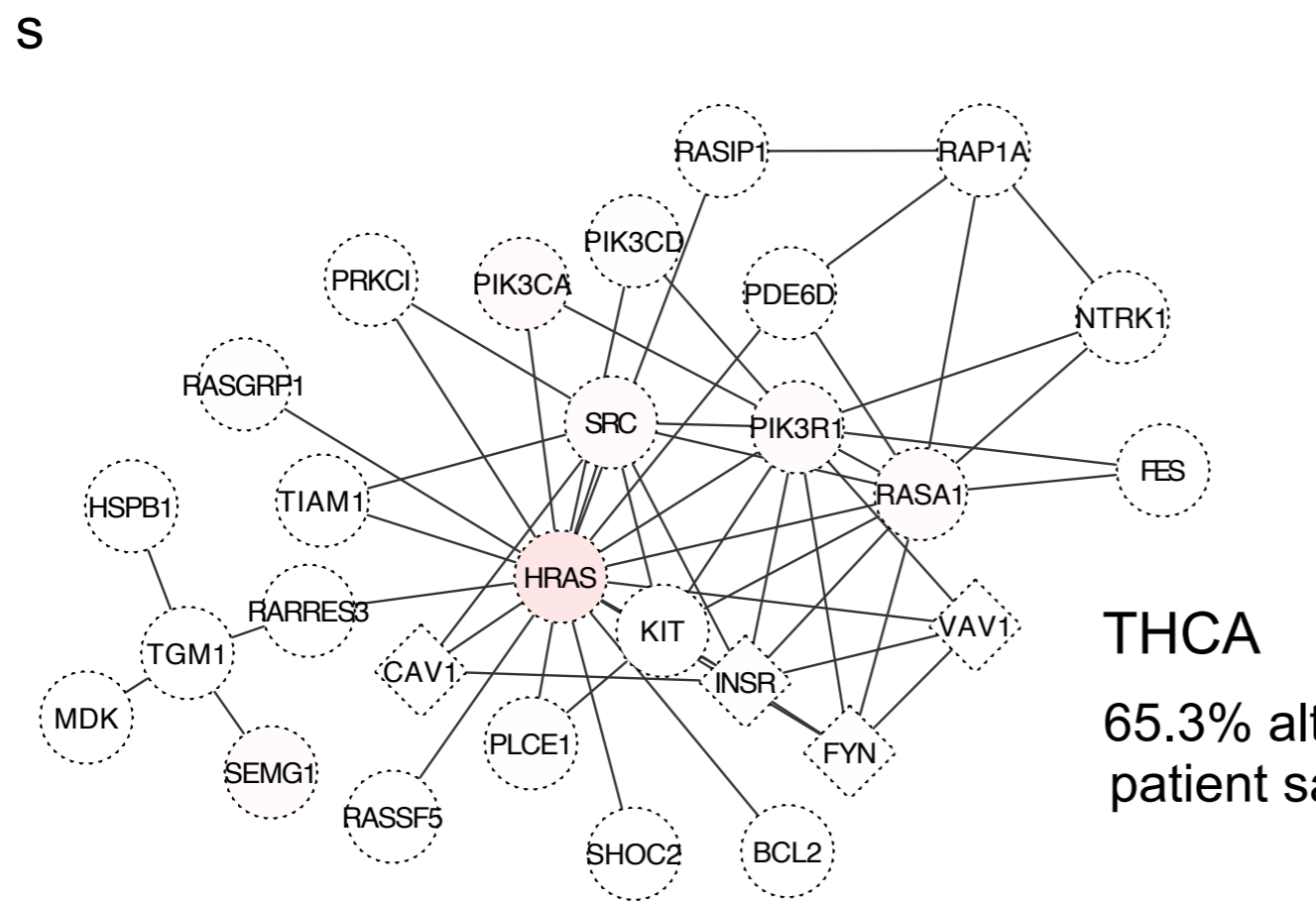
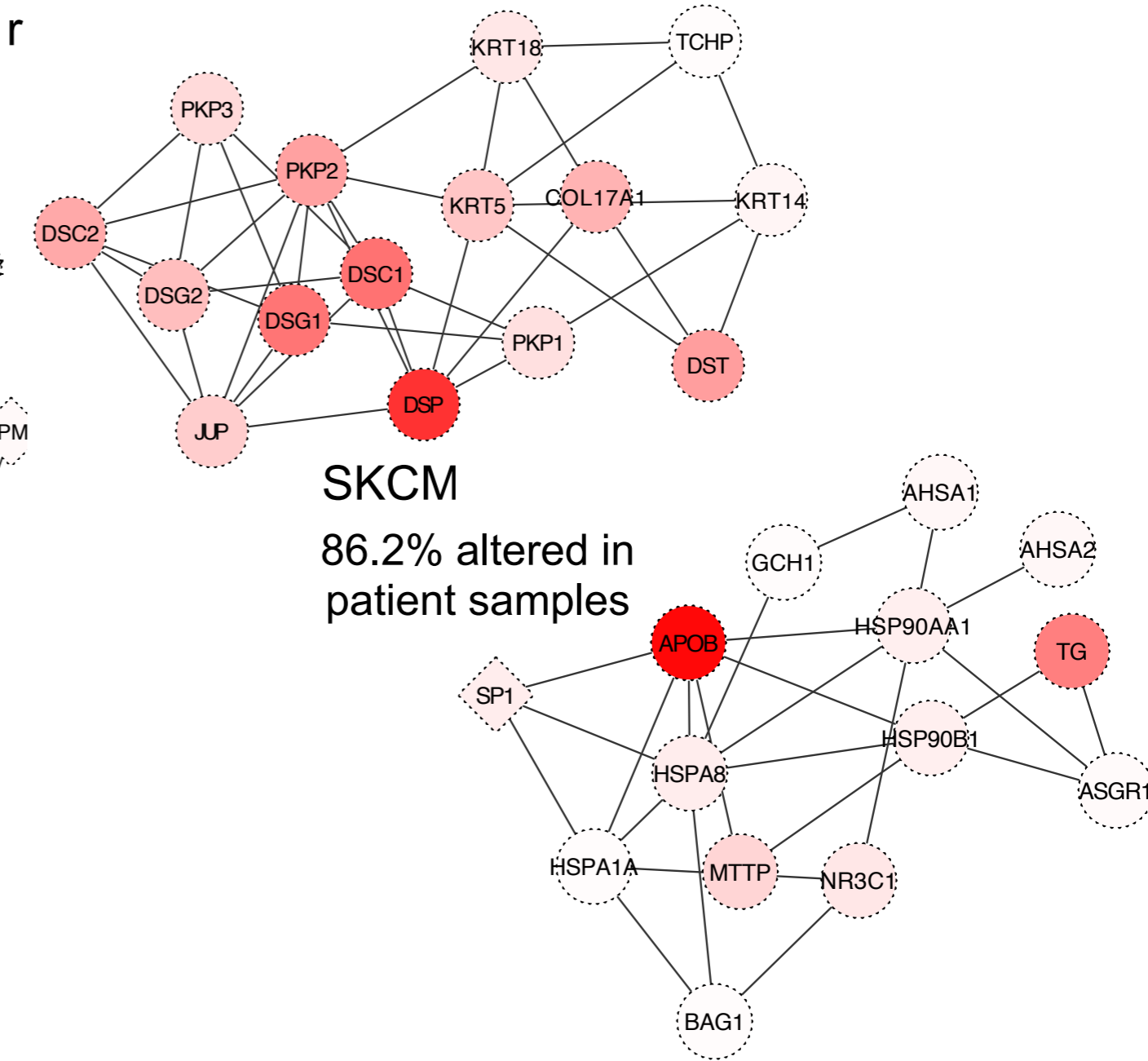
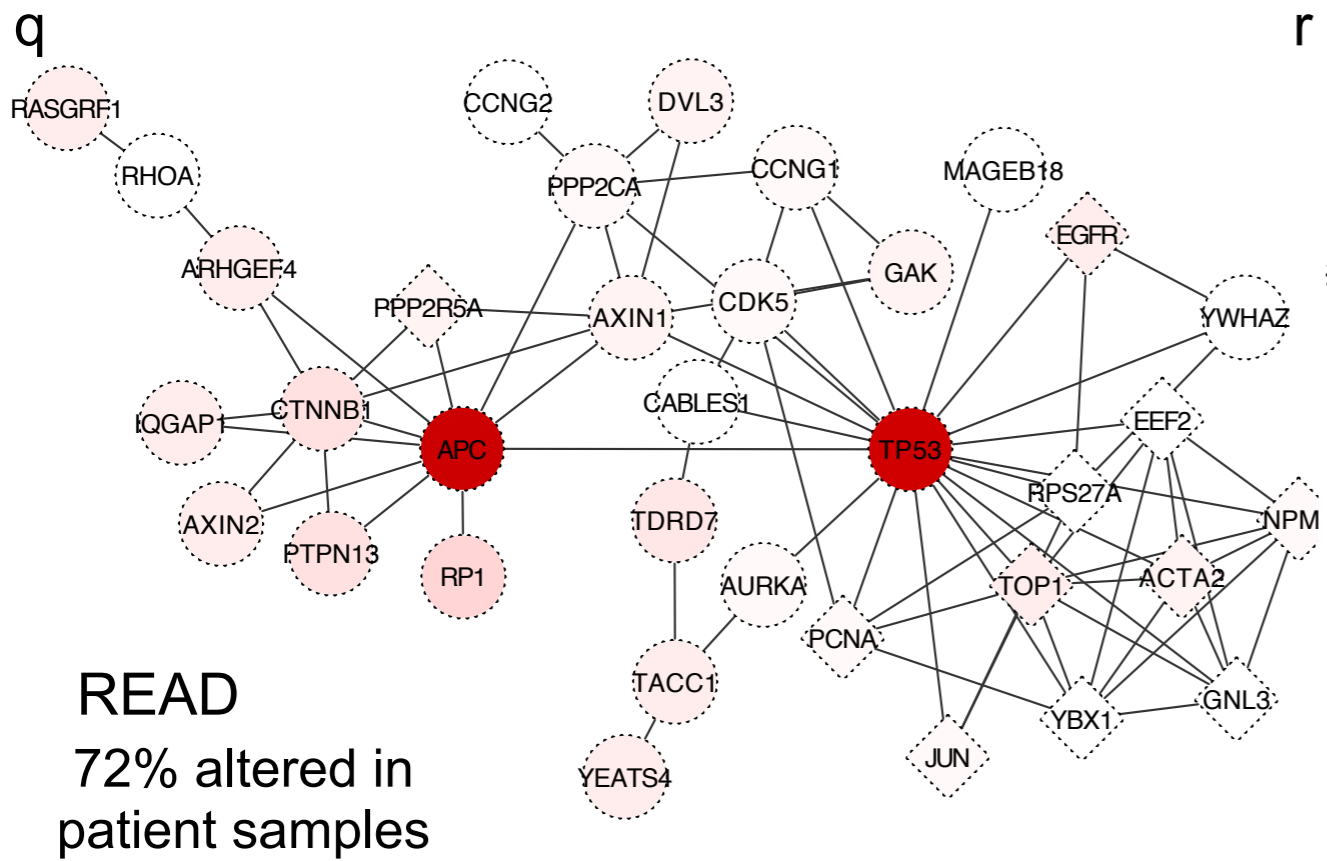


p

PRAD

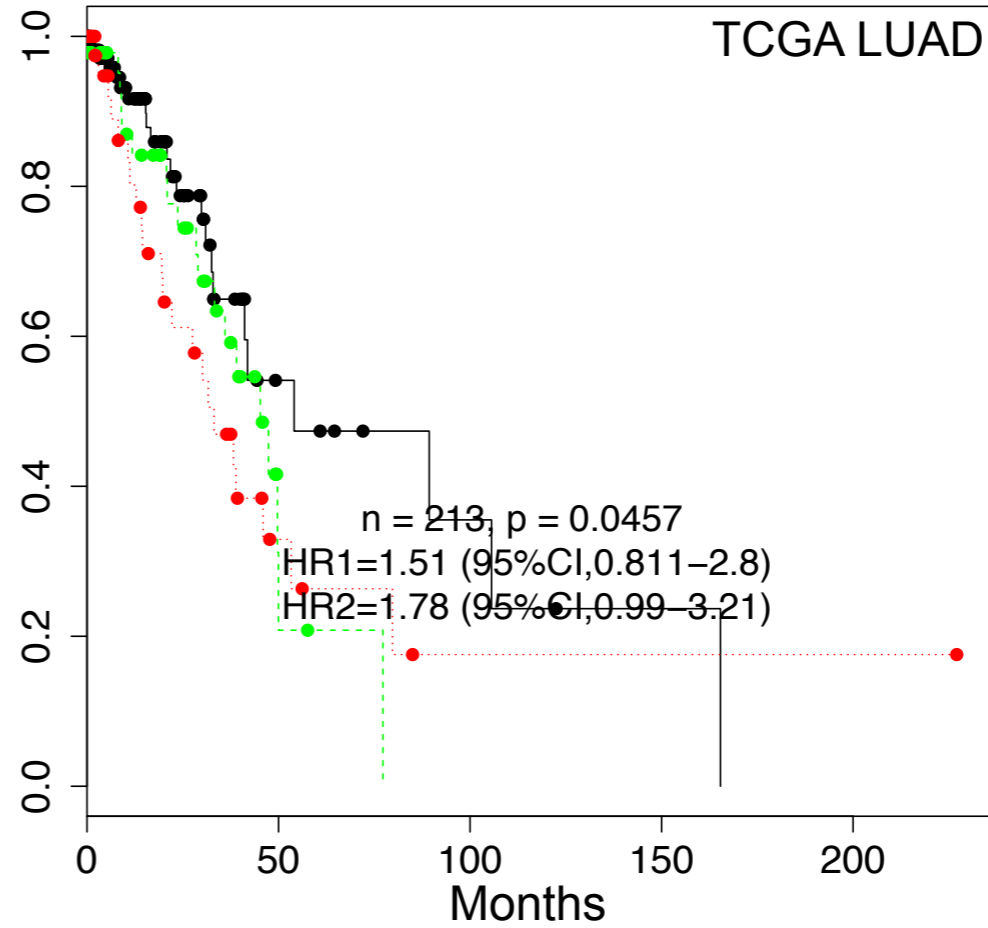
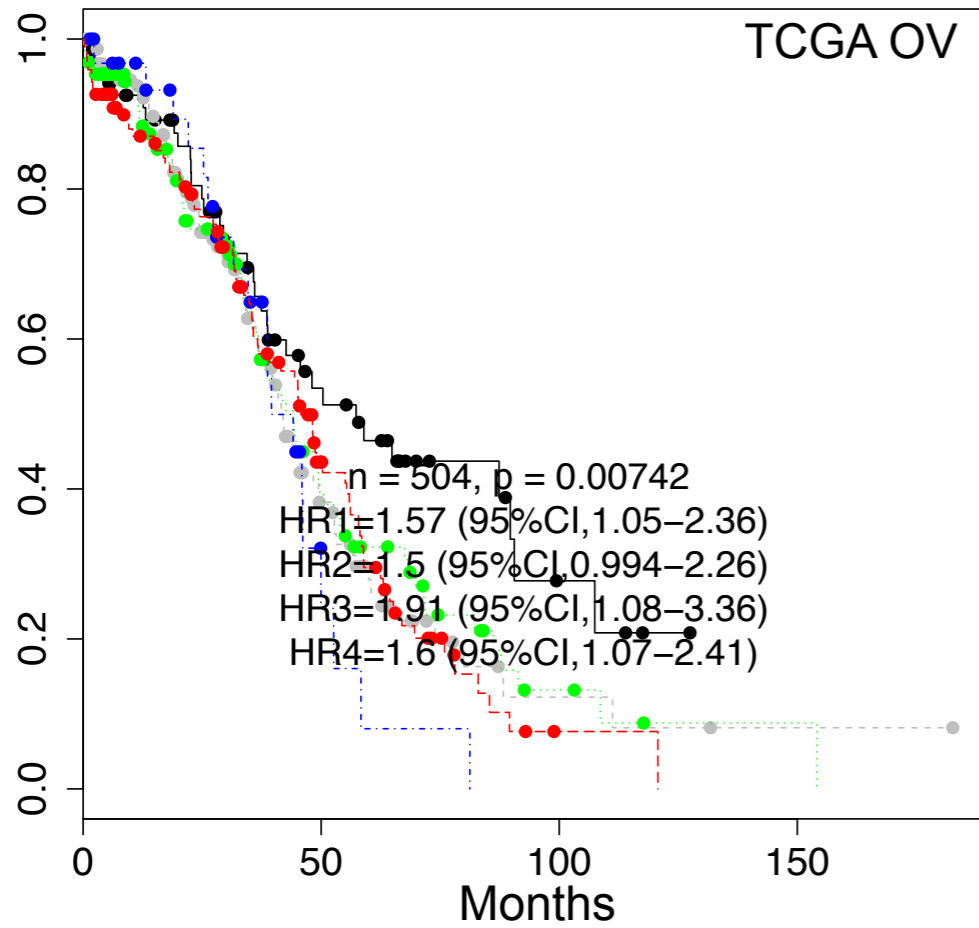
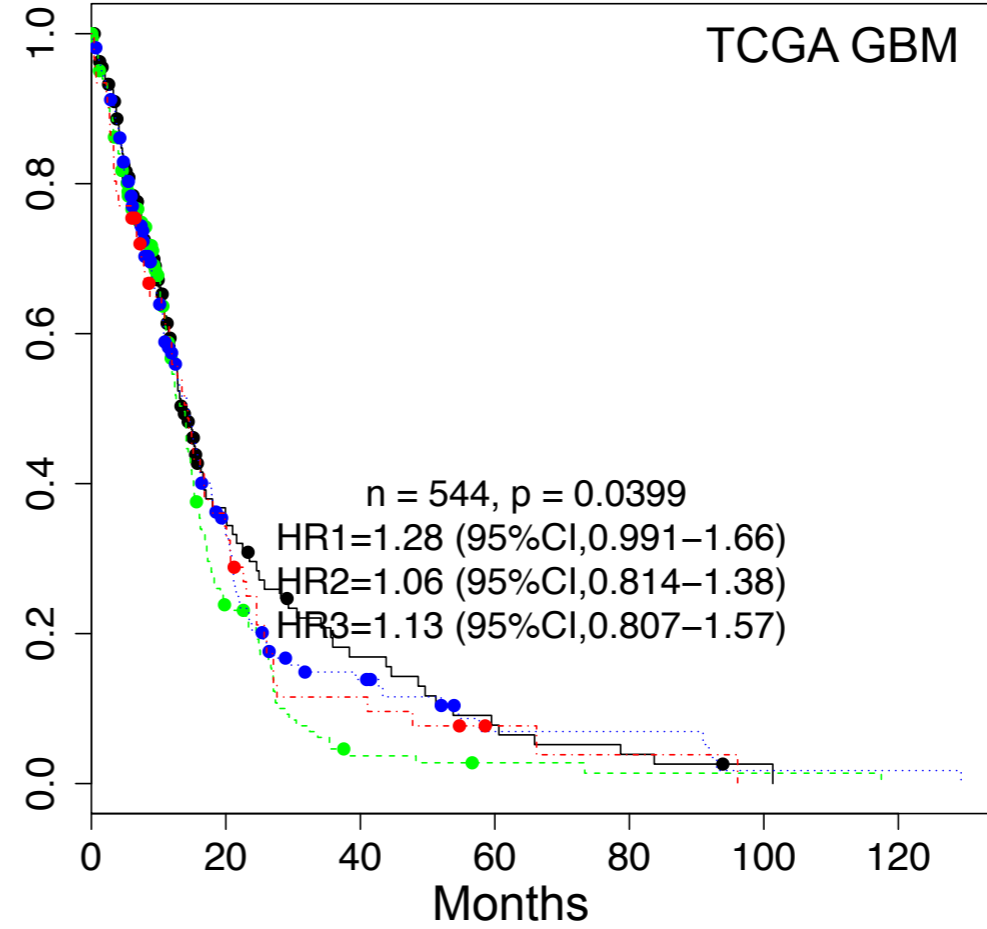
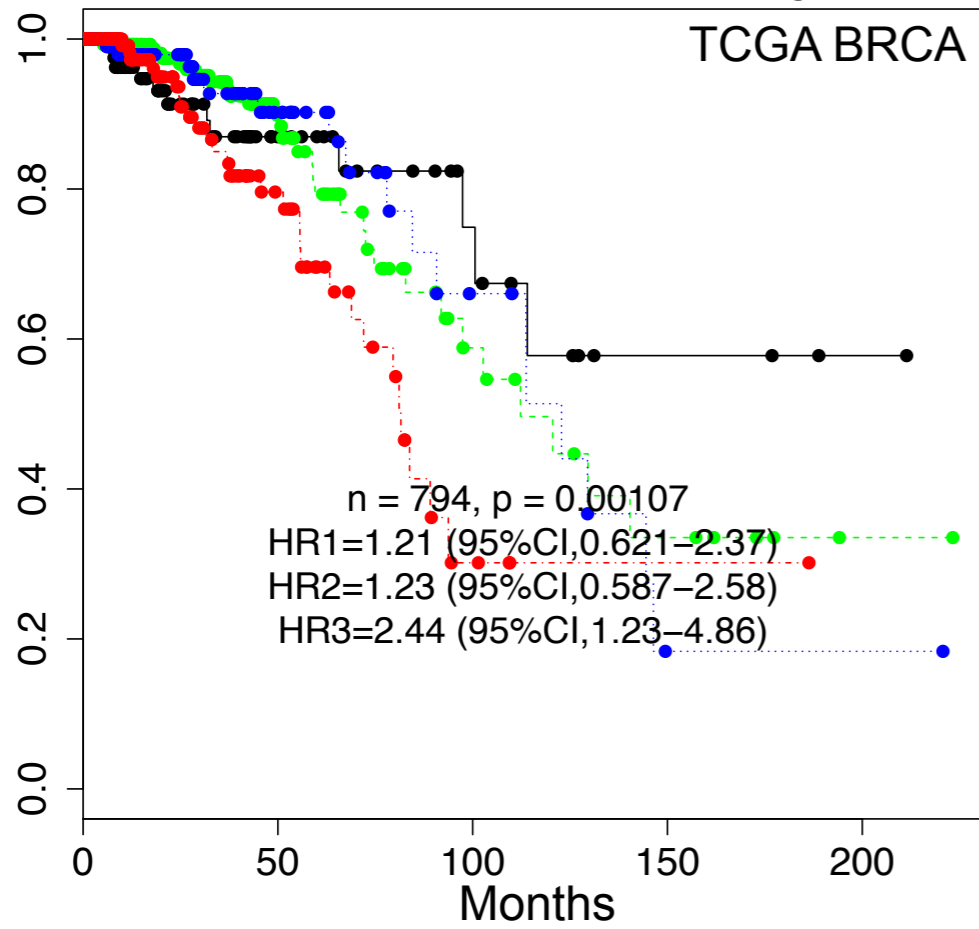
48% altered in patient samples







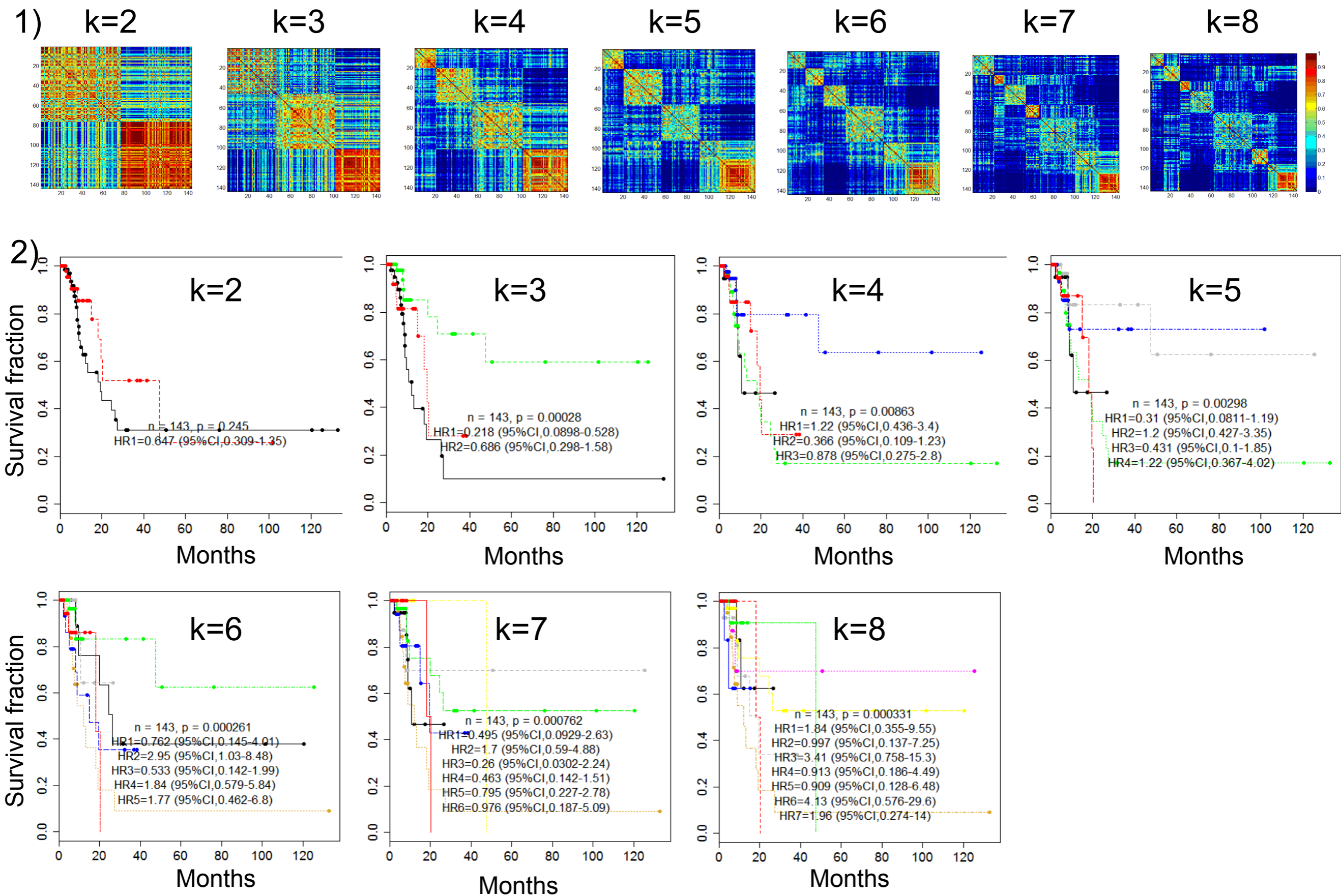
# Supplementary Figure 2



# Supplementary Figure 3

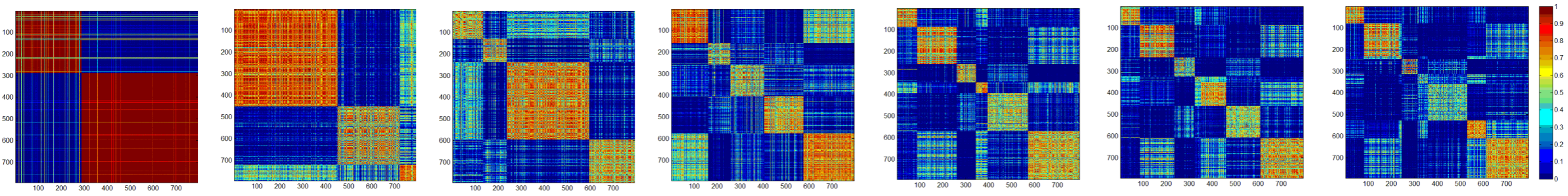
## KM plots

# a. TCGA BLCA

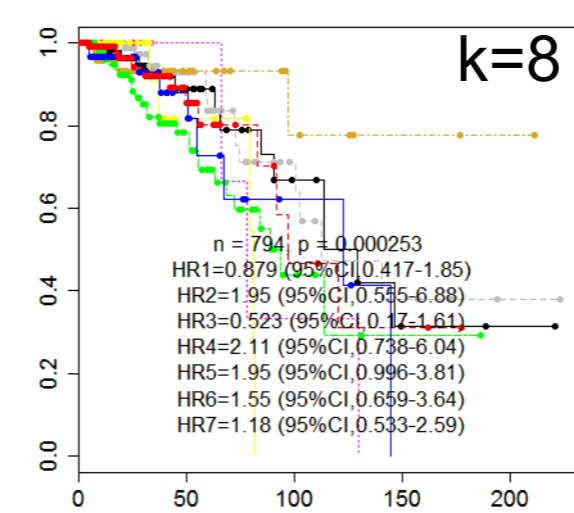
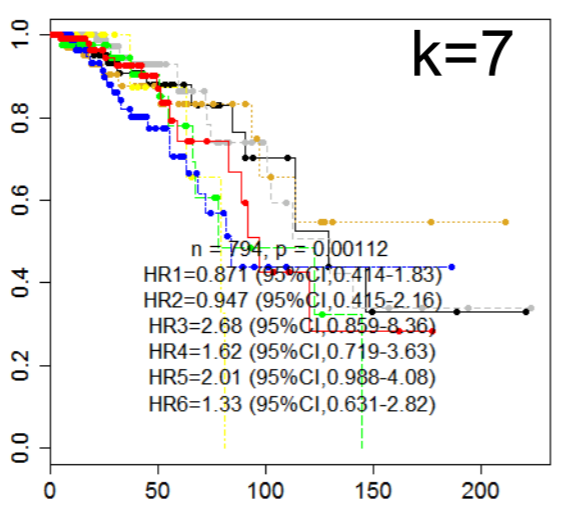
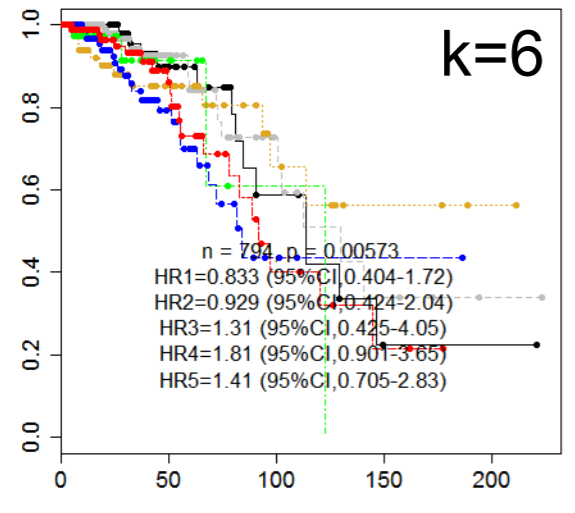
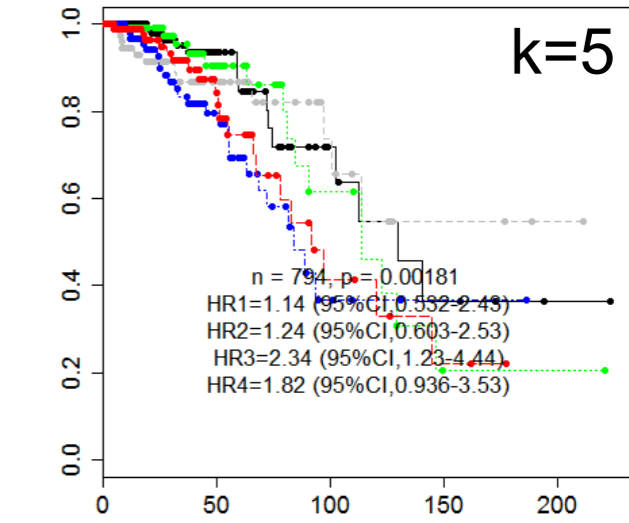
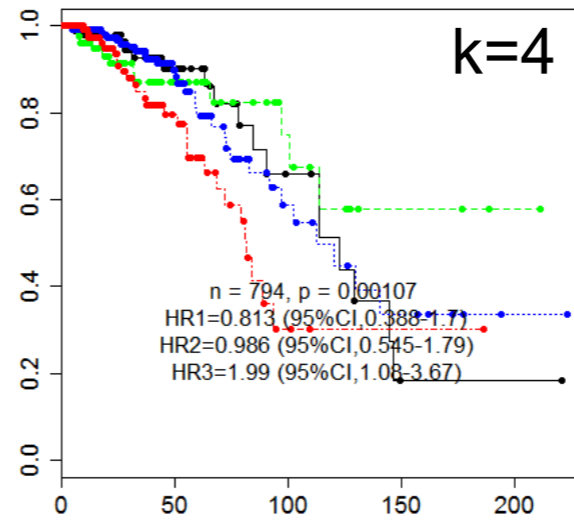
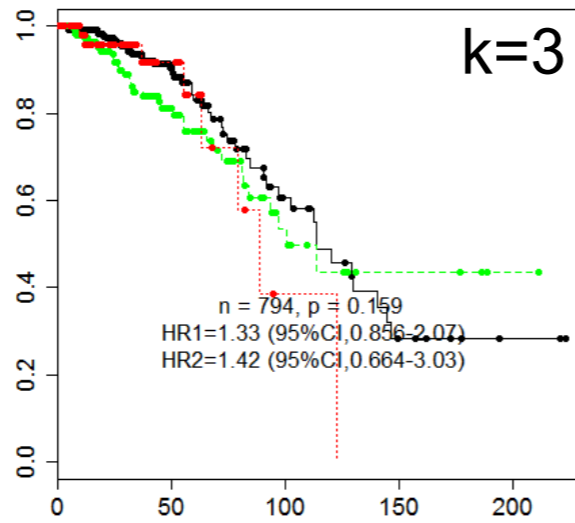
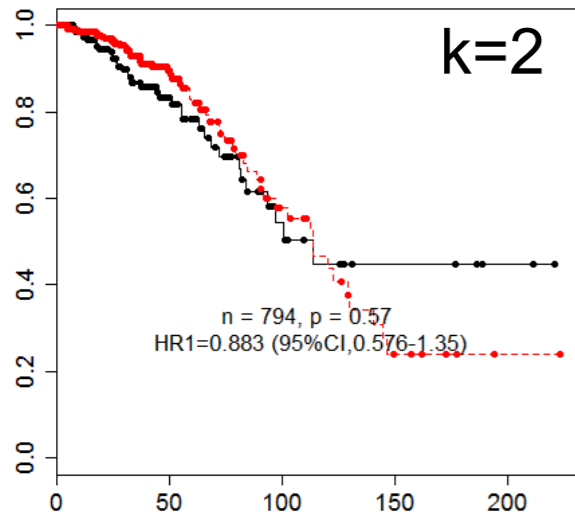


# b. TCGA BRCA

1) **k=2**      **k=3**      **k=4**      **k=5**      **k=6**      **k=7**      **k=8**



2)



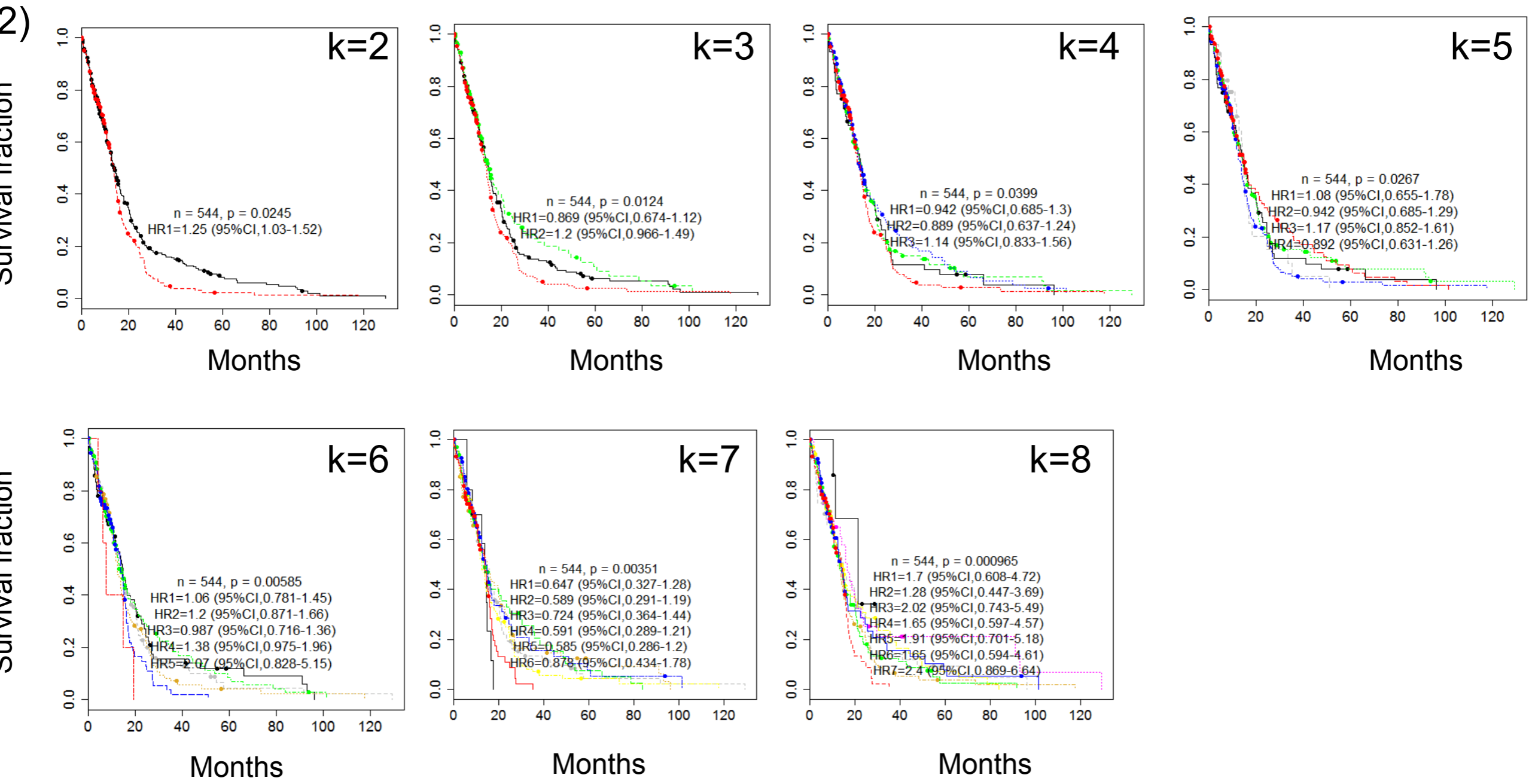
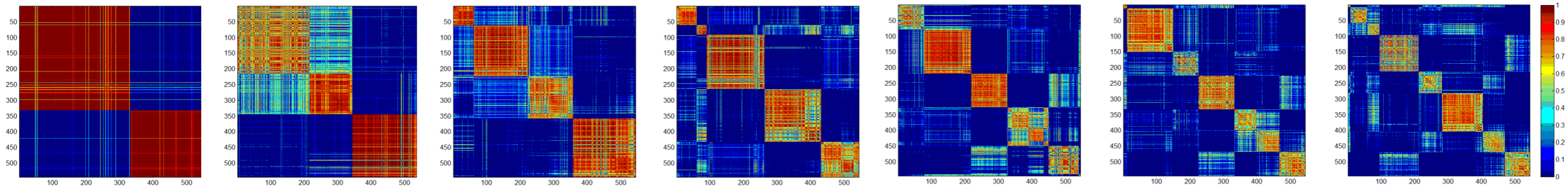
Months

Months

Months

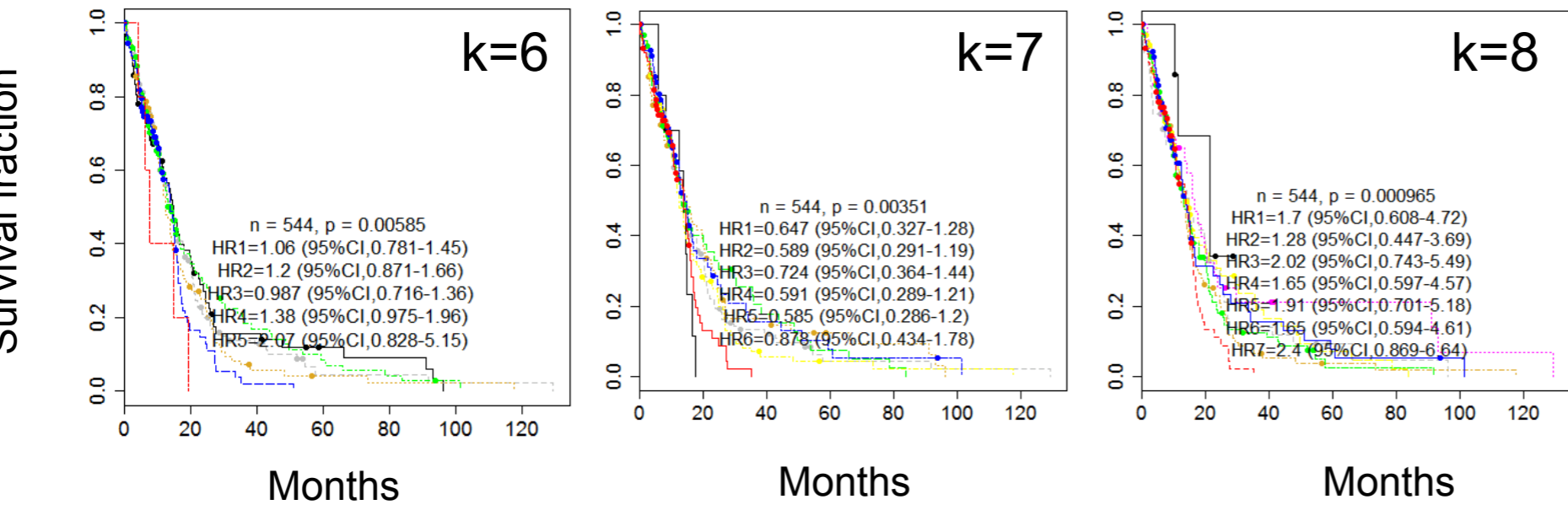
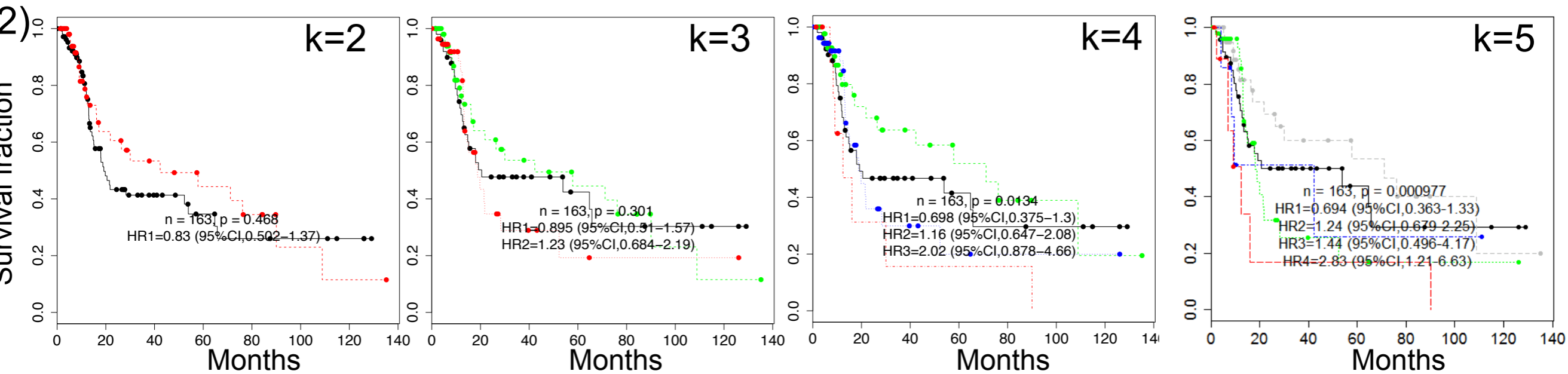
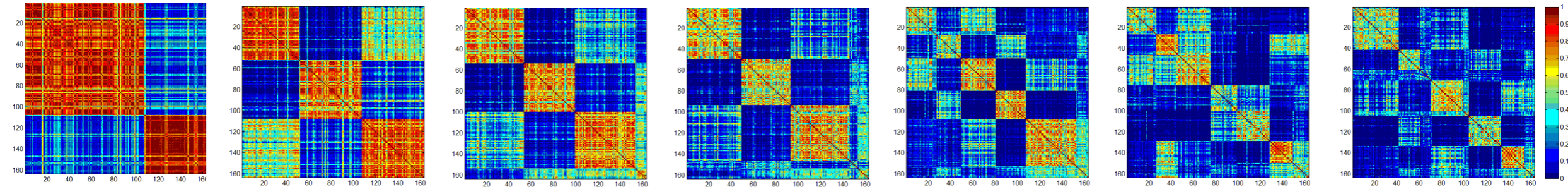
# c. TCGA GBM

1) k=2      k=3      k=4      k=5      k=6      k=7      k=8



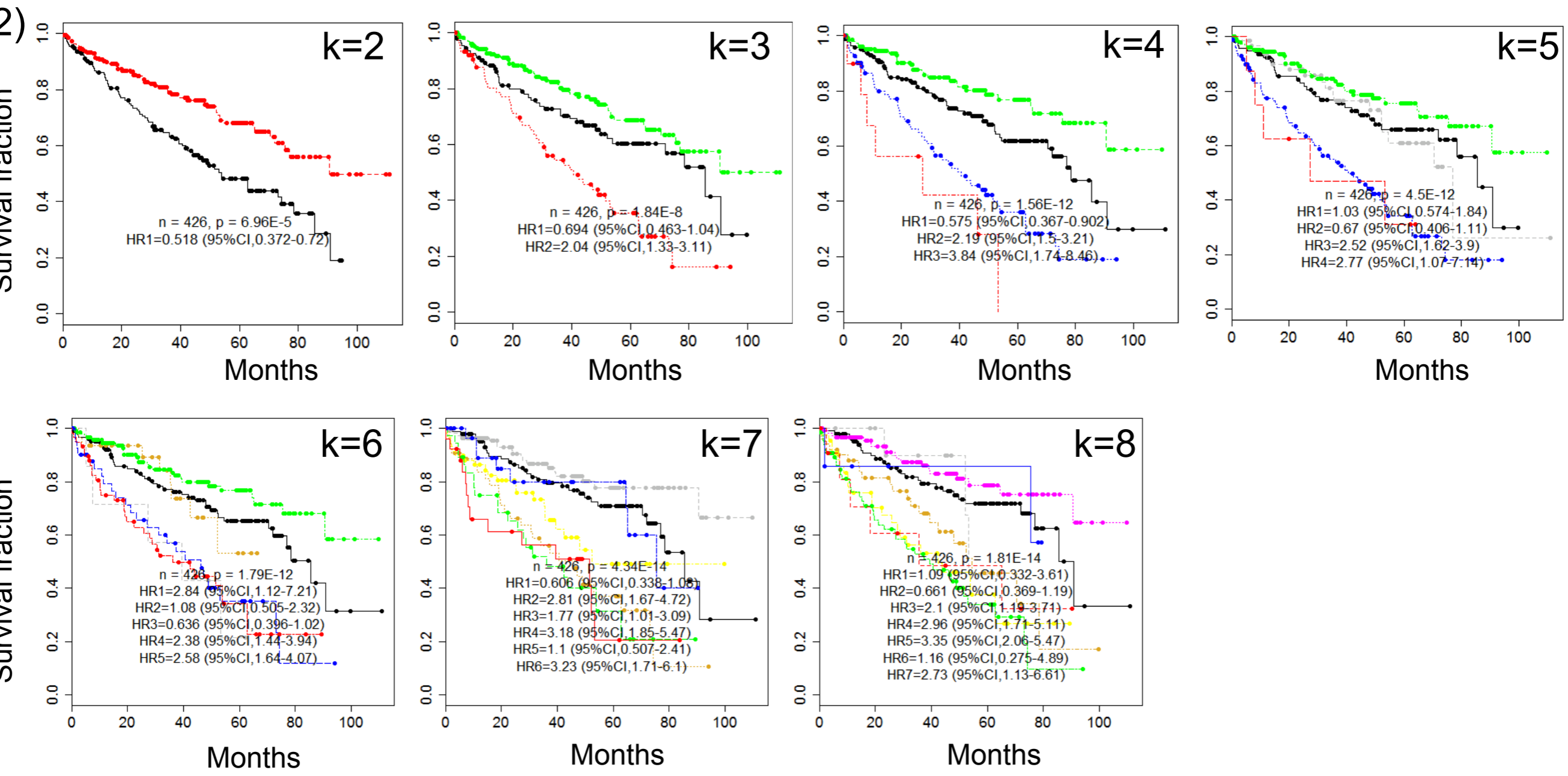
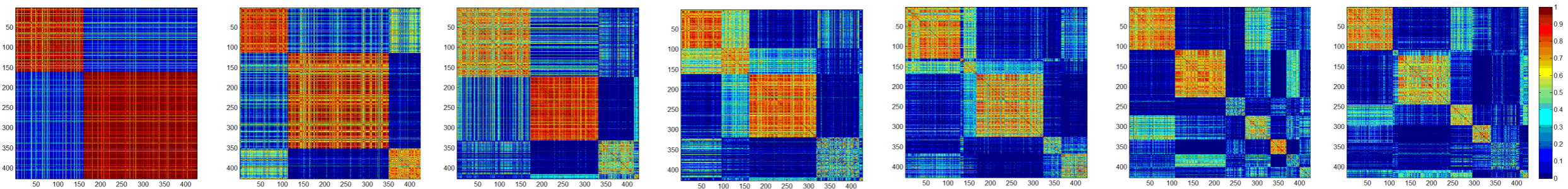
# d.TCGA HNSC

1) k=2      k=3      k=4      k=5      k=6      k=7      k=8



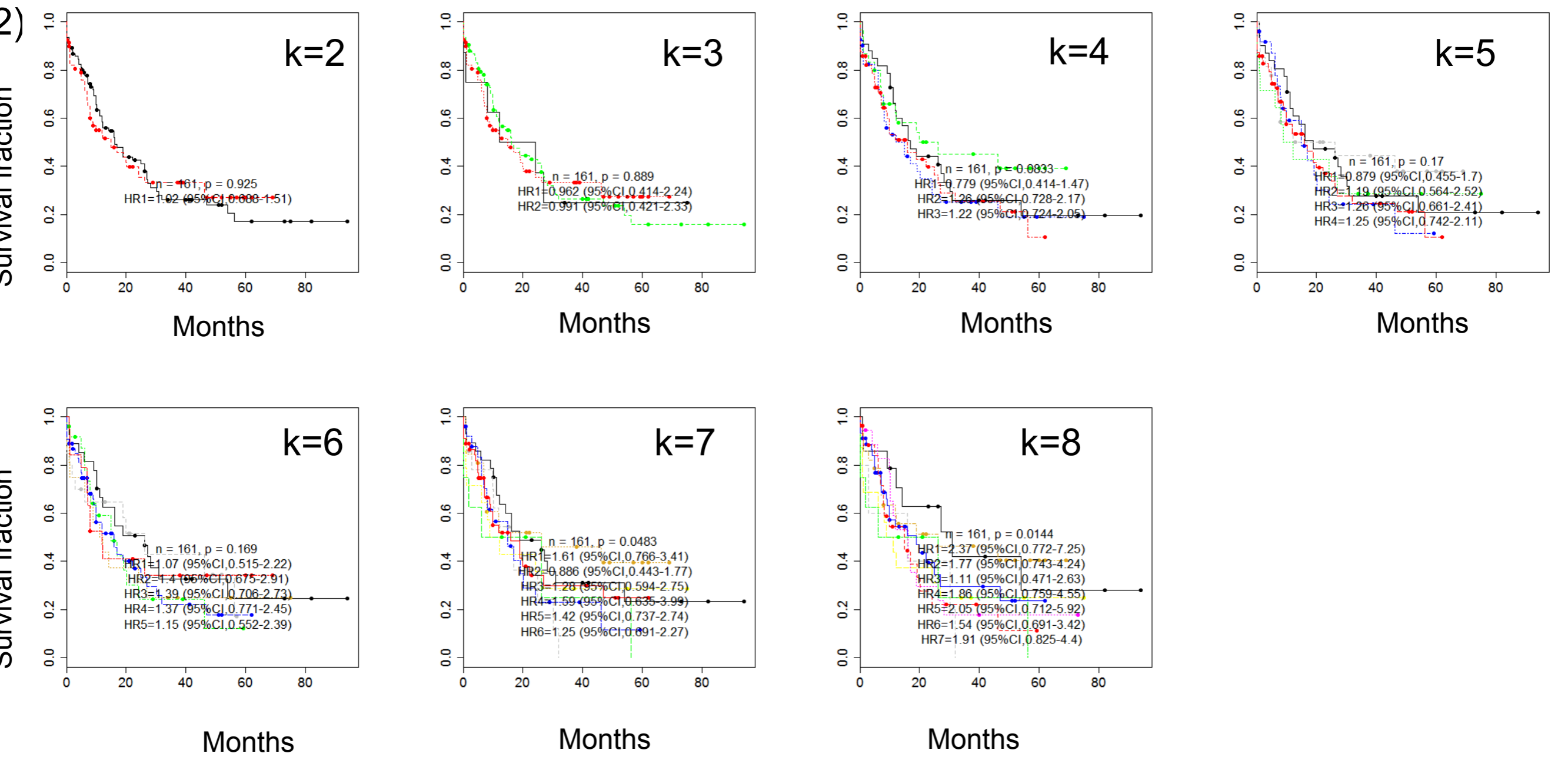
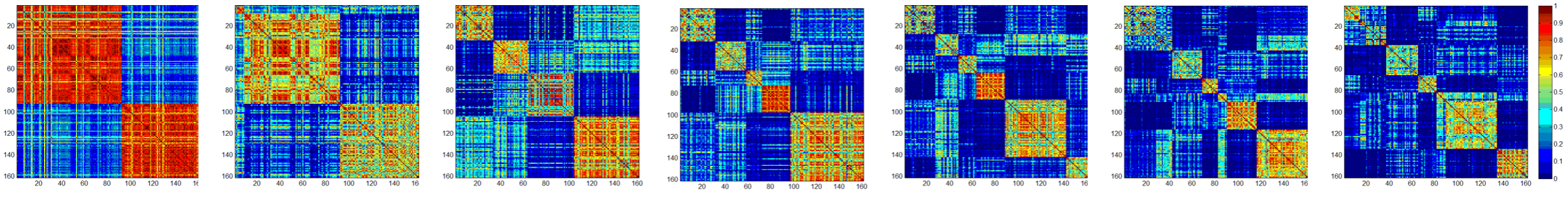
# e. TCGA KIRC

1) k=2      k=3      k=4      k=5      k=6      k=7      k=8



# f. TCGA LAML

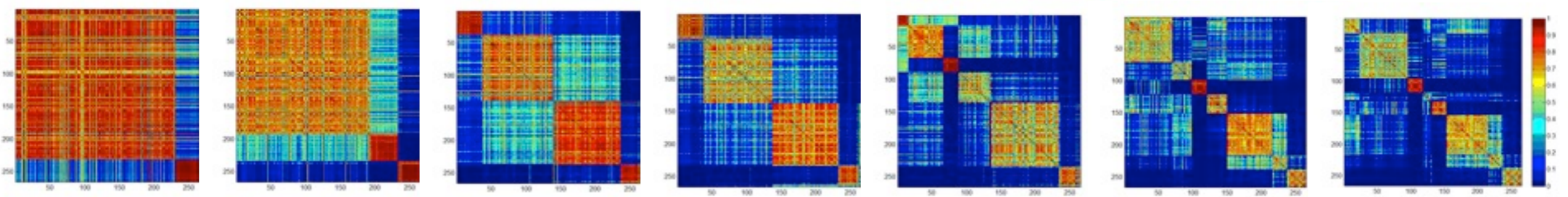
1) **k=2**      **k=3**      **k=4**      **k=5**      **k=6**      **k=7**      **k=8**





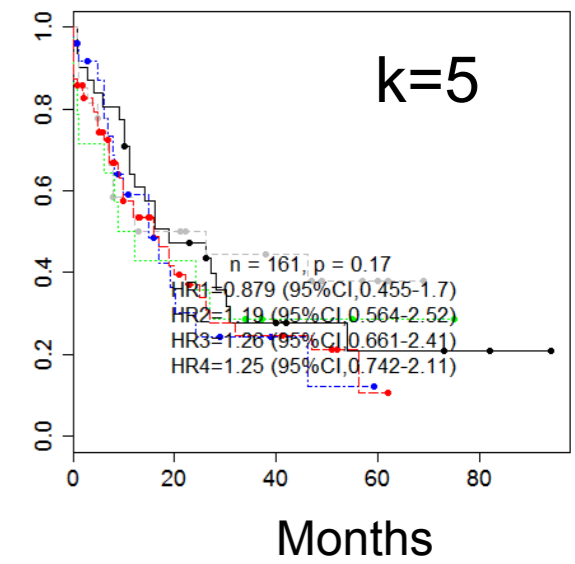
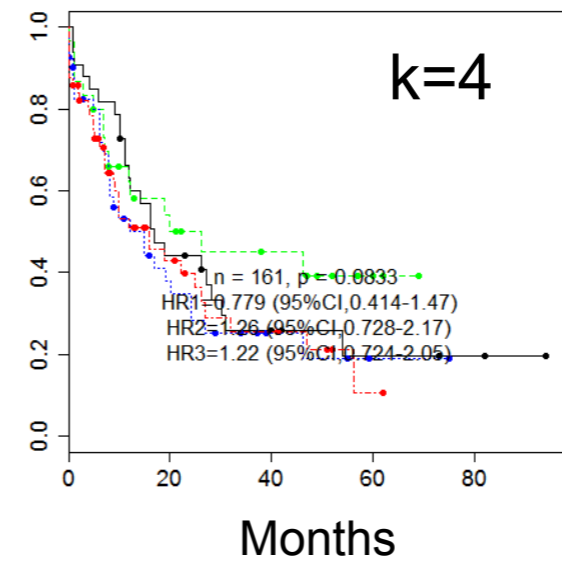
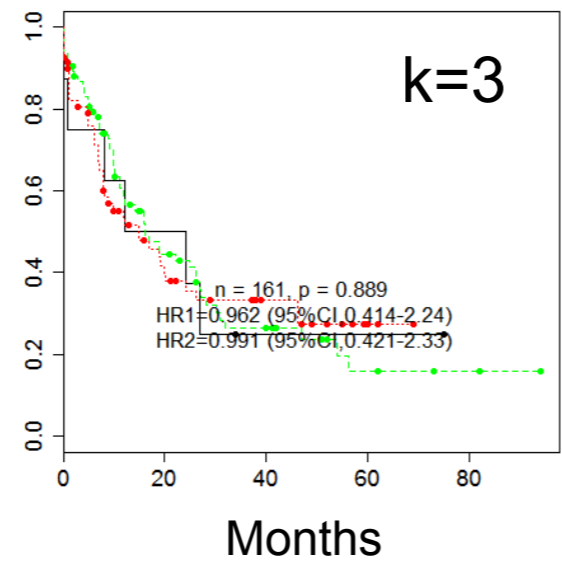
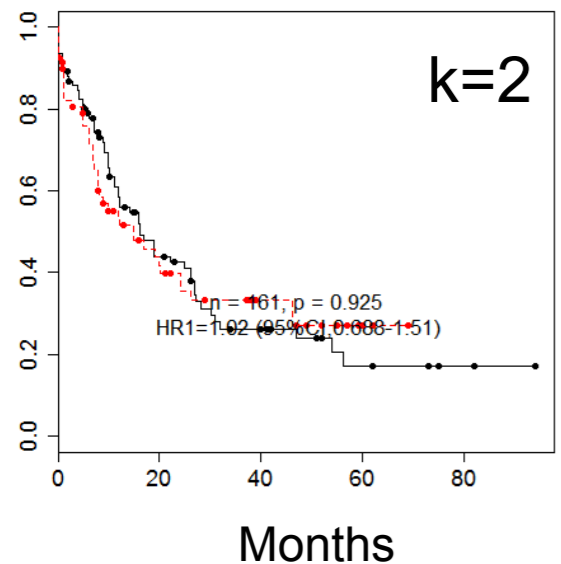
# g. TCGA LGG

1)  $k=2$        $k=3$        $k=4$        $k=5$        $k=6$        $k=7$        $k=8$

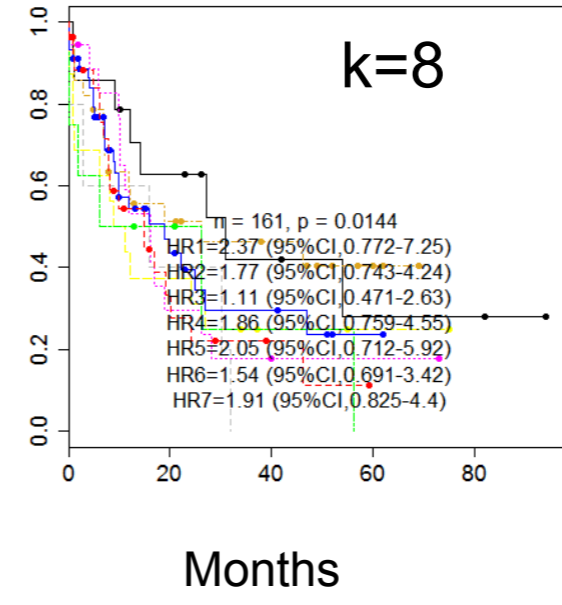
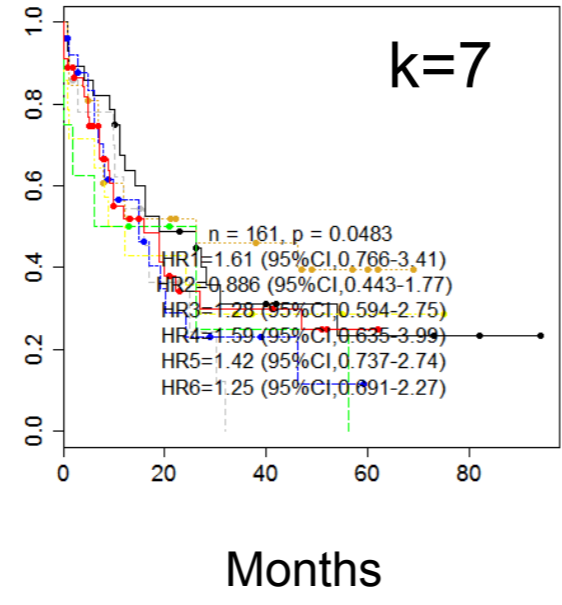
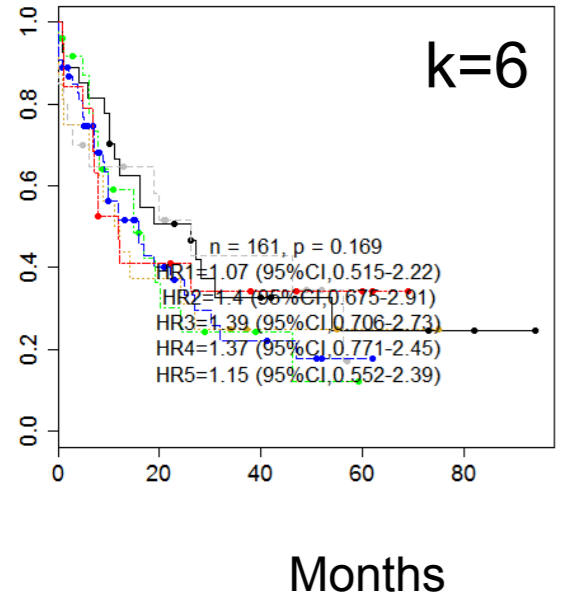


2)

Survival fraction



Survival fraction



# h.TCGA LUAD

1) k=2

k=3

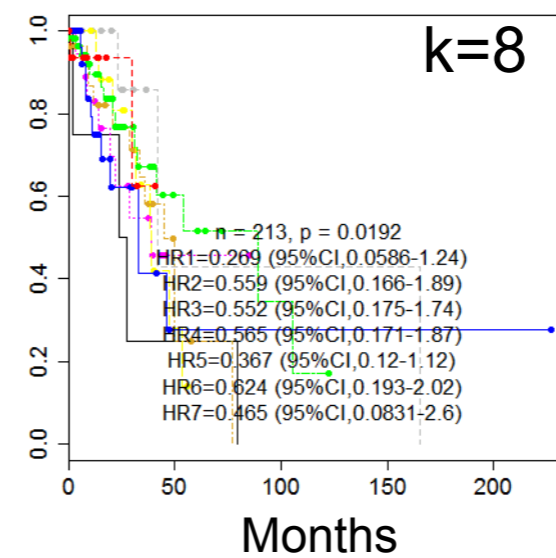
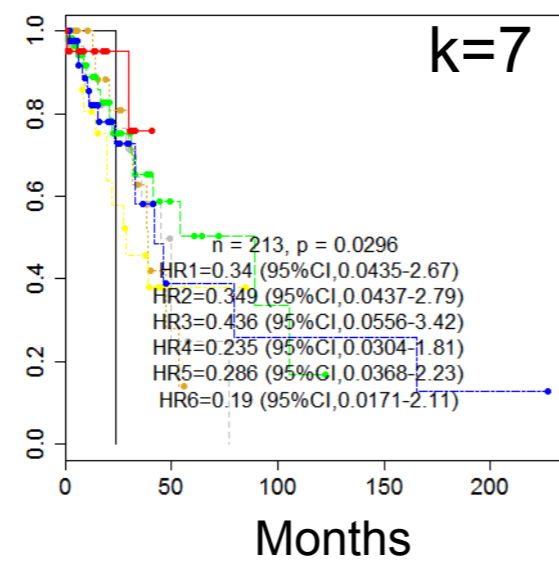
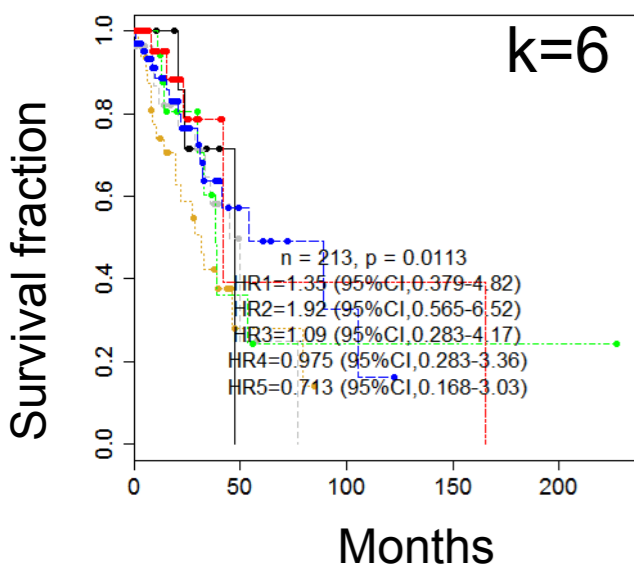
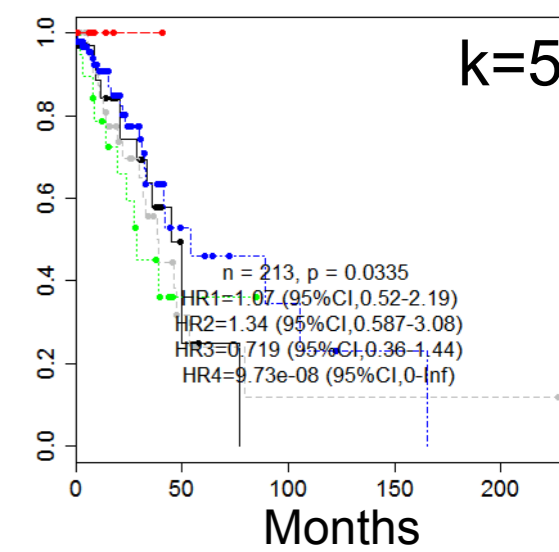
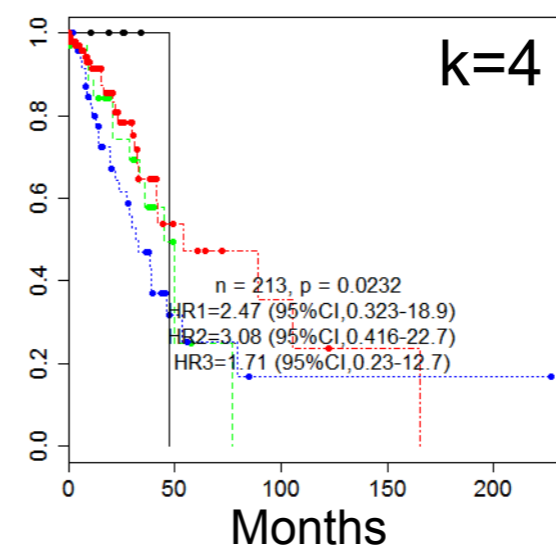
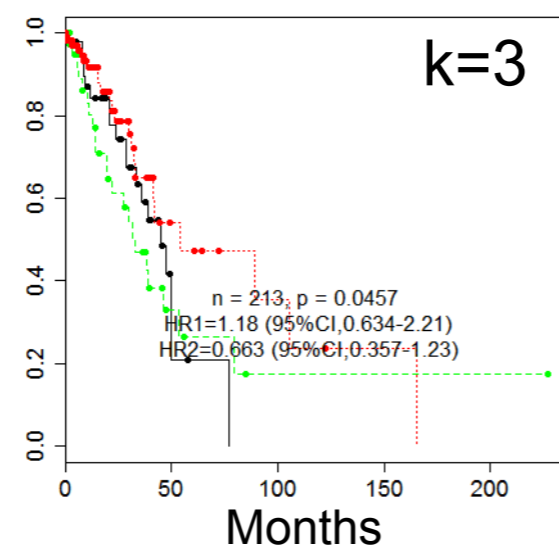
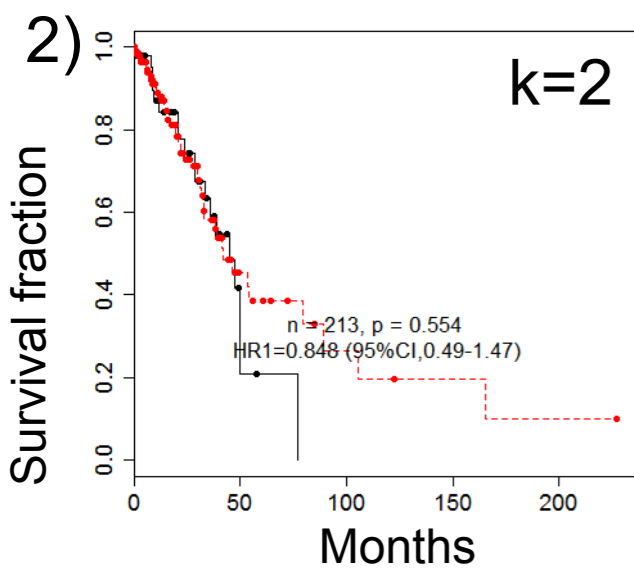
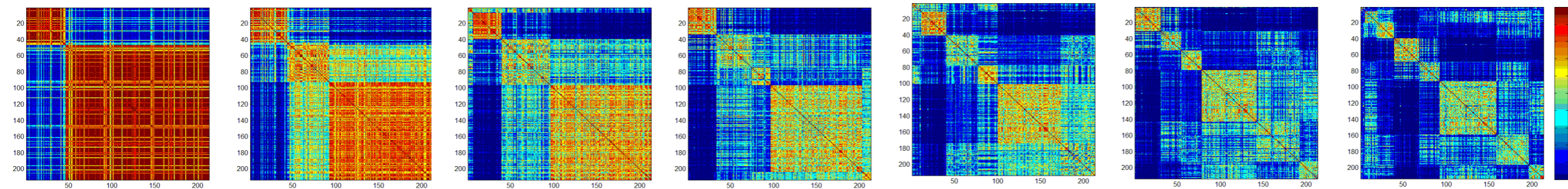
k=4

k=5

k=6

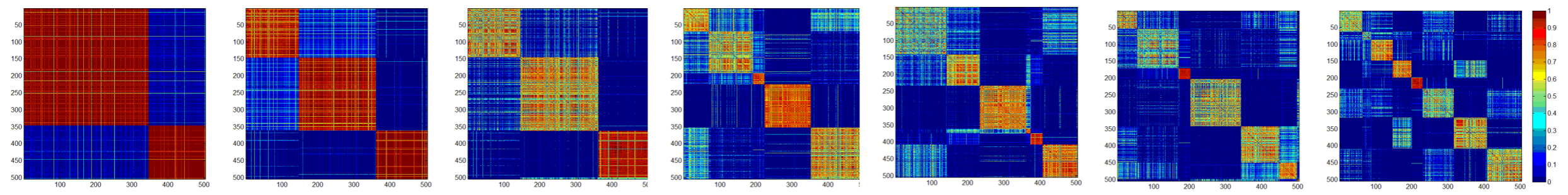
k=7

k=8

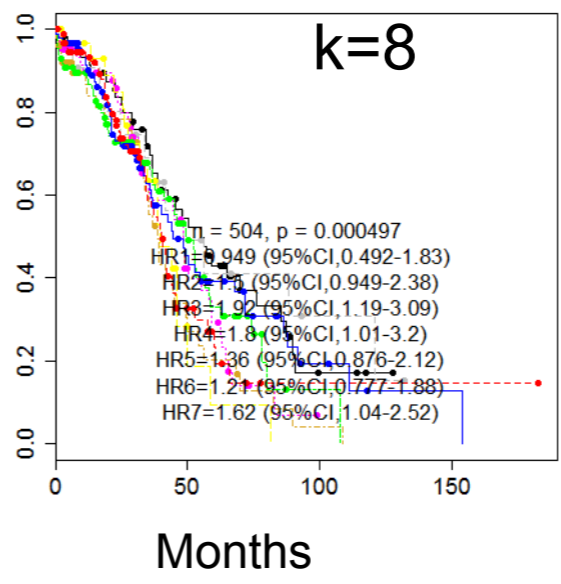
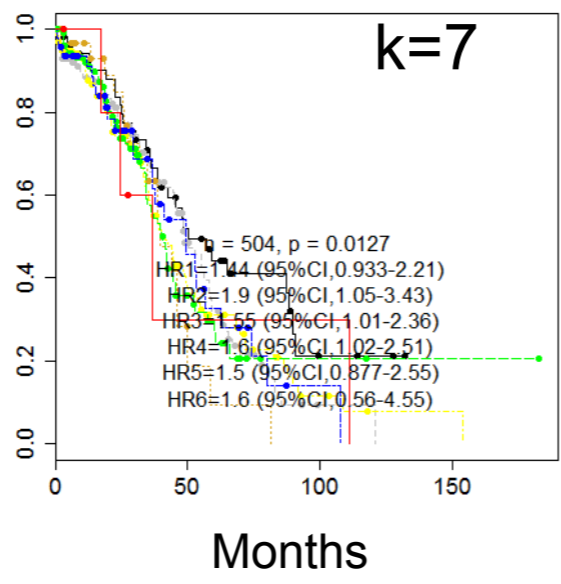
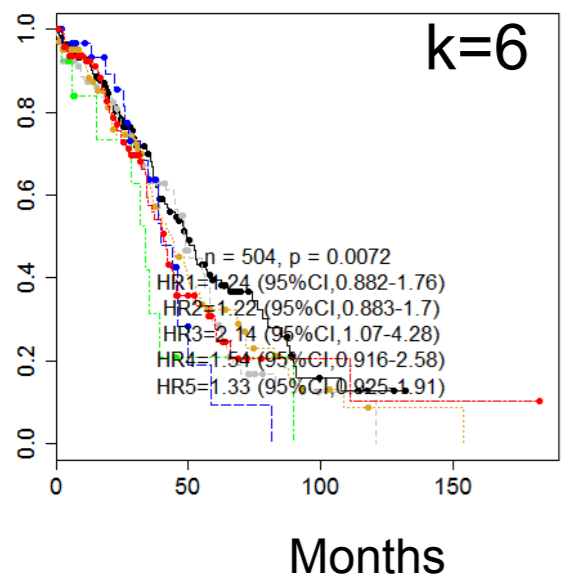
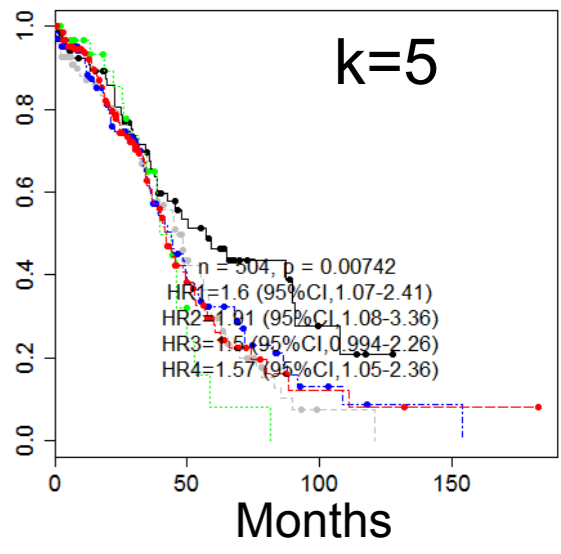
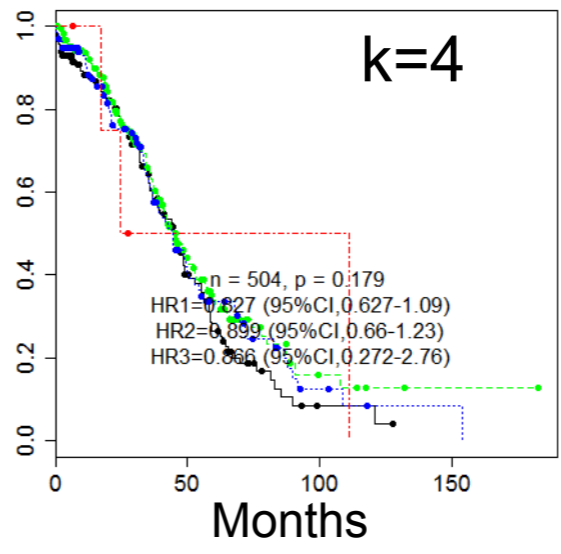
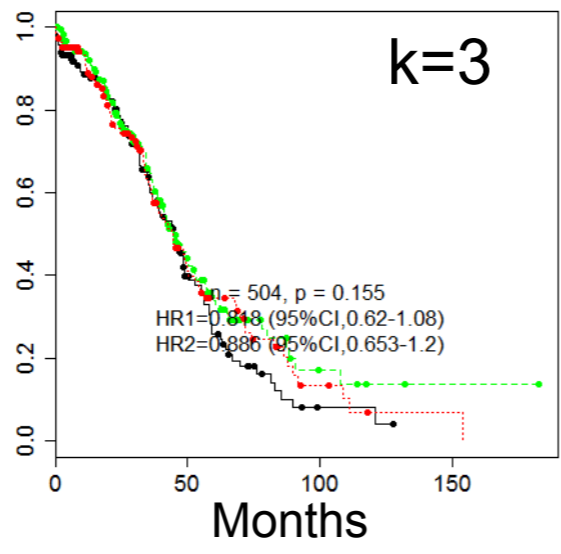
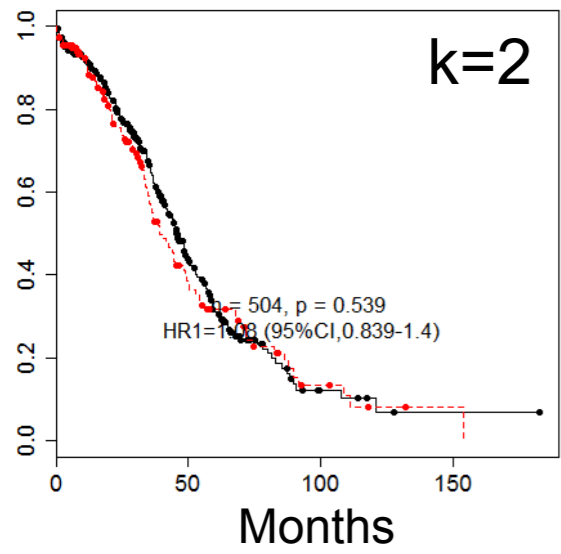


# i. TCGA OV

1) k=2      k=3      k=4      k=5      k=6      k=7      k=8

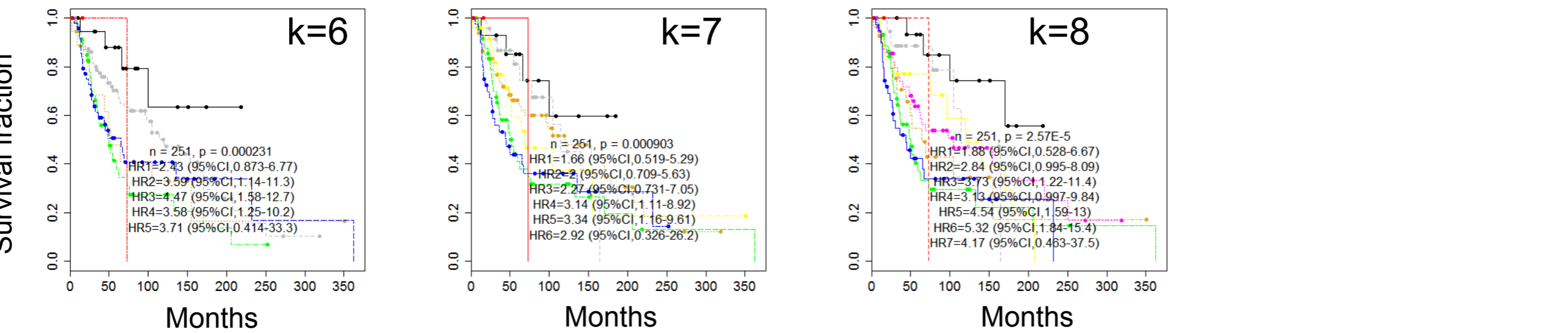
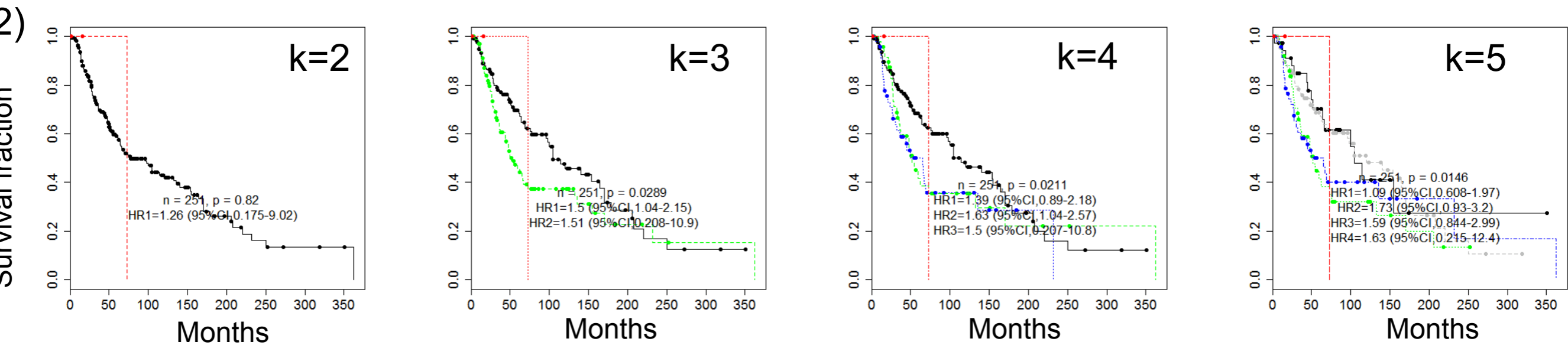
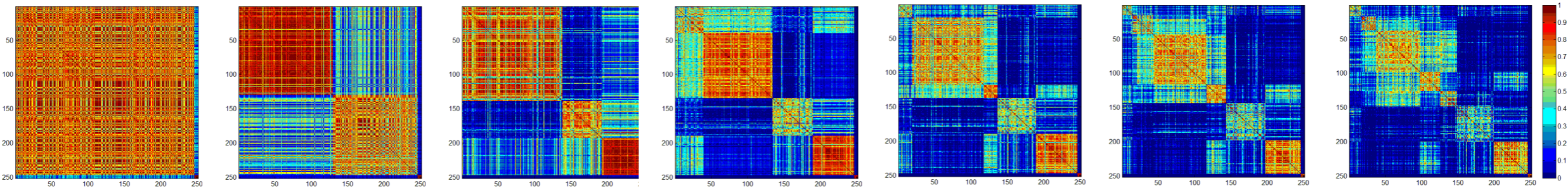


2)



# j.TCGA SKCM

1)  $k=2$        $k=3$        $k=4$        $k=5$        $k=6$        $k=7$        $k=8$

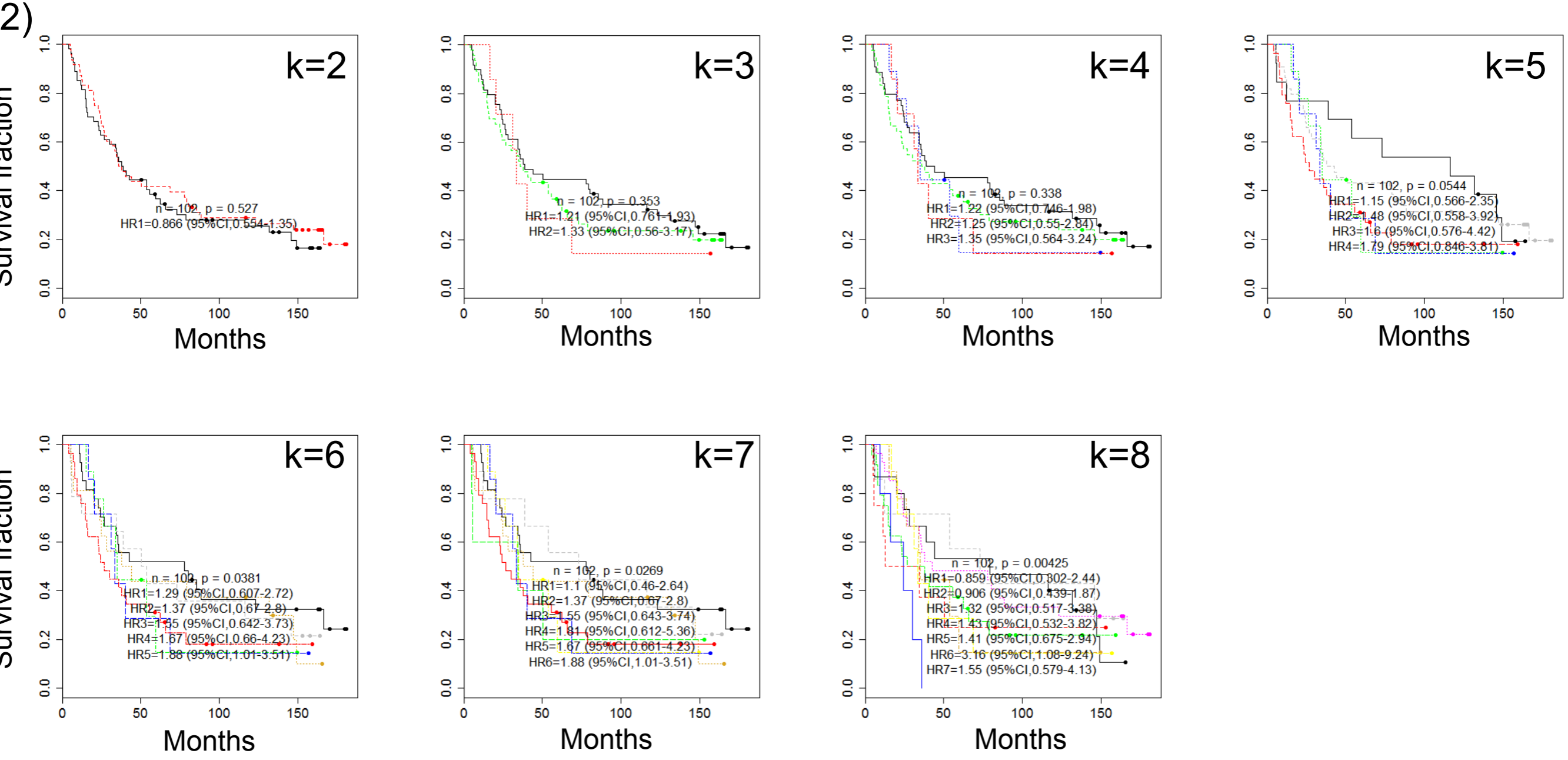
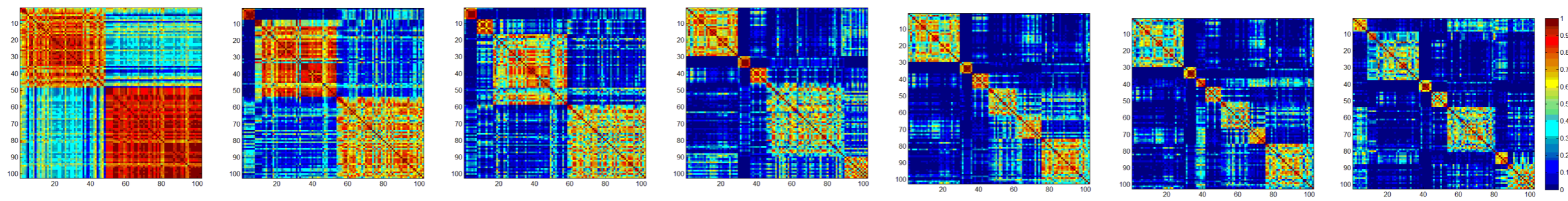


# Supplementary Figure 4

## Independent dataset

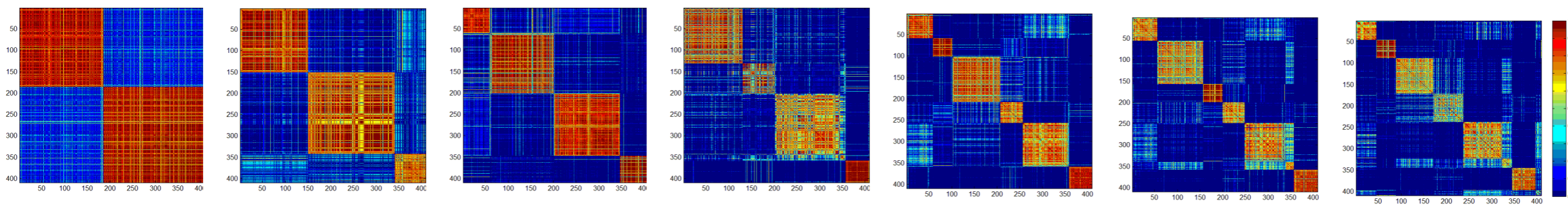
# a. UTSW HNSC

1) k=2      k=3      k=4      k=5      k=6      k=7      k=8



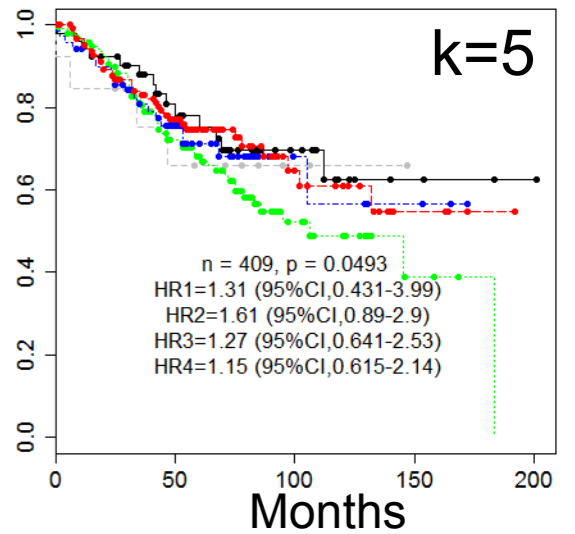
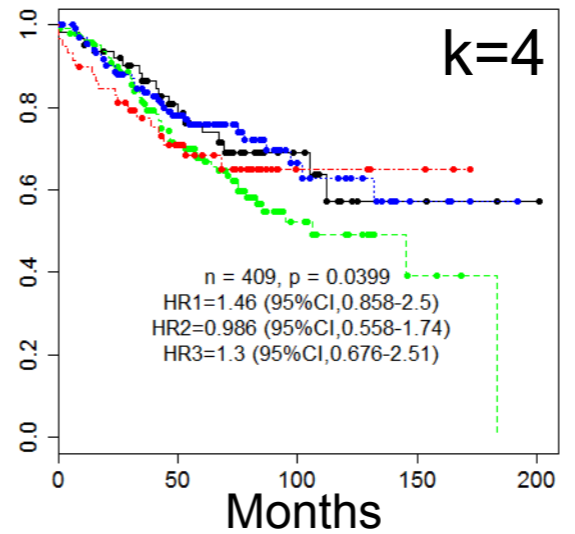
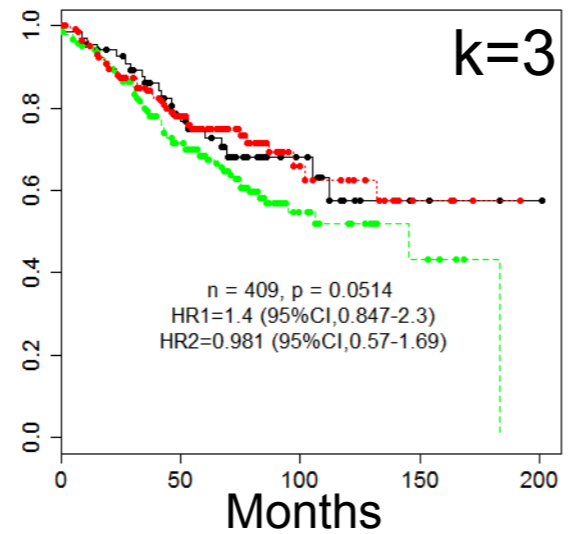
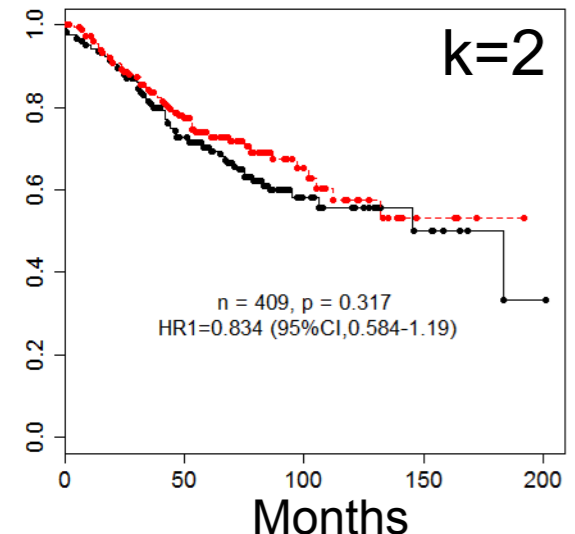
# b. Colon cancer (GSE39582)

a      k=2                      k=3                      k=4                      k=5                      k=6                      k=7                      k=8

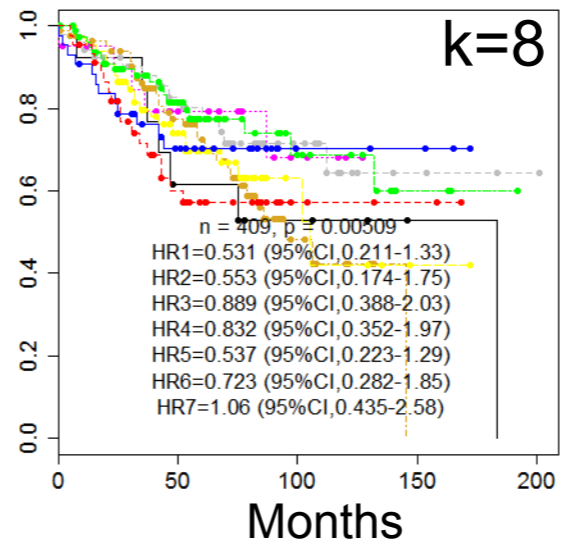
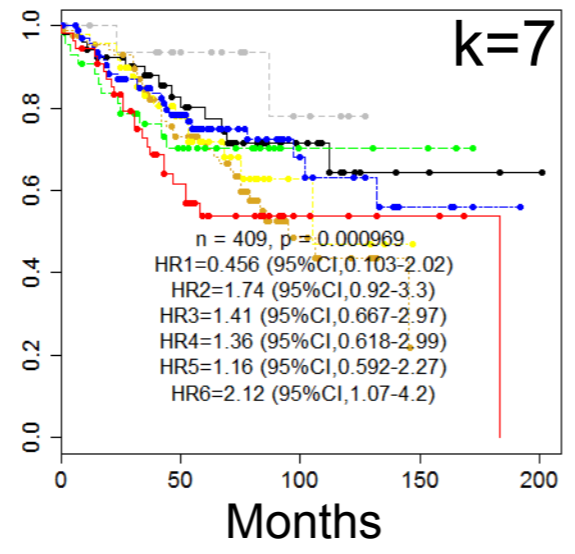
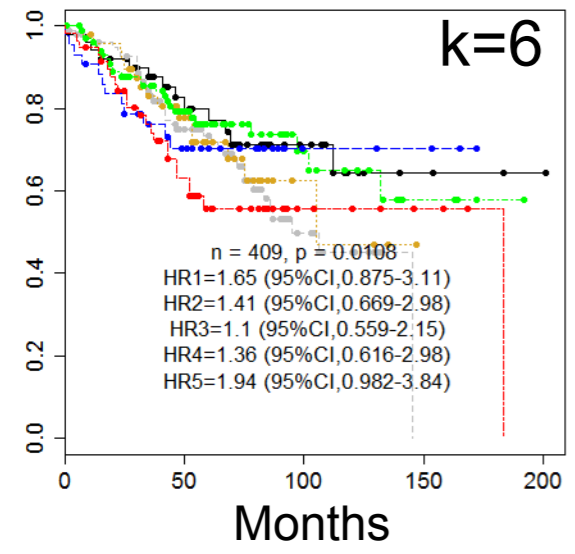


b

Survival fraction

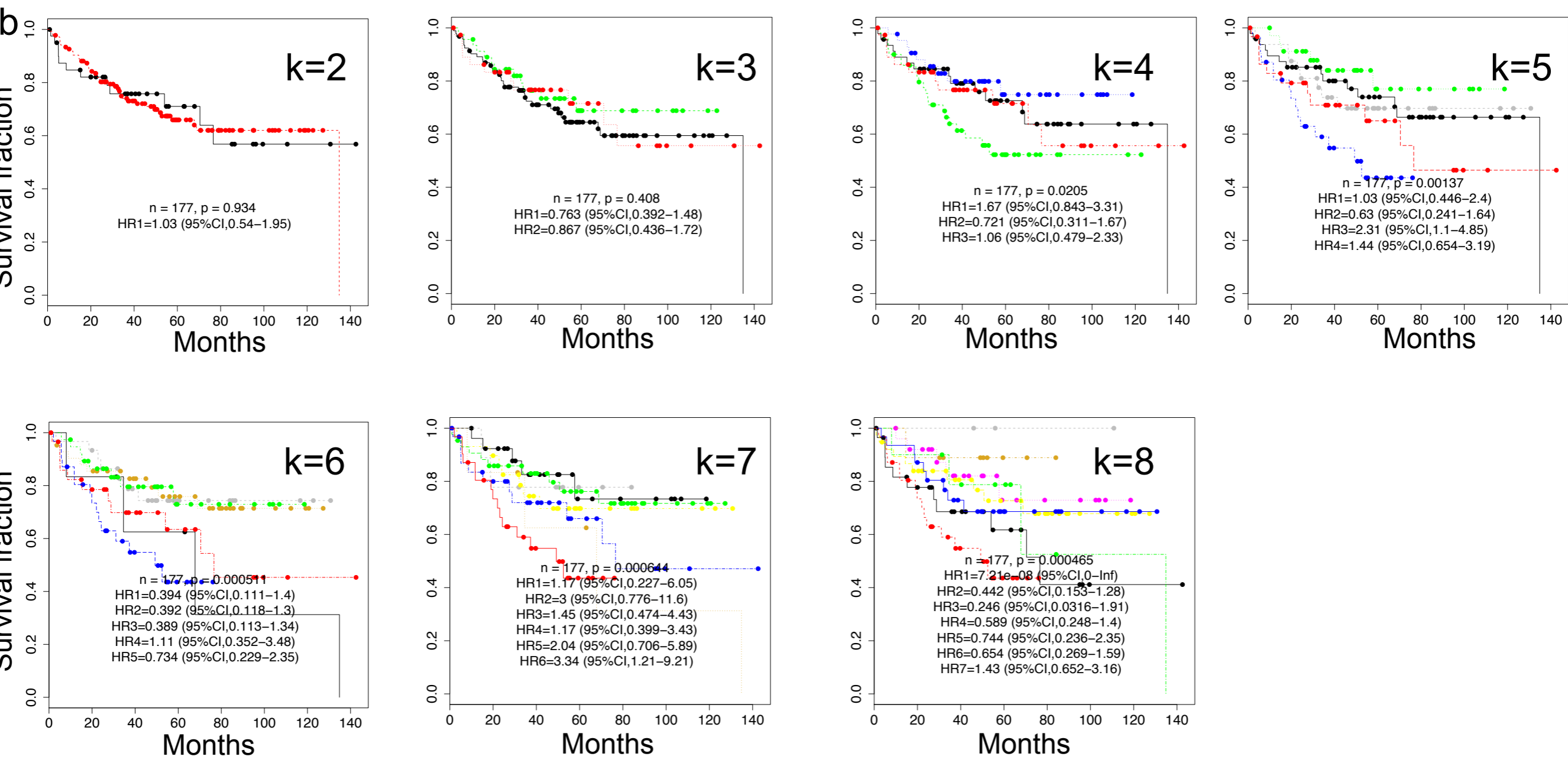
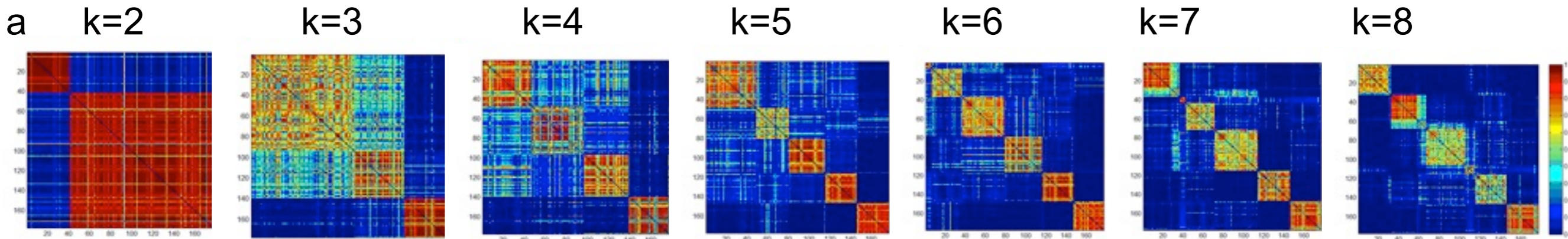


Survival fraction



Months

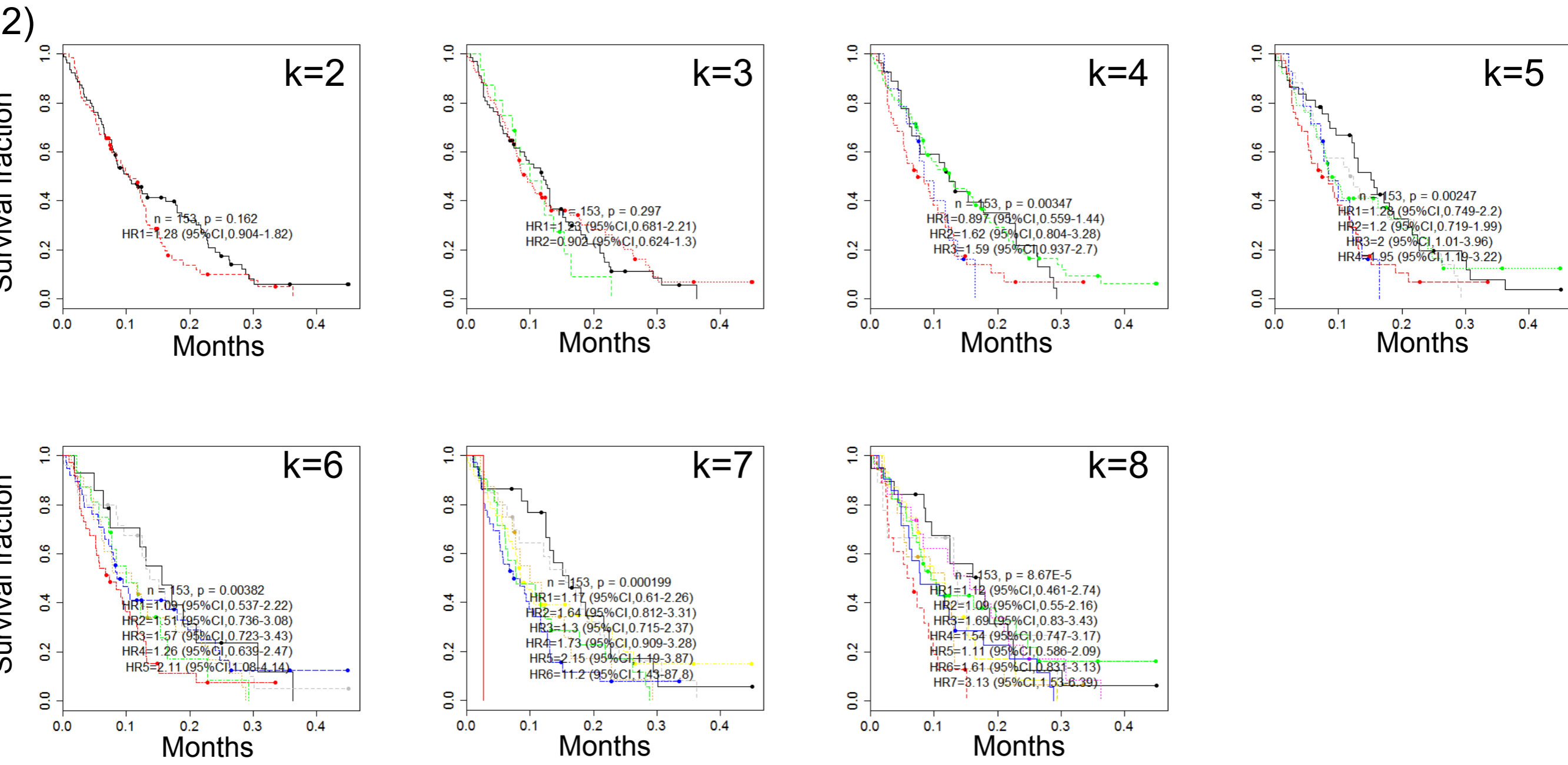
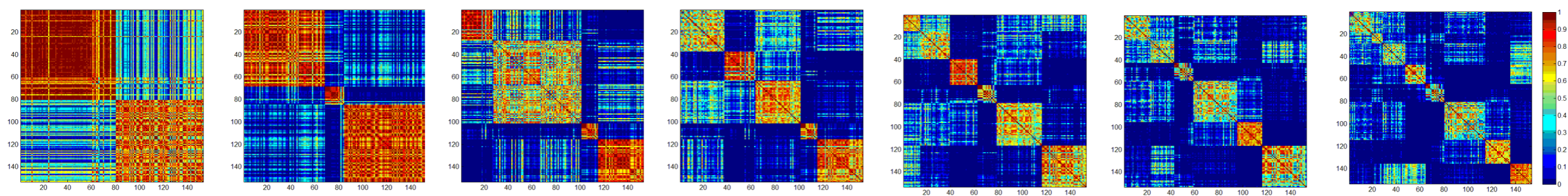
# b. Colon cancer (GSE17538)





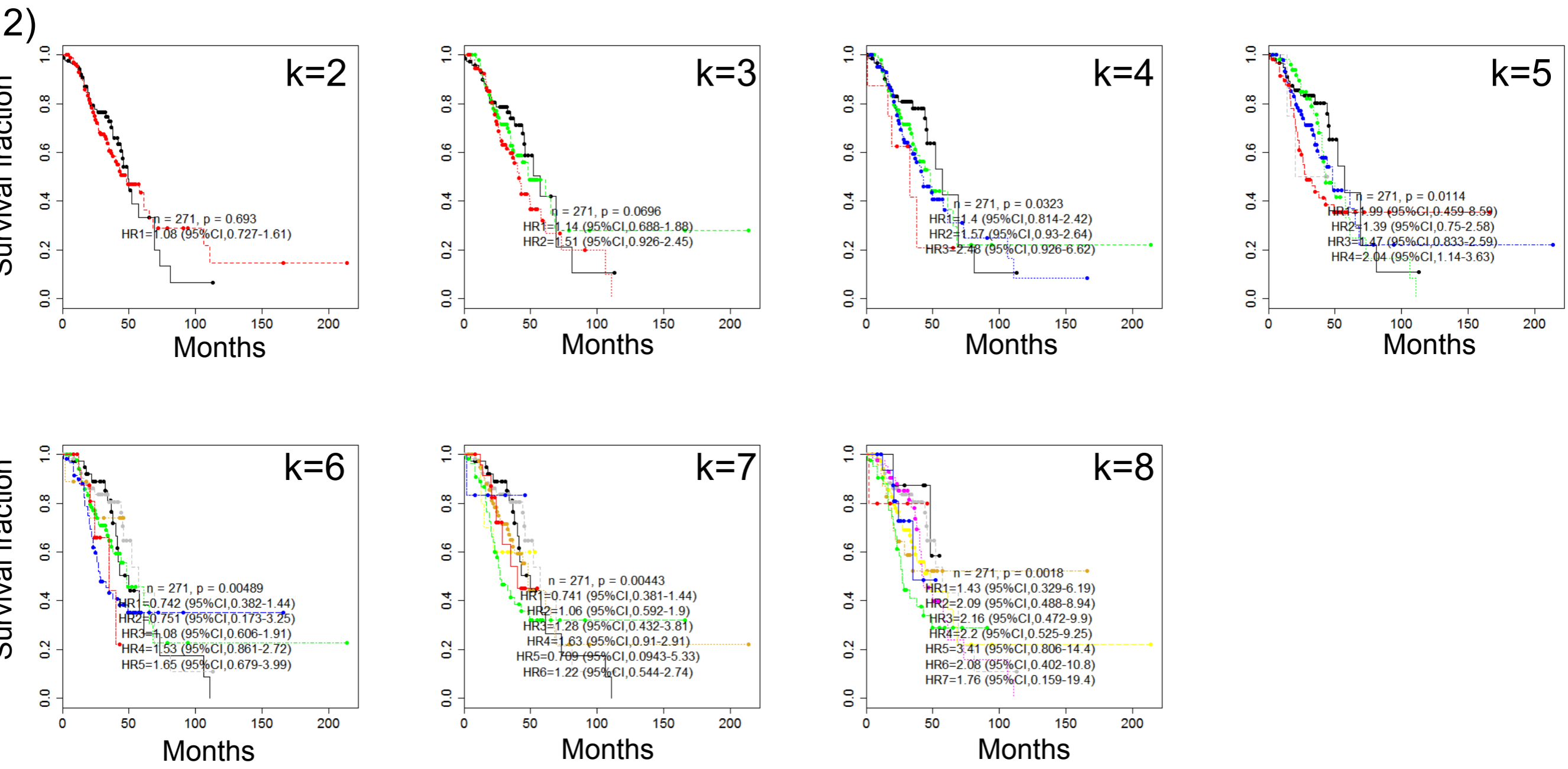
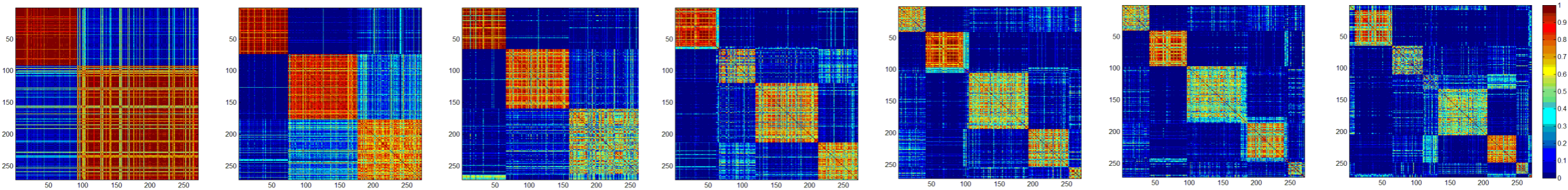
c. OV (Bonome et. al., Cancer Research 2008, )

1) k=2                      k=3                      k=4                      k=5                      k=6                      k=7                      k=8



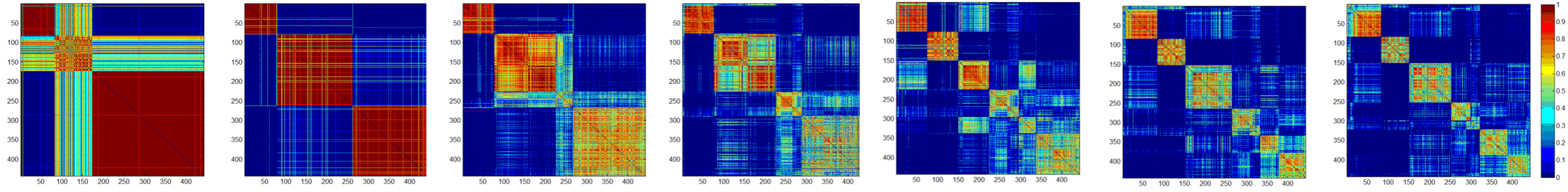
d. OV (Tothill et. al., Clinical Cancer Research, 2008)

1) k=2      k=3      k=4      k=5      k=6      k=7      k=8

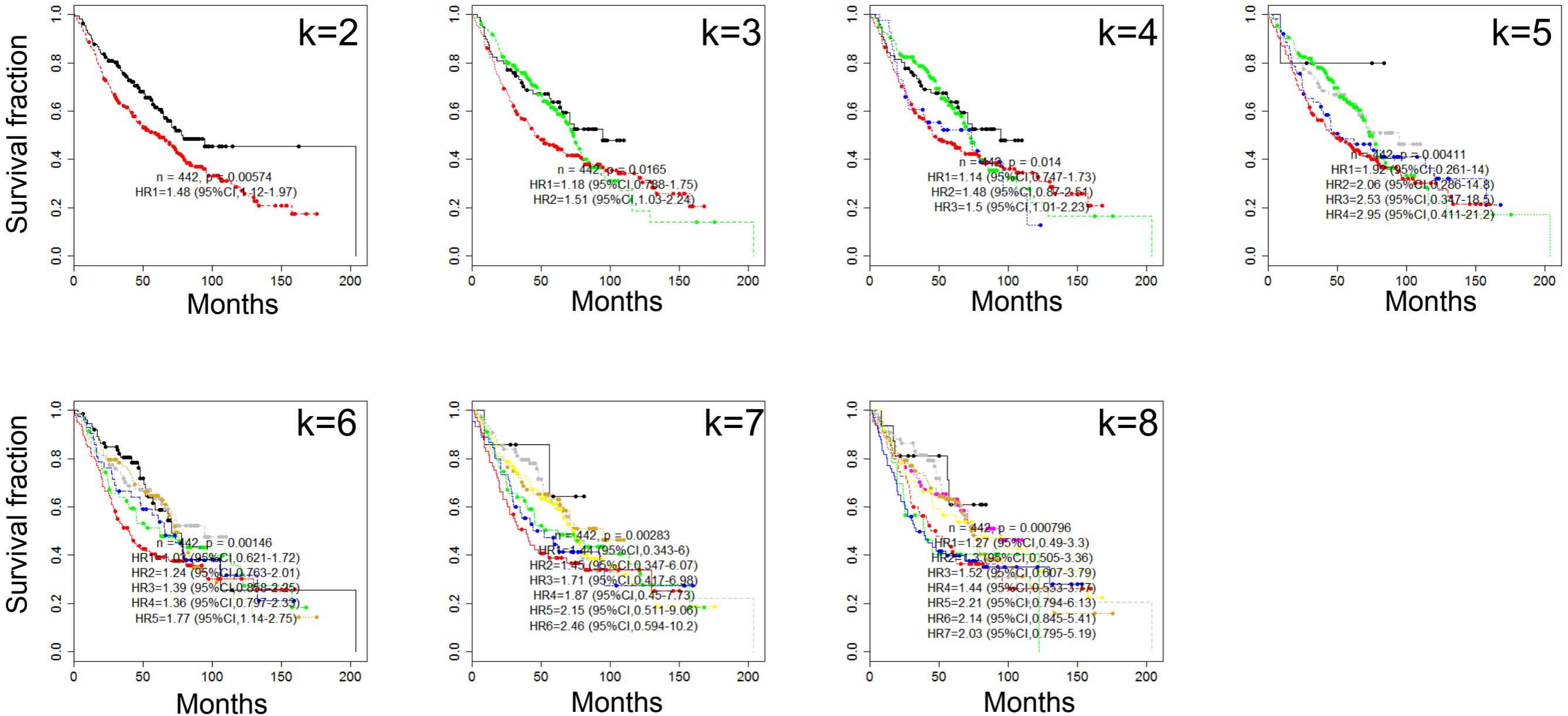


e. Lung adenocarcinoma (Shedden et. al., Nature Methods, 2008)

1) k=2      k=3      k=4      k=5      k=6      k=7      k=8



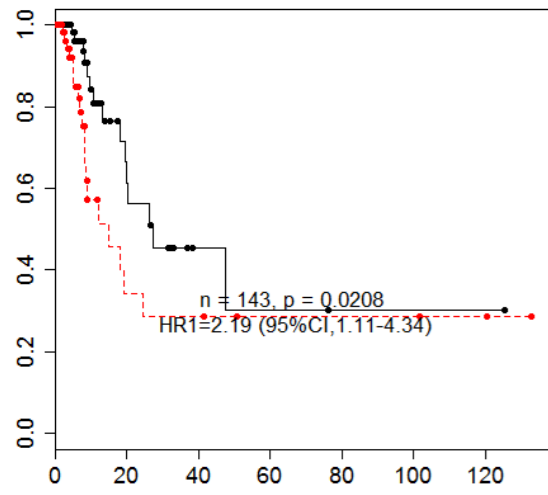
2)



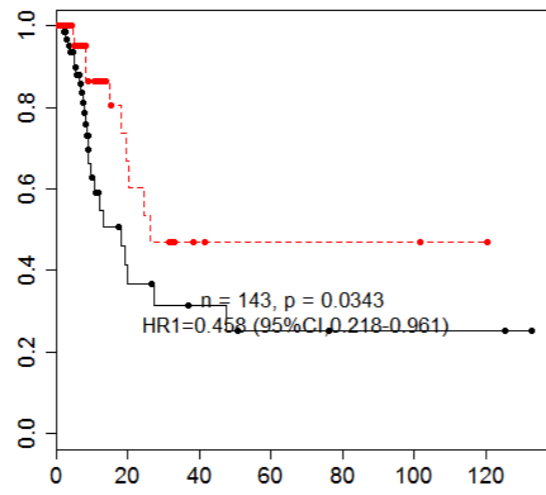
# Supplementary Figure 5

## Single gene analysis

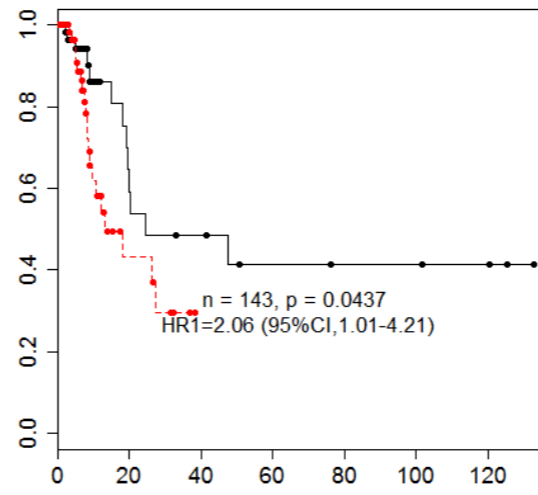
# a. TCGA BLCA



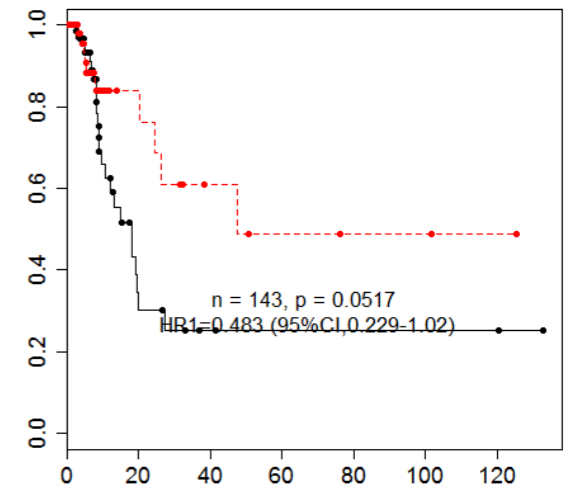
TP53



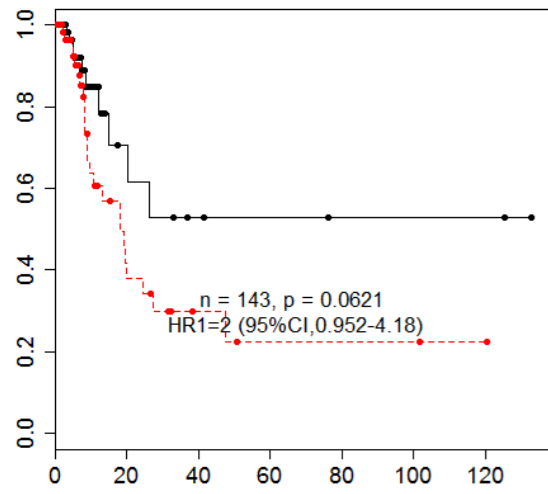
ASCC2



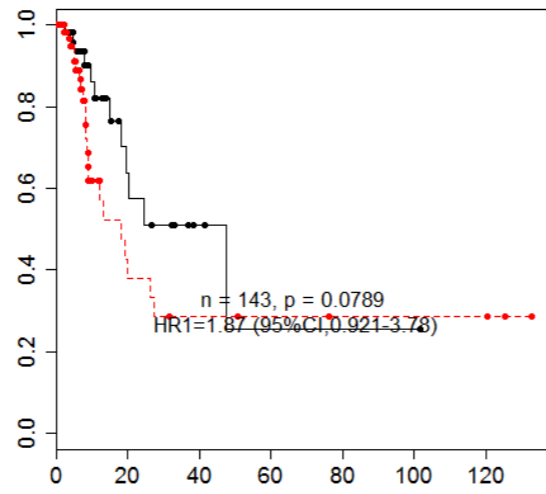
WRN



FHL2

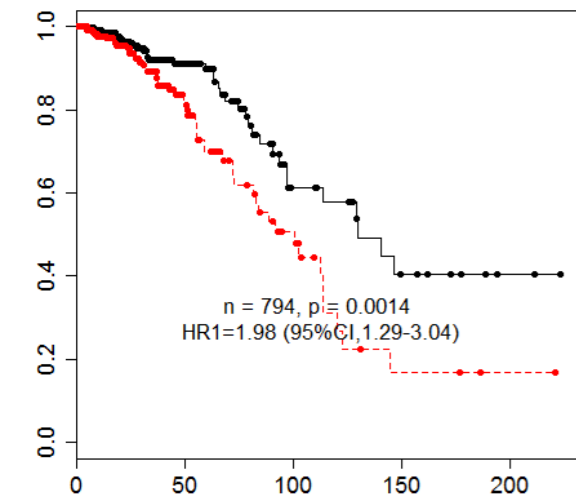


NCOA6

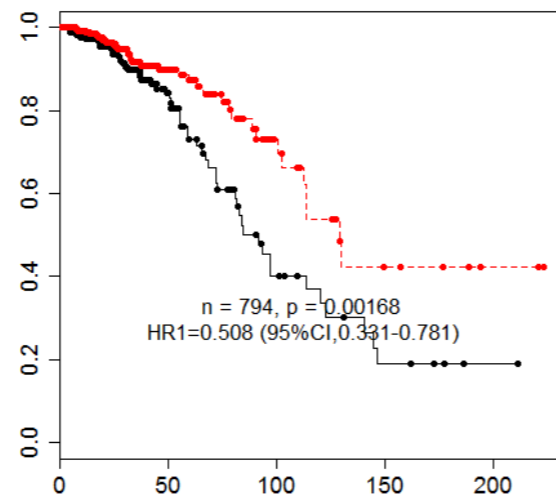


MSH6

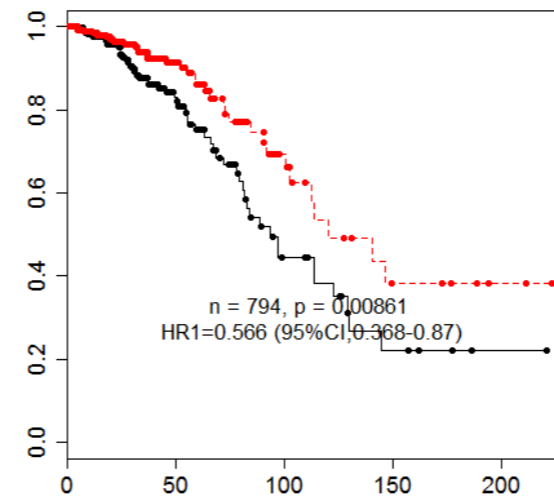
## b. TCGA BRCA



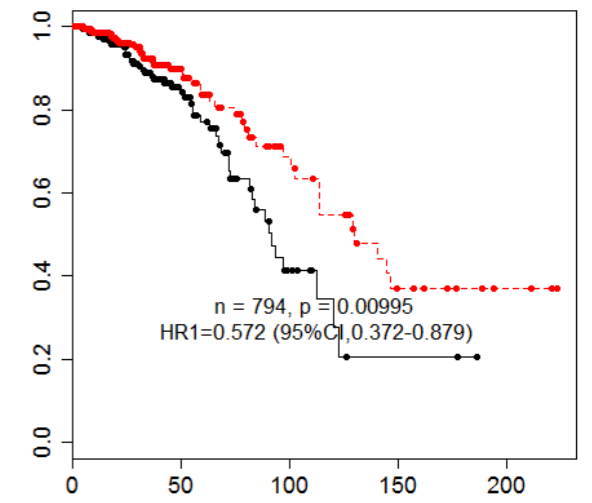
ZHX1



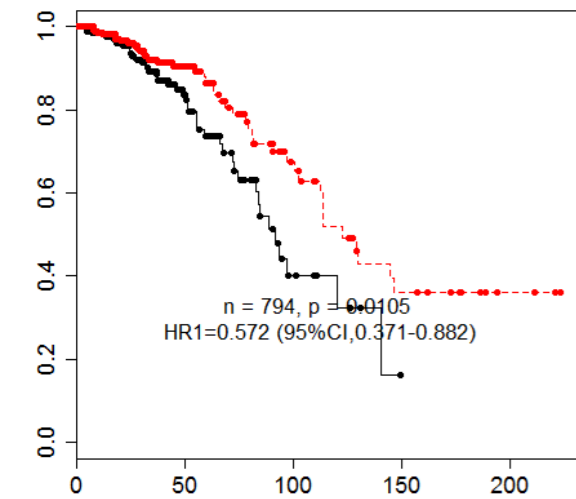
NFKBIA



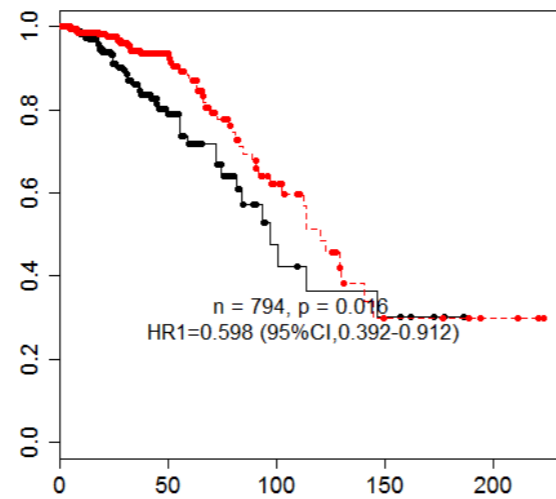
ACTA2



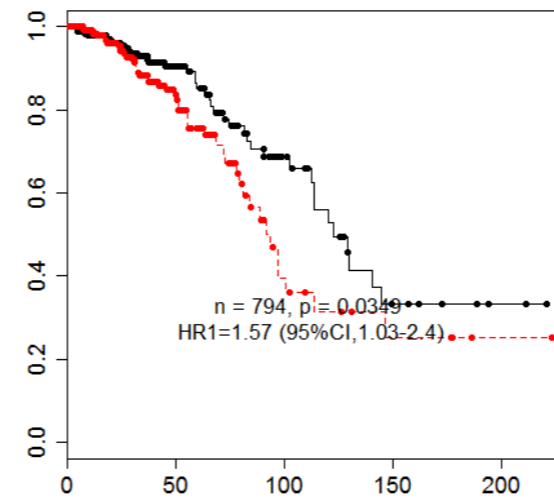
PIN1



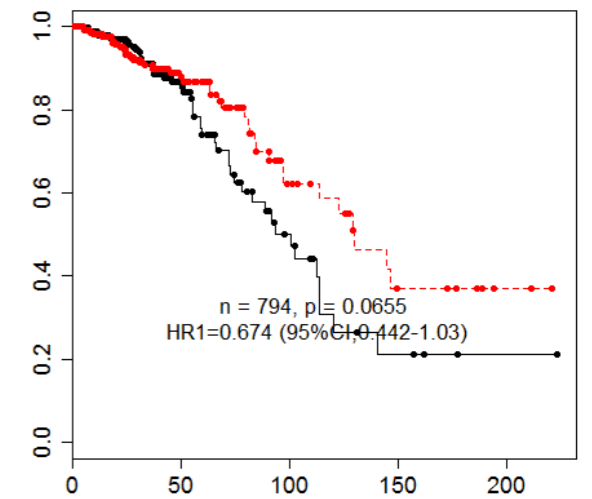
THAP8



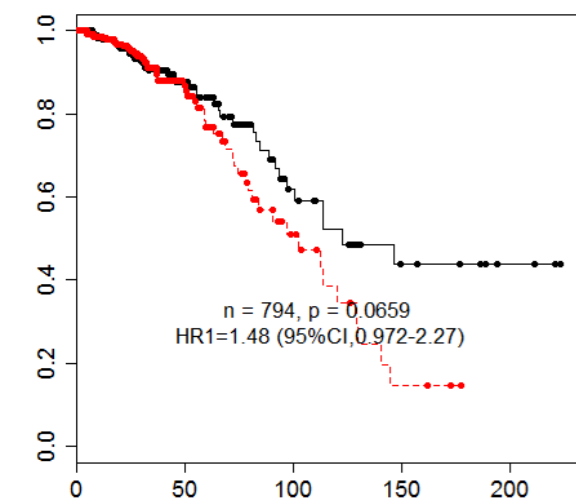
NAT6



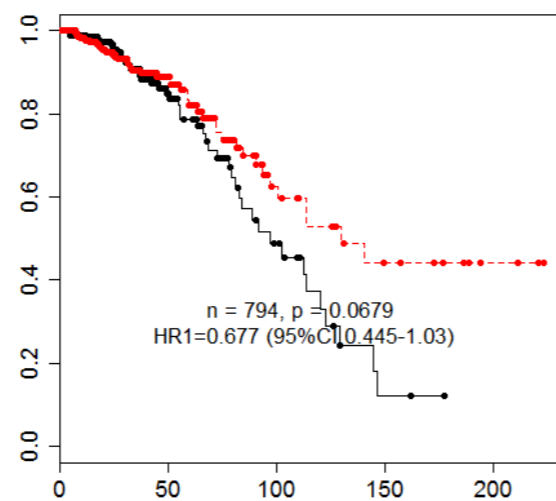
BRCA2



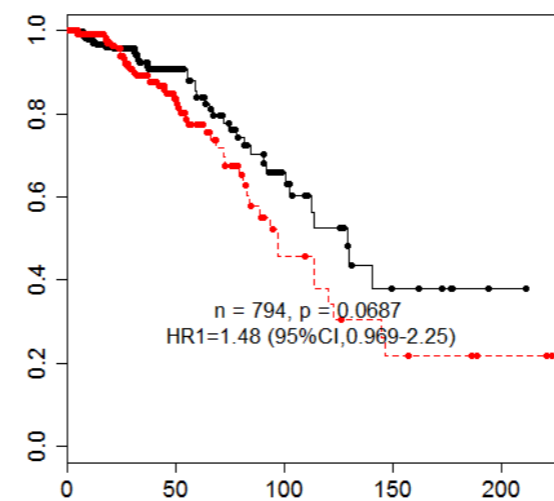
CDK5



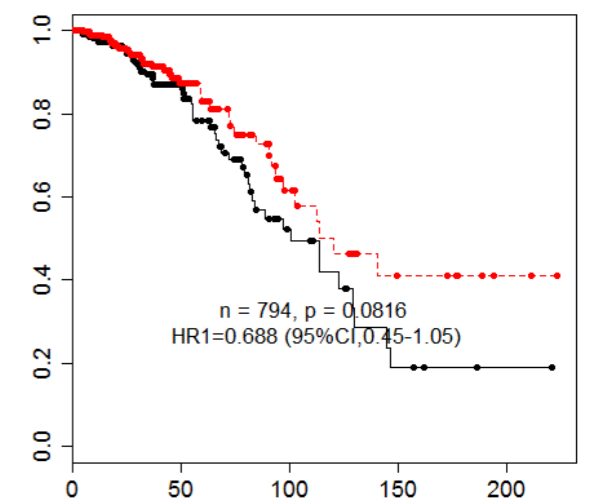
CHD3



YBX1

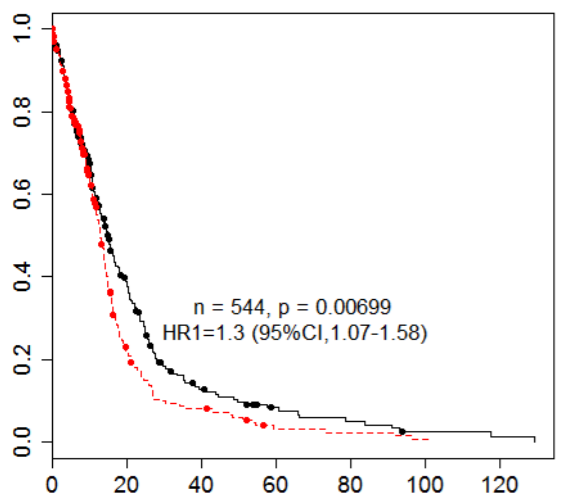


TDRD7

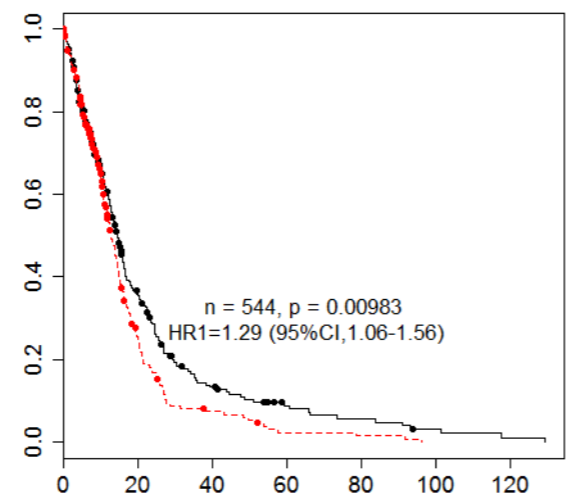


VIM

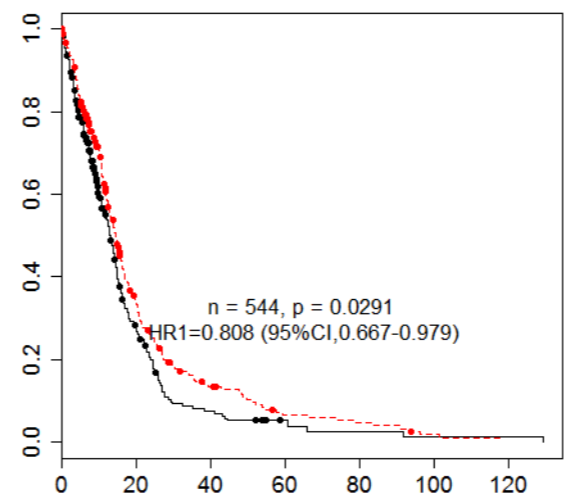
# c. TCGA GBM



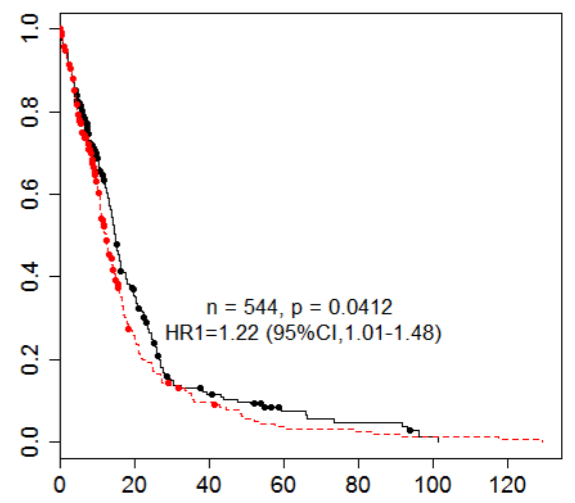
TEP1



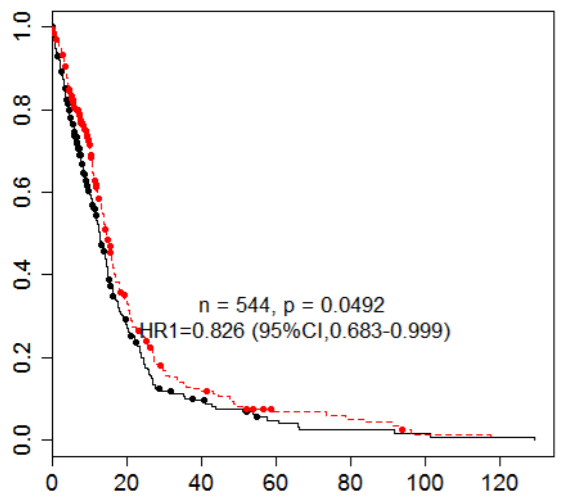
PTRF



EEF2

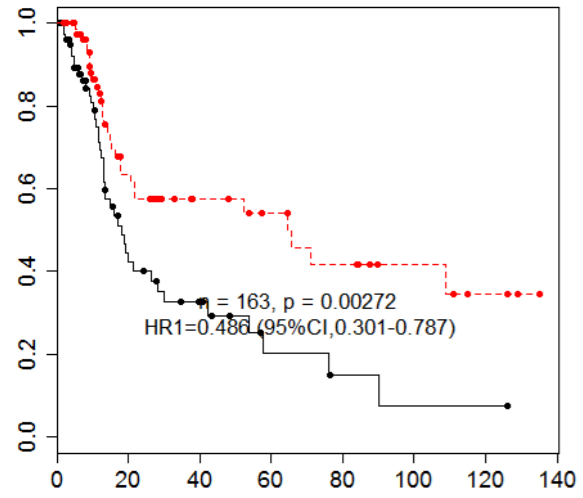


GNL3

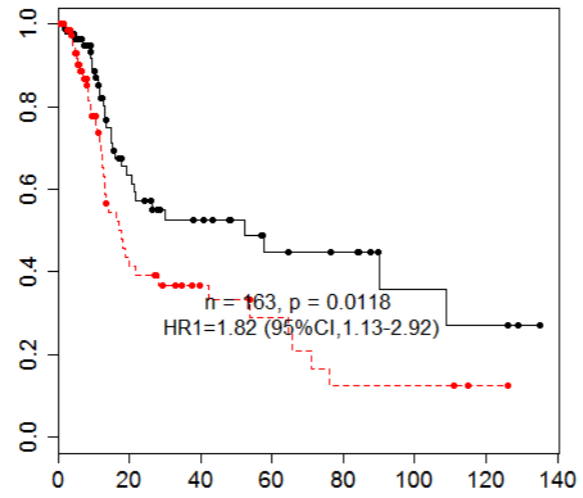


TOP1

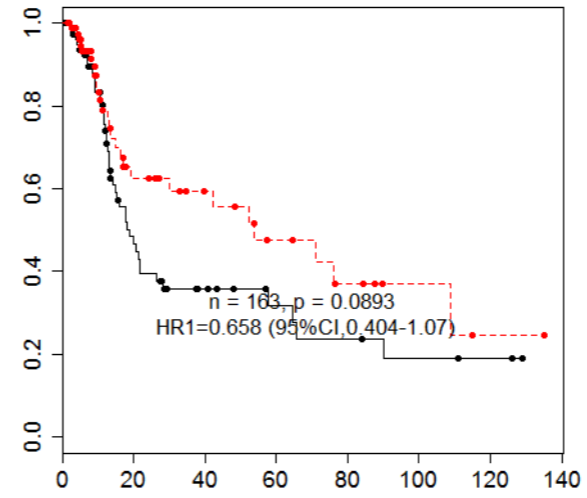
# d. TCGA HNSC



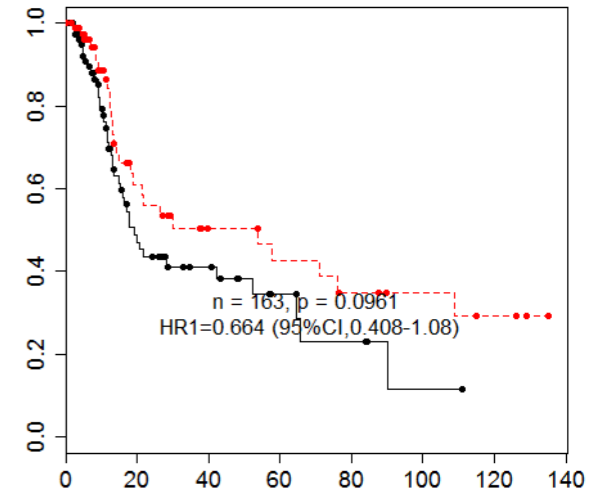
CABLES1



TOP1



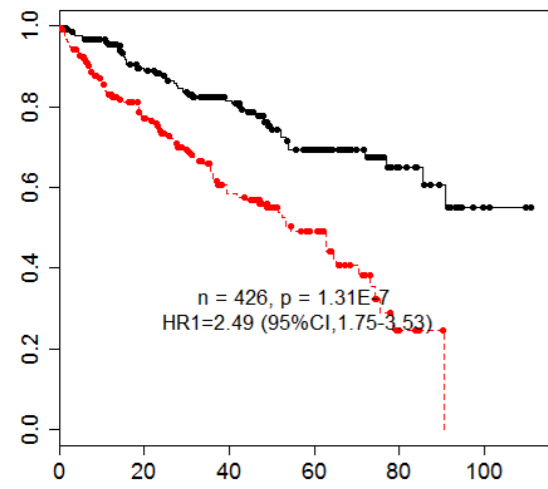
EEF2



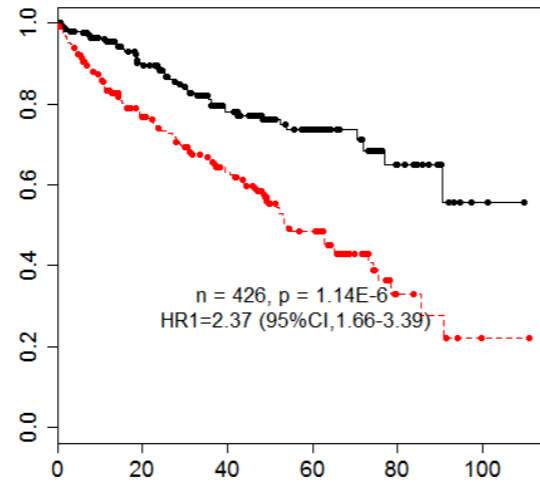
HIPK2



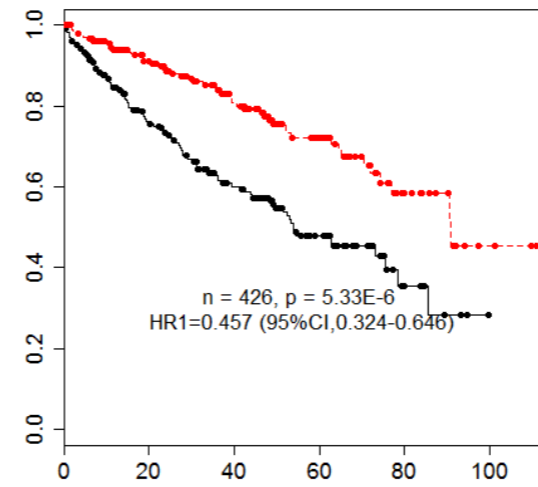
# e. TCGA KIRC



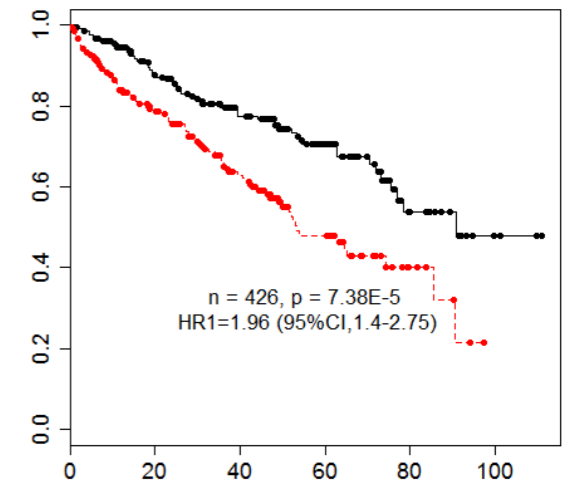
**GNL3**



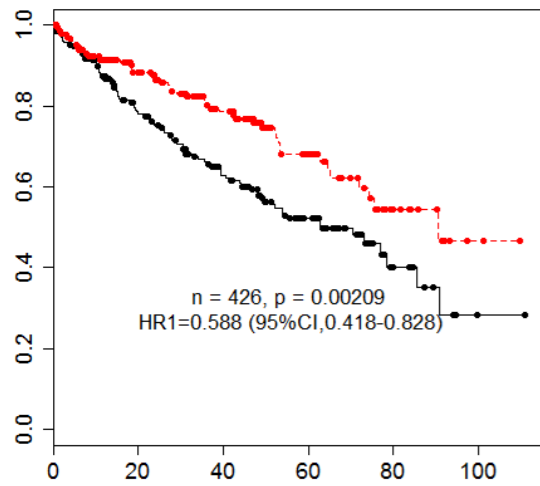
**SERPINA3**



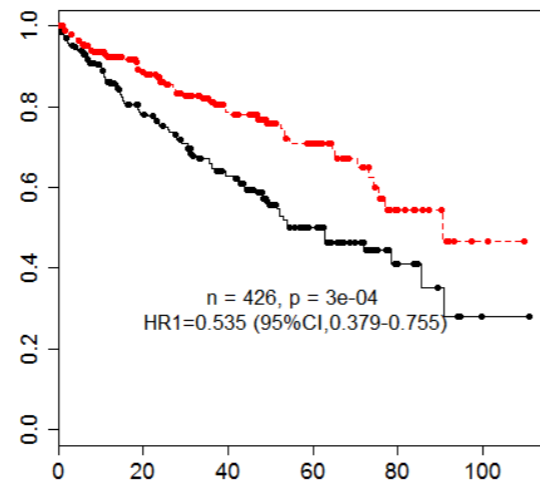
**HIPK2**



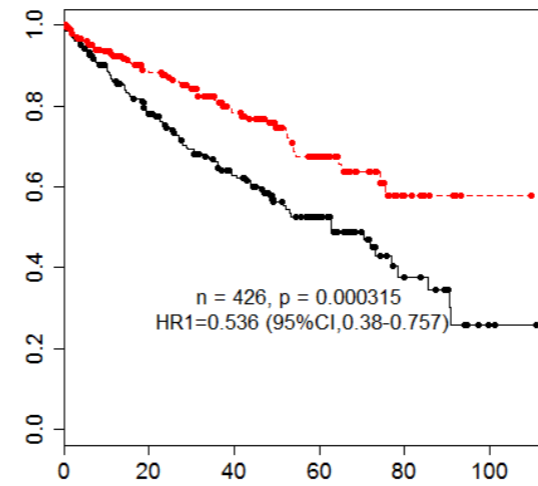
**BEST1**



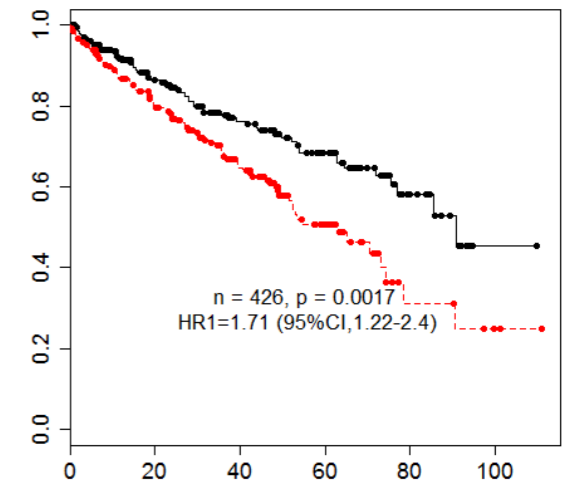
**CREBBP**



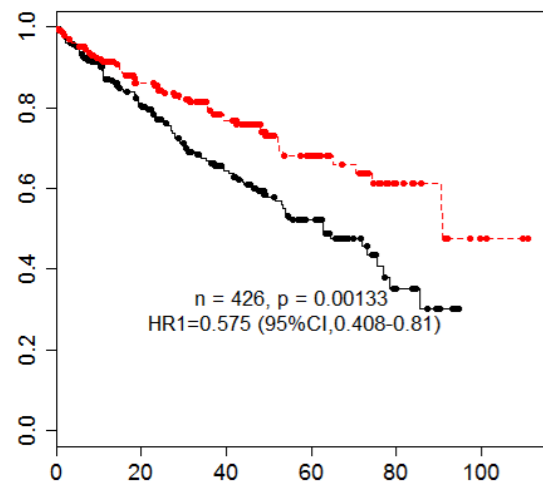
**EP300**



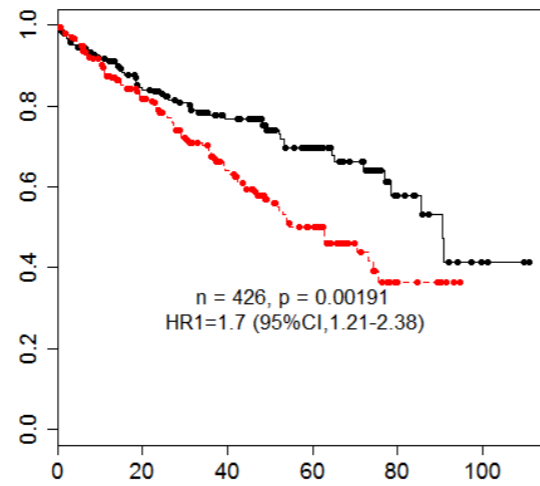
**APP**



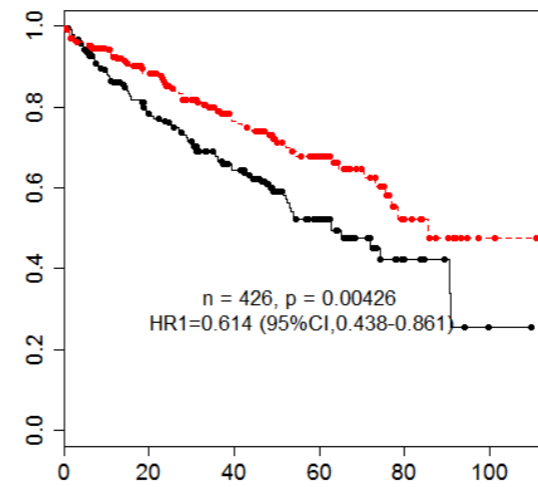
**MDM4**



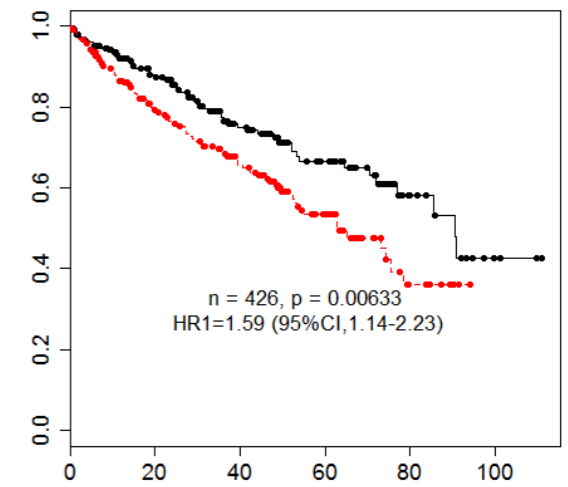
**PPP1R13B**



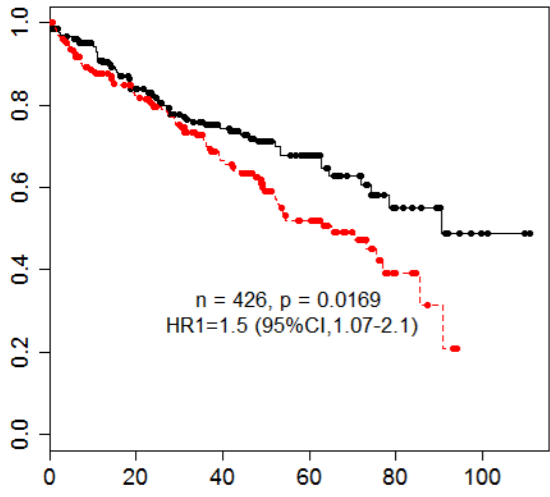
**MMP7**



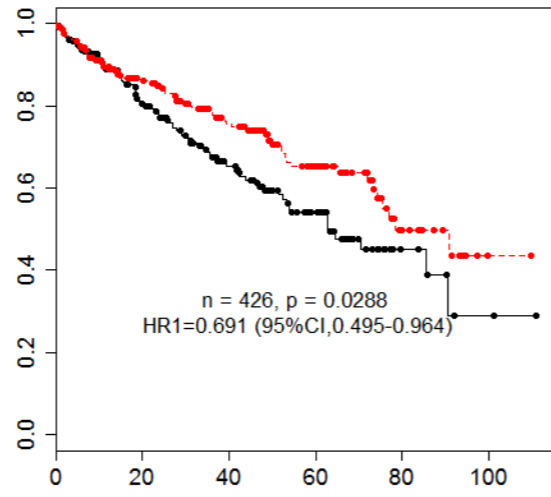
**USP33**



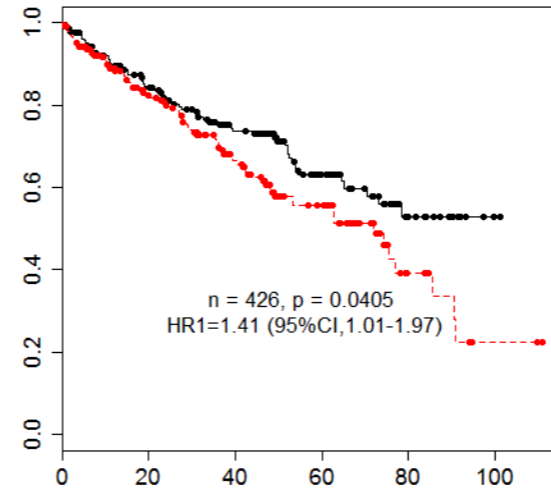
**MMP9**



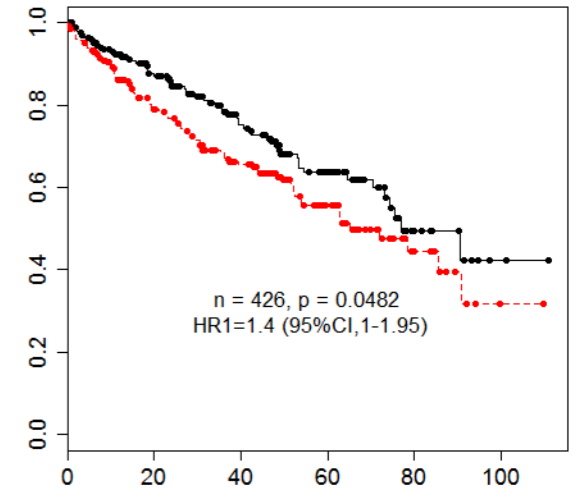
TP53



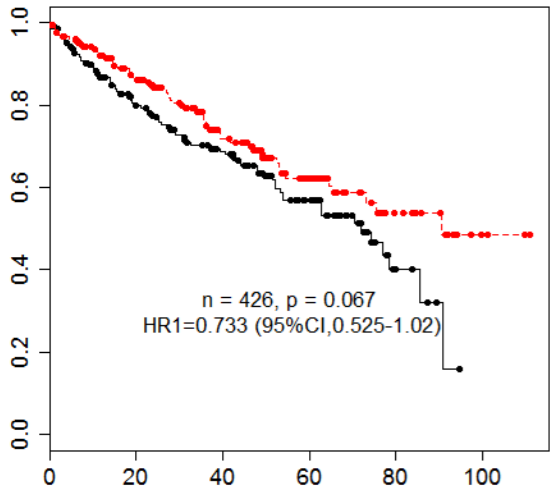
CTSG



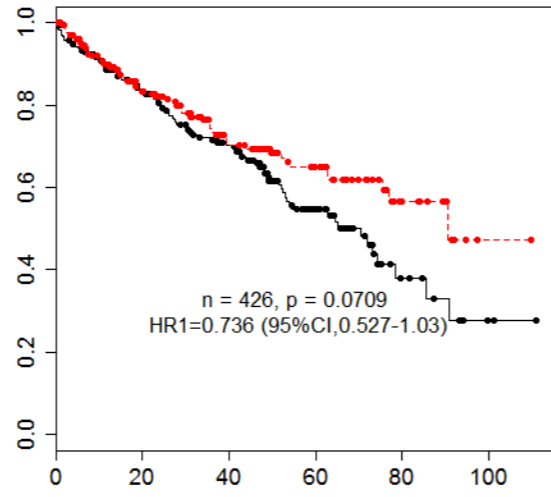
SERPINA1



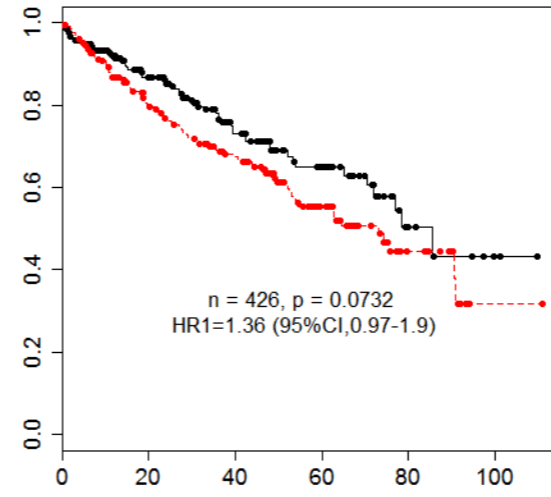
TCEB2



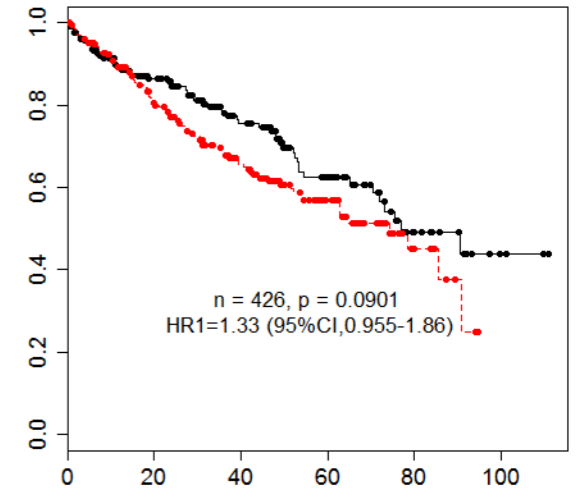
SP1



THBS1

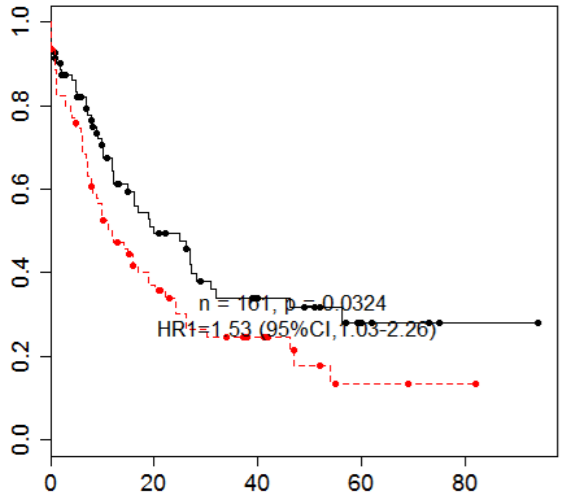


WT1

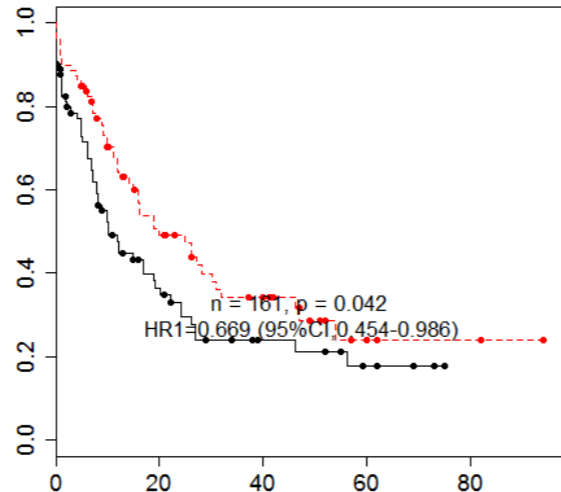


TCEB1

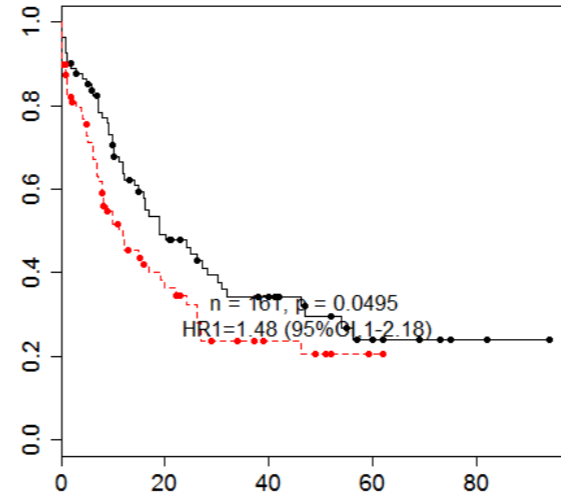
# f. TCGA LAML



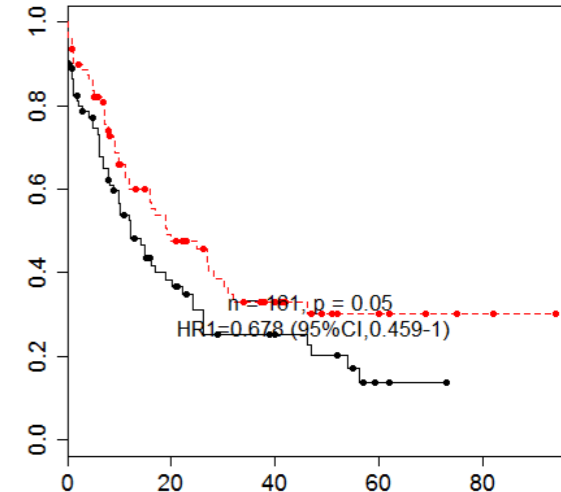
SHFM1



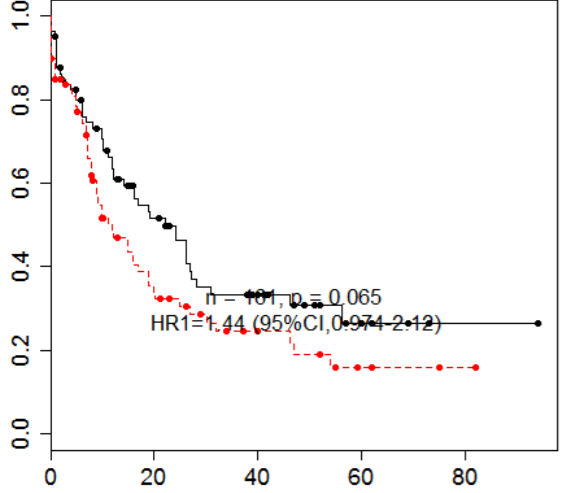
NPM1



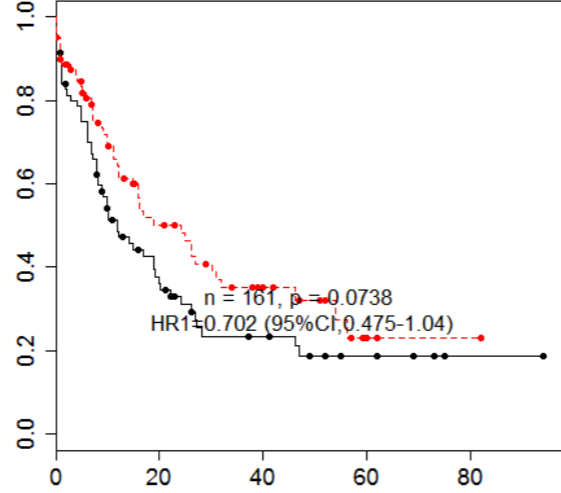
HSPA1A



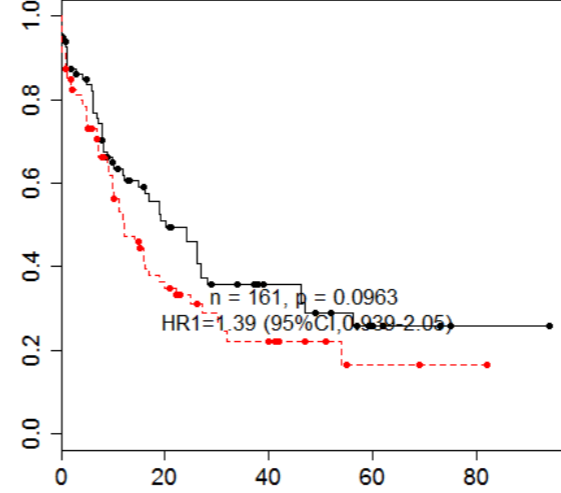
TP53



NQO1

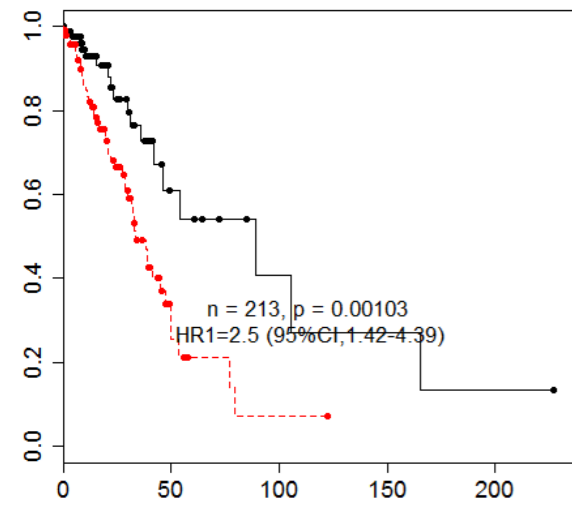


WDR16

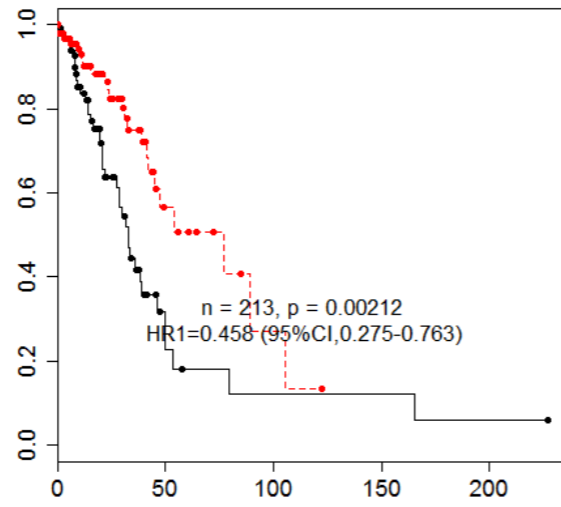


PARP1

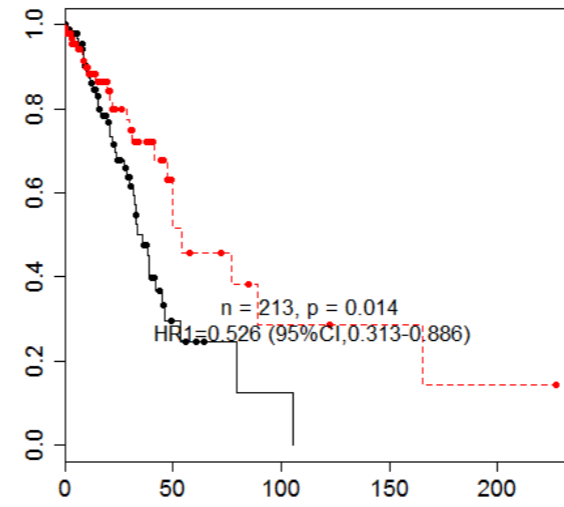
# g. TCGA LUAD



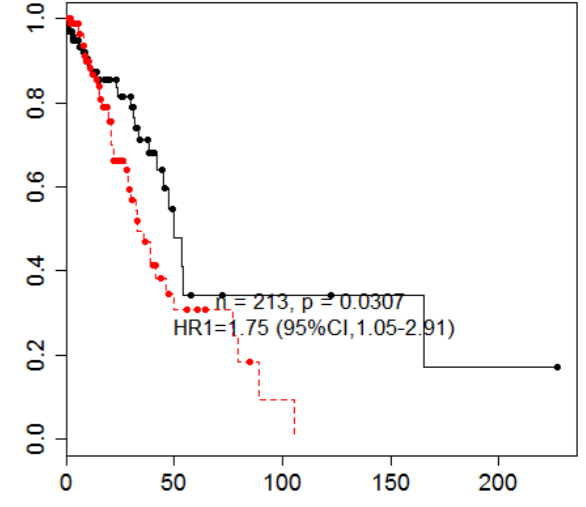
YBX1



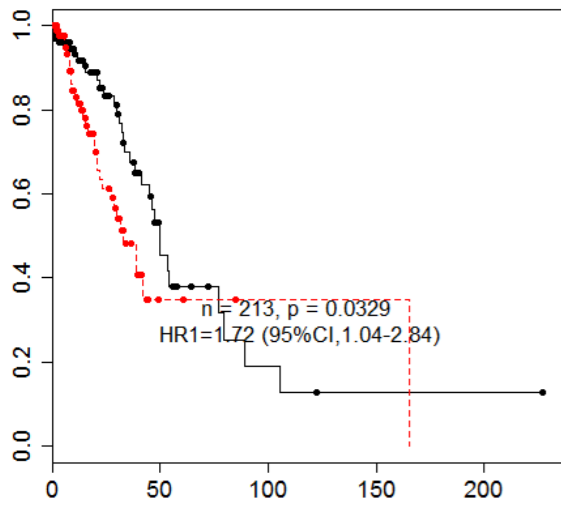
CABLES1



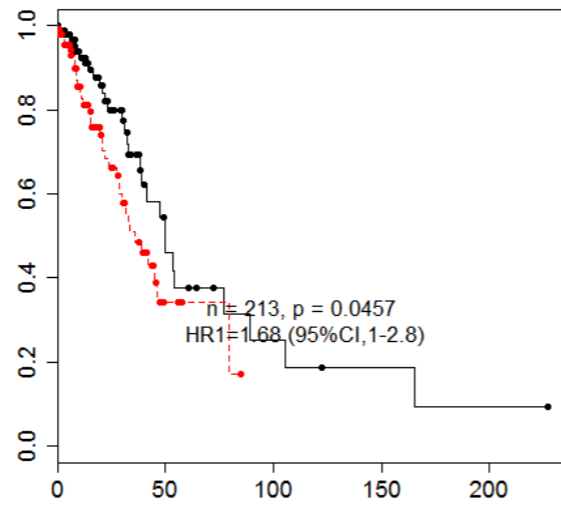
DDX5



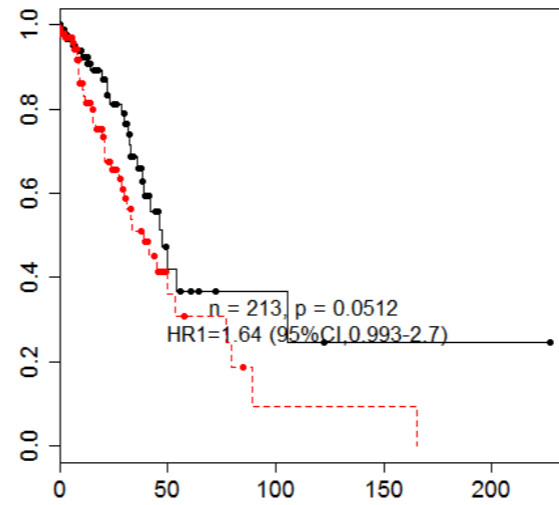
HSPA1A



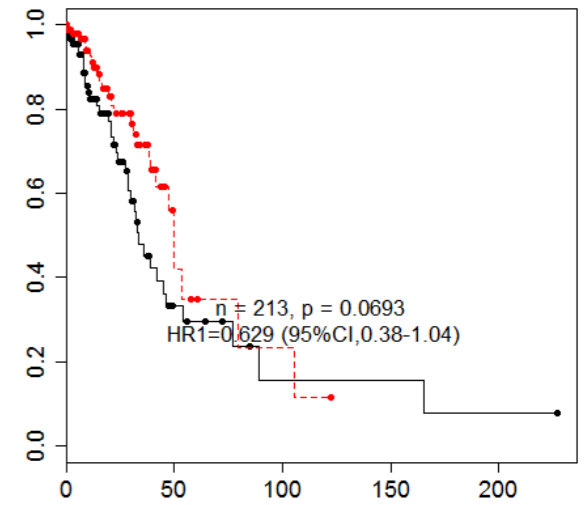
UPF3B



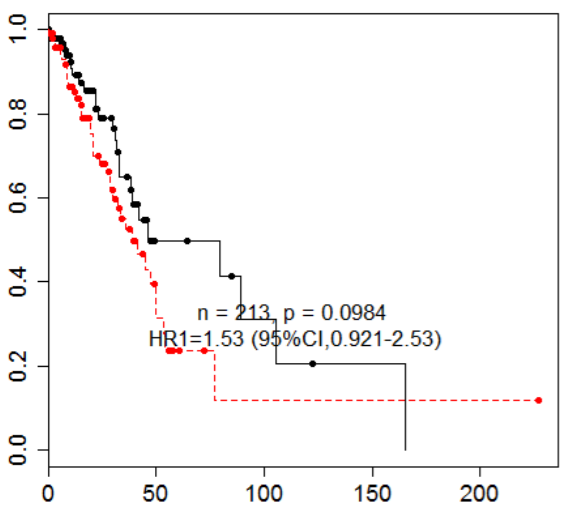
AURKA



YWHAZ

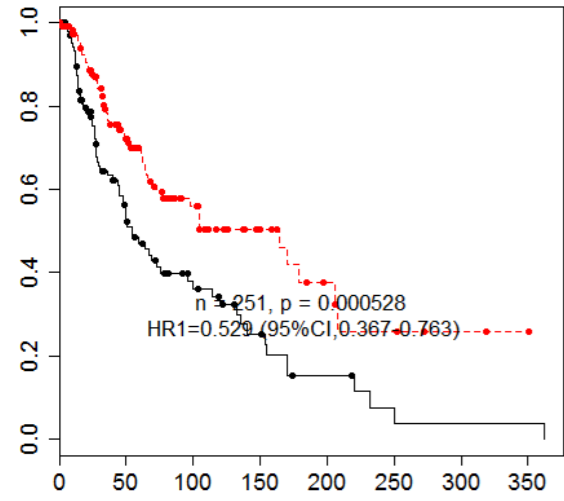


ACTC1

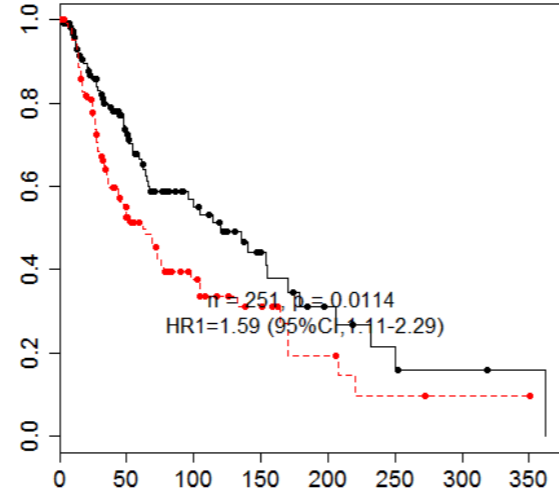


NPM1

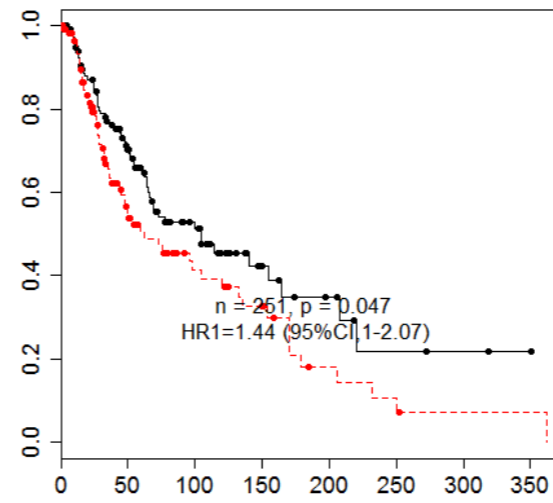
# h.TCGA SKCM



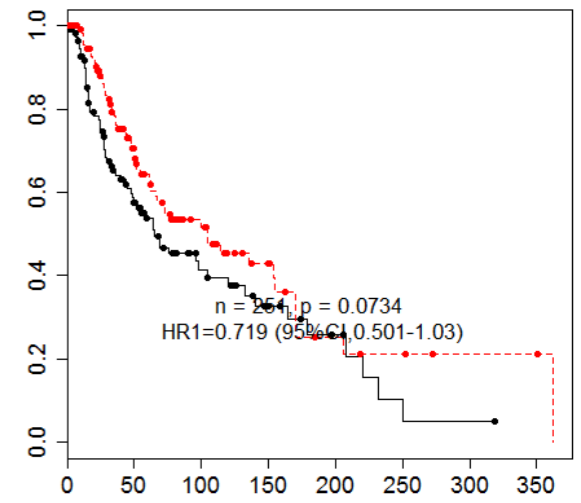
GCH1



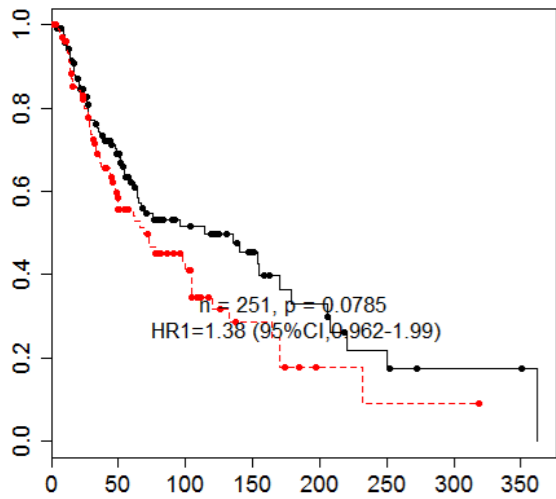
DSC1



AHSA1

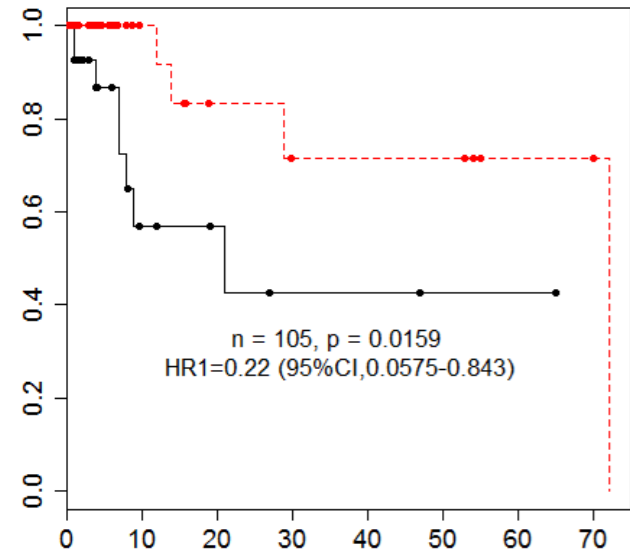


KRT18

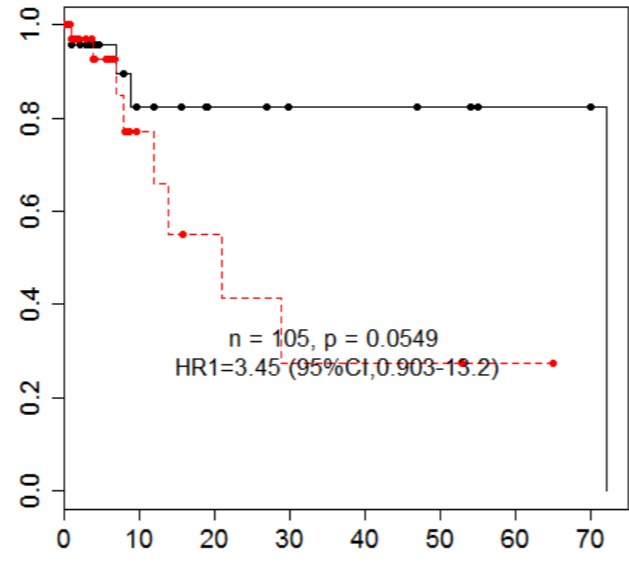


COL17A1

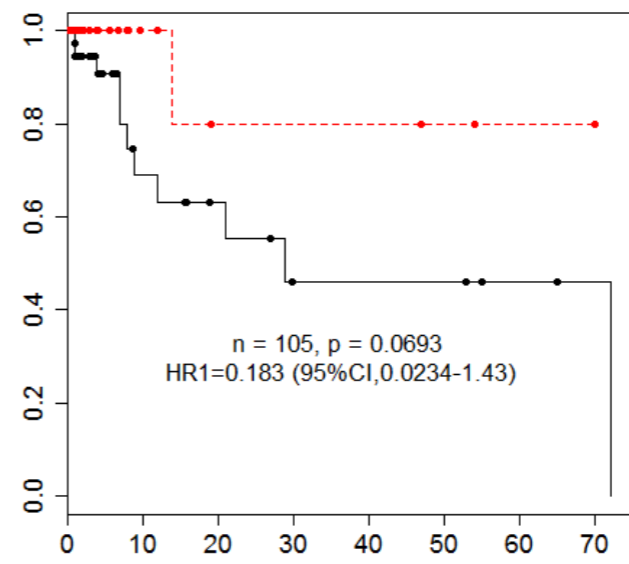
# i. TCGA STAD



HSF1

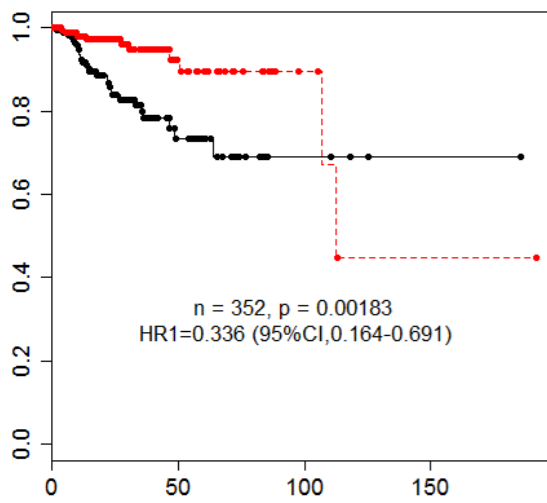


ASCC2

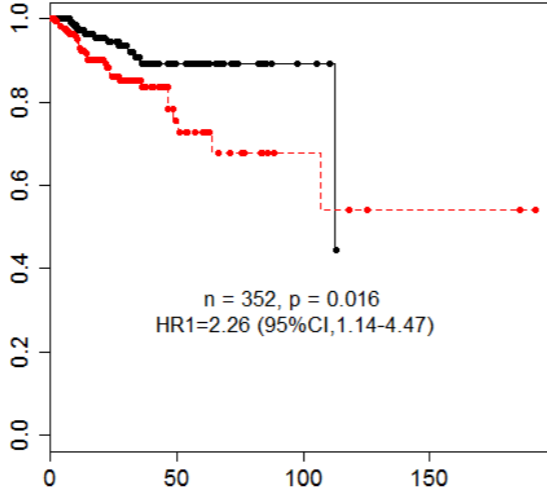


ACTA2

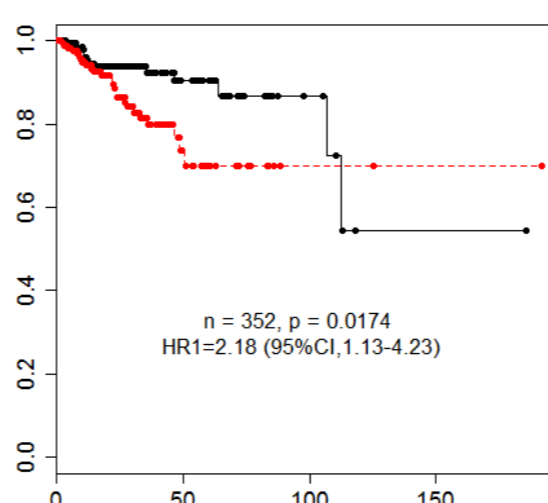
# j. TCGA UCEC



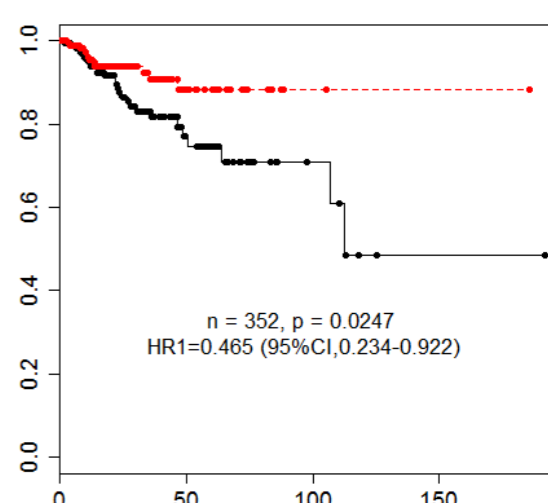
SPN



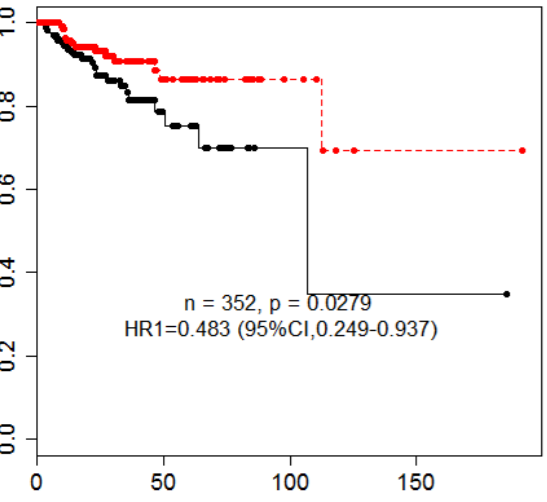
CDK6



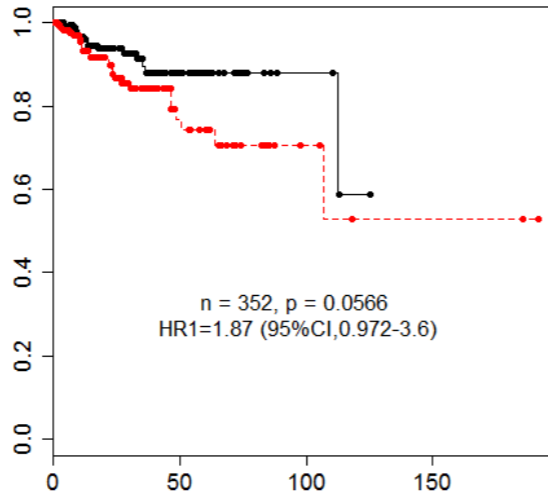
PIK3CA



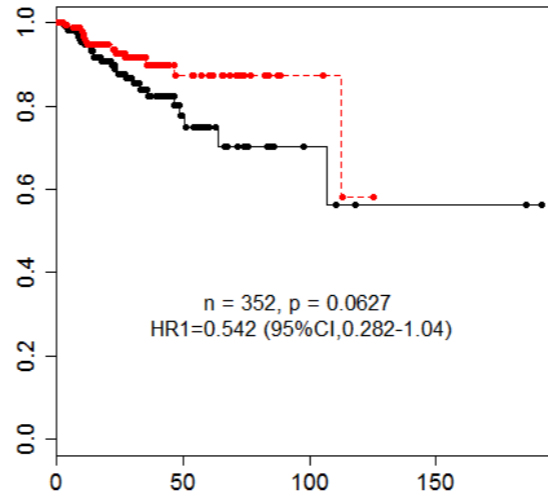
PIK3R1



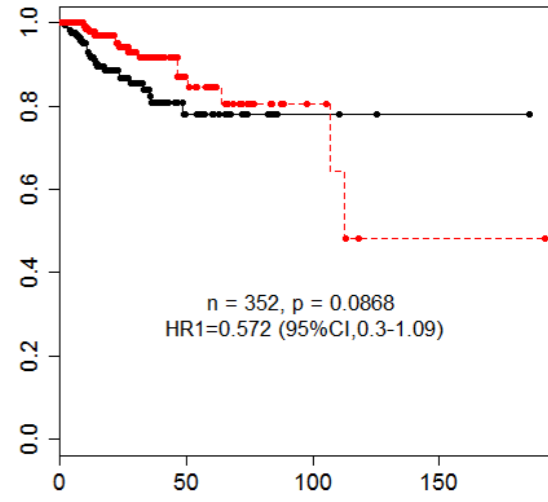
CDKN2D



FYN



RP1



PIK3CD