

Co-regulated gene expression of splicing factors as drivers of cancer progression

Esmee Koedoot¹, Marcel Smid², John A. Foekens², John W.M. Martens², Sylvia E. Le Dévédec¹, Bob van de Water^{1,3}

¹ Division of Drug Discovery and Safety, LACDR, Leiden University, the Netherlands

² Department of Medical Oncology and Cancer Genomics Netherlands, Erasmus MC Cancer Institute, Erasmus University Medical Center, Rotterdam, the Netherlands

³ To whom correspondence should be addressed

³Address of correspondence:

Bob van de Water
Division of Drug Discovery and Safety
Leiden Academic Centre for Drug Research
Leiden University
2300 RA Leiden
The Netherlands
Tel: +31-71-5276223
Fax: +31-71-5274277
b.water@lacdr.leidenuniv.nl

Running Title: Co-regulated expression of splicing factors in cancer

Supplementary Figure Legends

Supplementary Figure 1. Change in splicing factor RNA expression levels comparing normal, tumor tissue and metastatic tissue. **A)** Unsupervised clustering of log2 fold change in splicing factor expression levels comparing tumor to normal tissue (n = 116, Euclidean distance, complete linkage). **B)** Unsupervised clustering of log2 fold change in splicing factor expression levels comparing tumor to metastatic tissue (n = 7, Euclidean distance, complete linkage).

Supplementary Figure 2. Unsupervised clustering of PC of splicing factor RNA expression levels. **A)** Unsupervised clustering (Euclidean distance, complete linkage) of the PC of splicing factor RNA expression levels based on TCGA RNA sequencing data as in Fig. 1J, displaying all splicing factors. **B)** Unsupervised clustering (Euclidean distance, complete linkage) of the PC of splicing factor RNA expression levels based on BASIS RNA sequencing data as in Fig. 1K, displaying all splicing factors.

Supplemental Figure 3. TCGA cluster correlations. **A)** Hierarchical clustering (Euclidean distance, complete linkage) of the correlation of splicing factor expression levels in TCGA RNA sequencing data (red = high positive correlation, green = high negative correlation). Similar clustering as in Figure 1J, but now with all clusters numbered. **(B-G)** Hierarchical clustering (Euclidean distance, complete linkage) of the correlation of splicing factor expression levels of all cluster combinations shown in figure A.

Supplementary Figure 4. Pearson Correlation of RNA expression levels of cluster 1 and cluster 2 splicing factors in 867 primary breast tumors of untreated patients (MA-867).

Supplementary Figure 5. Distribution of cluster 1 and 2 splicing factors over different spliceosome sub-complexes. **A)** Cluster 1 splicing factor distribution over core and non-core categories. **B)** Same as in A but now for Cluster 2 splicing factors. **C)** Cluster 1 splicing factor

distribution over different spliceosomal sub-complexes. **D)** Same as in C, but now for Cluster 2 splicing factors.

Supplementary Figure 6. High cluster 2 splicing factor expression levels are related to a more aggressive breast tumor phenotype. **A)** Unsupervised clustering SF RNA expression in primary breast tumors of BASIS RNA sequencing data, annotated with known driver gene mutations. SF levels were first log₂ normalized. Next median patient levels were equalized to 0 per SF. **B)** Log₂ fold change in expression of cluster 1 and cluster 2 splicing factors comparing ER negative to ER positive primary breast tumors using the TCGA RNA sequencing dataset **C)** Hierarchical clustering of SF expression levels in primary breast tumors of TCGA RNA sequencing data. SF levels were first log₂ normalized, followed by equalizing median patient levels to 0 per SF. **D)** Cluster 1 and cluster 2 splicing factor expression levels in primary breast tumors with different pleomorphism scores. Per SF, fold changes were calculated compared to score 1. **E)** Cluster 1 and cluster 2 splicing factor expression in PAM50 breast cancer subtypes. Groups are compared using a student's t-test. * P<0.05, ** P<0.01, ***P<0.001.

Supplementary Figure 7. Log₂ difference in cluster 1 and 2 expression levels comparing **A)** primary tumor tissue with normal tissue and **B)** primary tumor tissue with metastatic tissue. Dots represent different patients.

Supplementary Figure 8. A) Association of the expression of cluster 1 and cluster 2 splicing factors with breast cancer overall survival. Per patient, mean expression of all factors within one cluster was calculated. Based on this mean expression, the patient cohort was median-split in low and high expression of cluster 1 or 2 splicing factors and survival curves for overall survival, relapse-free survival and metastasis-free survival were generated. Breast cancer survival curves of cluster 1 and cluster 2 splicing factors in estrogen receptor (ER) negative **(B)** and ER positive **(C)** subtypes. Per patient, mean expression of all factors within

one cluster was calculated. Based on this mean expression, the patient cohort was median-split in low and high expression of cluster 1 or 2 splicing factors and survival curves for overall survival, relapse-free survival and metastasis-free survival were generated. **D)** Hazard ratio (HR) for distant metastasis formation for cluster 1 and cluster 2 splicing factors in all (left) and estrogen receptor (ER) positive (right) tumors in the MA-867 dataset.

Supplementary Figure 9. A) Examples of expression levels of multiple isoforms of one gene in human breast cancer patient primary tumor RNA sequencing data. **B)** Density plot of difference in isoform length comparing cluster 1 and cluster 2 isoforms of the same gene. The length of the isoform in cluster 2 was subtracted from the length of the same gene isoform in cluster 1. Isoform clusters are based on Figure 4B. **C)** PCs of SMAD3 isoforms with Suppressor-SF DHX9 and Enhancer-SF FAM50A. **D)** Same as in C, but for NFKB2 isoforms. **E)** Same as in C, but for MCL1 isoforms. **F)** Same as in C, but for HNRNPA1 isoforms.

Supplementary Figure 10. Zoom in of isoform splicing factor clustering in Fig. 4B.

Supplementary Figure 11. Unsupervised clustering (Euclidean distance, complete linkage) of PC of Suppressor- and Enhancer-SFs with genes known to be involved in different stages or regulatory pathways of cell cycle. Red = highly positively correlated, green = highly negatively correlated.

Supplementary Figure 12. Transcription factors ATF1 and CREB1 can bind promoter regions of Suppressor-SFs but not Enhancer-SFs. **A)** Transcription binding site enrichment for splicing factors in promoter regions of Enhancer- and Suppressor-SFs. **B)** Expression levels of ATF1, CREB1 and CREM in High-Suppressor Low-Enhancer and Low-Suppressor High-Enhancer primary breast tumors. **C)** Unsupervised clustering (Euclidean distance,

complete linkage) of PCs of CREB1, ATF1 and CREM with Suppressor- and Enhancer-SFs. Red = highly positively correlated, green = highly negatively correlated.

Supplementary Figure 13. Hierarchical clustering (Euclidean distance, complete linkage) of PCs of Suppressor- and Enhancer-SFs in different cancer types. Red = highly positively correlated, green = highly negatively correlated.

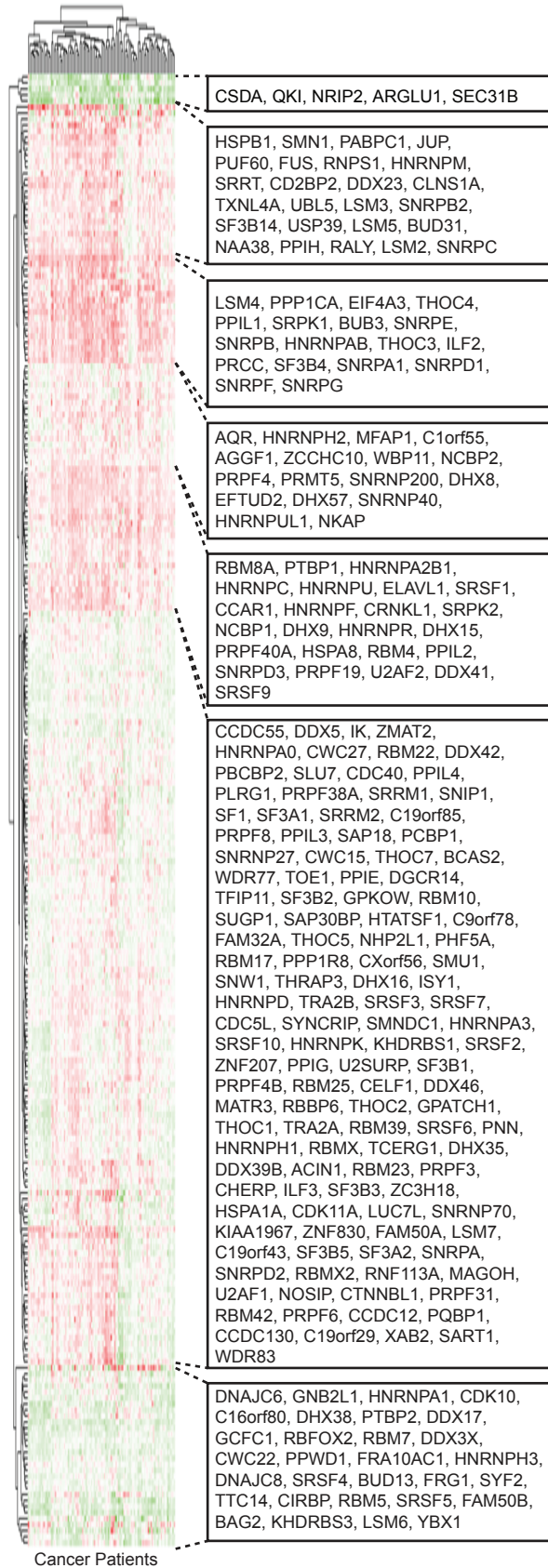
Supplementary Figure 14. Hierarchical clustering (Euclidean distance, complete linkage) of Suppressor- and Enhancer-SF PCs to genes involved in mitochondrial translation, cell cycle, M phase and respiratory electron transport in breast cancer, lung cancer, pancreas cancer and prostate cancer. Red = highly positively correlated, green = highly negatively correlated.

Supplementary Figure 15. Enhancer- and Suppressor-SF expression levels related to overall survival, post-progression survival and progression-free survival in ovarian cancer. Per ovarian cancer patient, mean expression of all Suppressor- and Enhancer-SFs was calculated. Based on these expression levels, the patient cohort was median-split and overall and post-progression survival plots were generated²⁵.

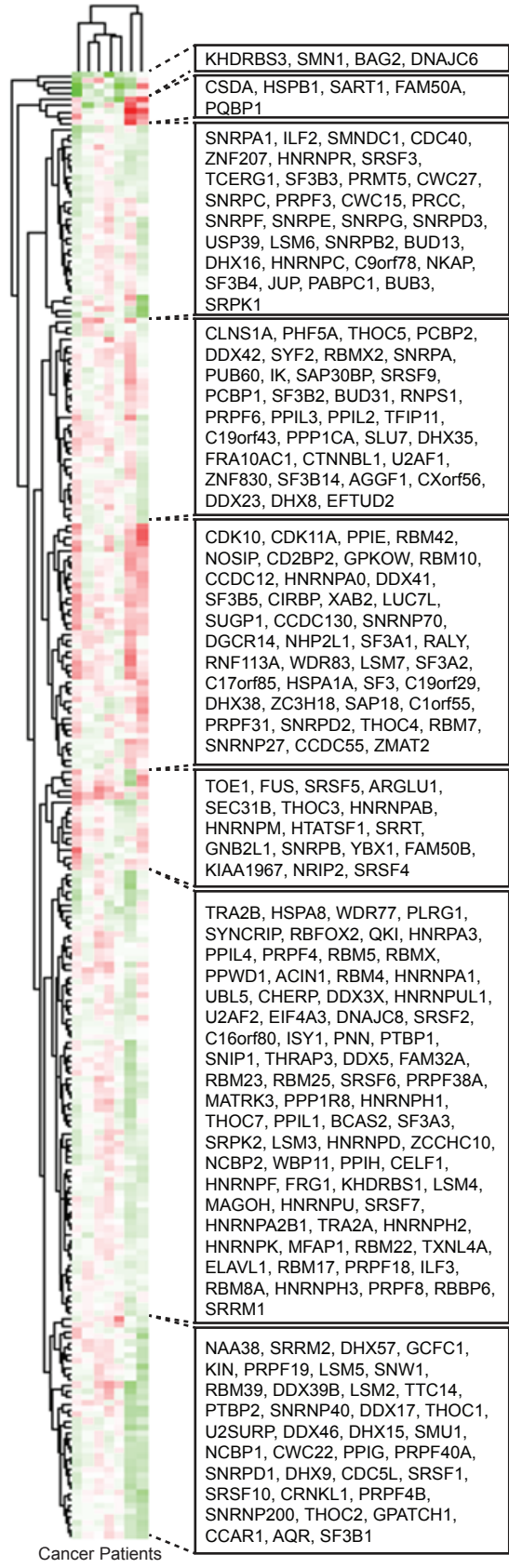
Supplementary Figure 16. RNA-protein spearman correlations of splicing factors. RNA sequencing data was derived from TCGA, protein expression data was derived from Mertins et al³⁴.

Supplementary Figure 17. Correlation of log2 fold change of splicing factor expression comparing normal to tumor tissue with 1) the non-GLM method using normalized counts data provided by the TCGA portal and 2) using raw sequencing counts and a GLM method.

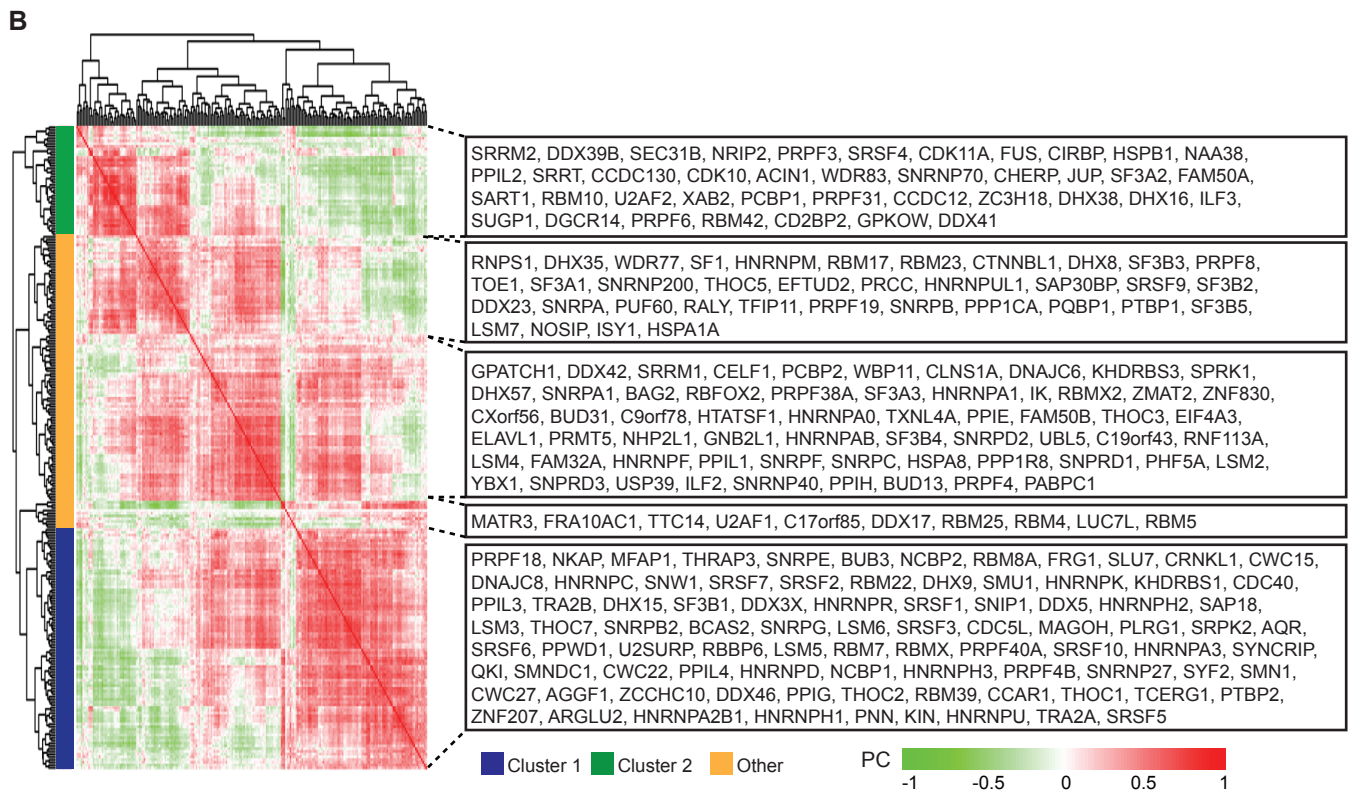
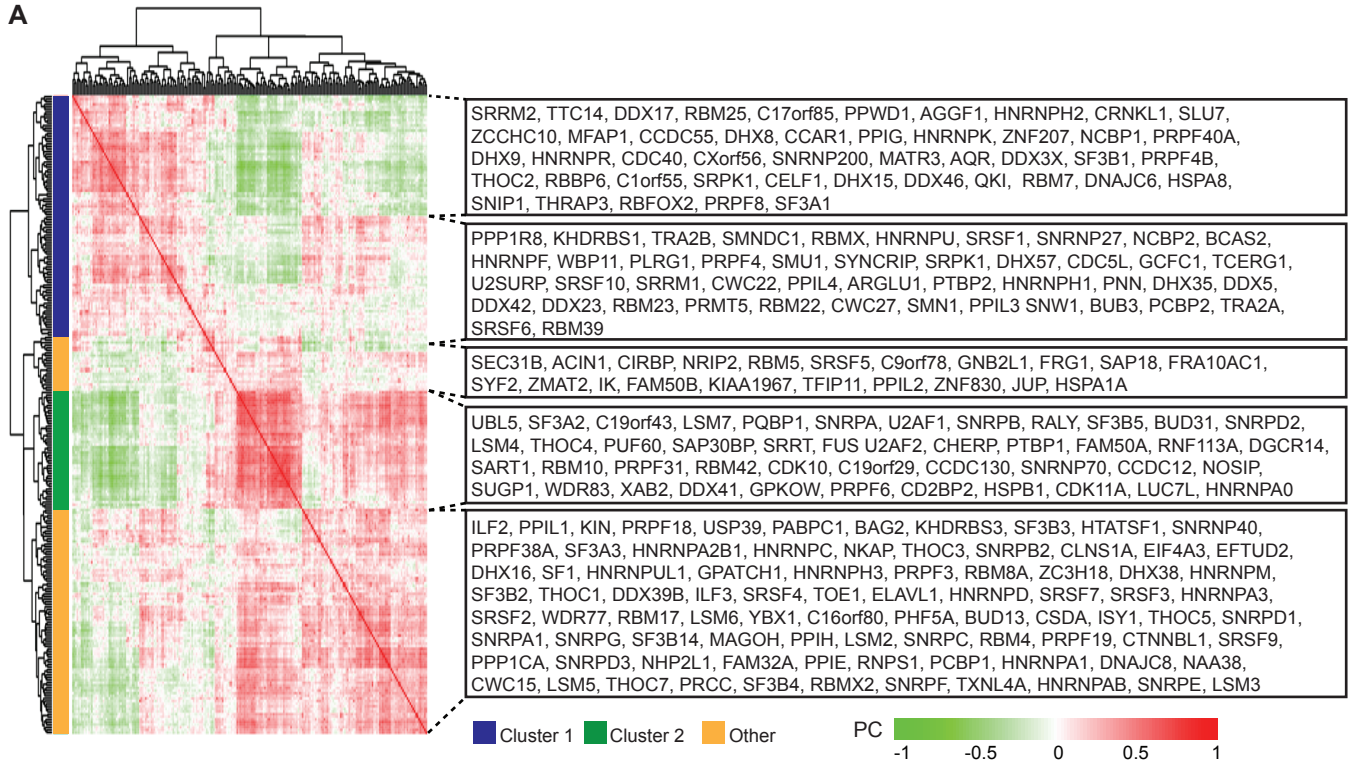
A



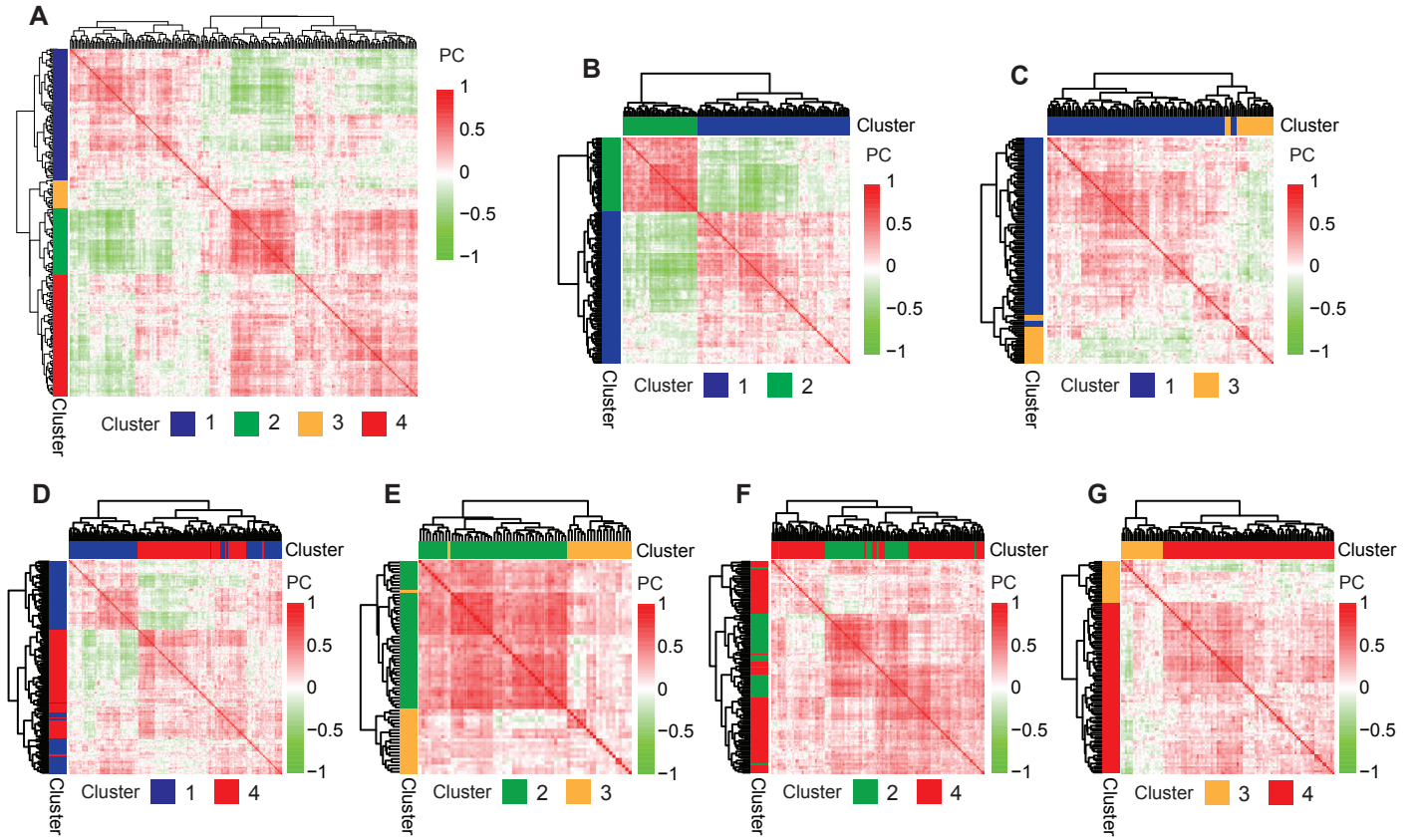
B



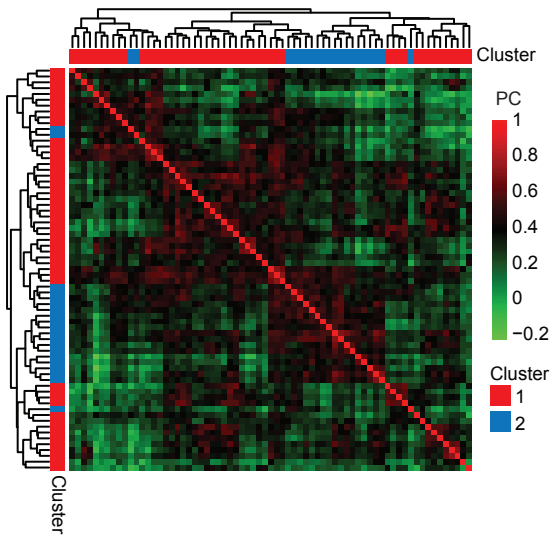
Supplementary Figure 2



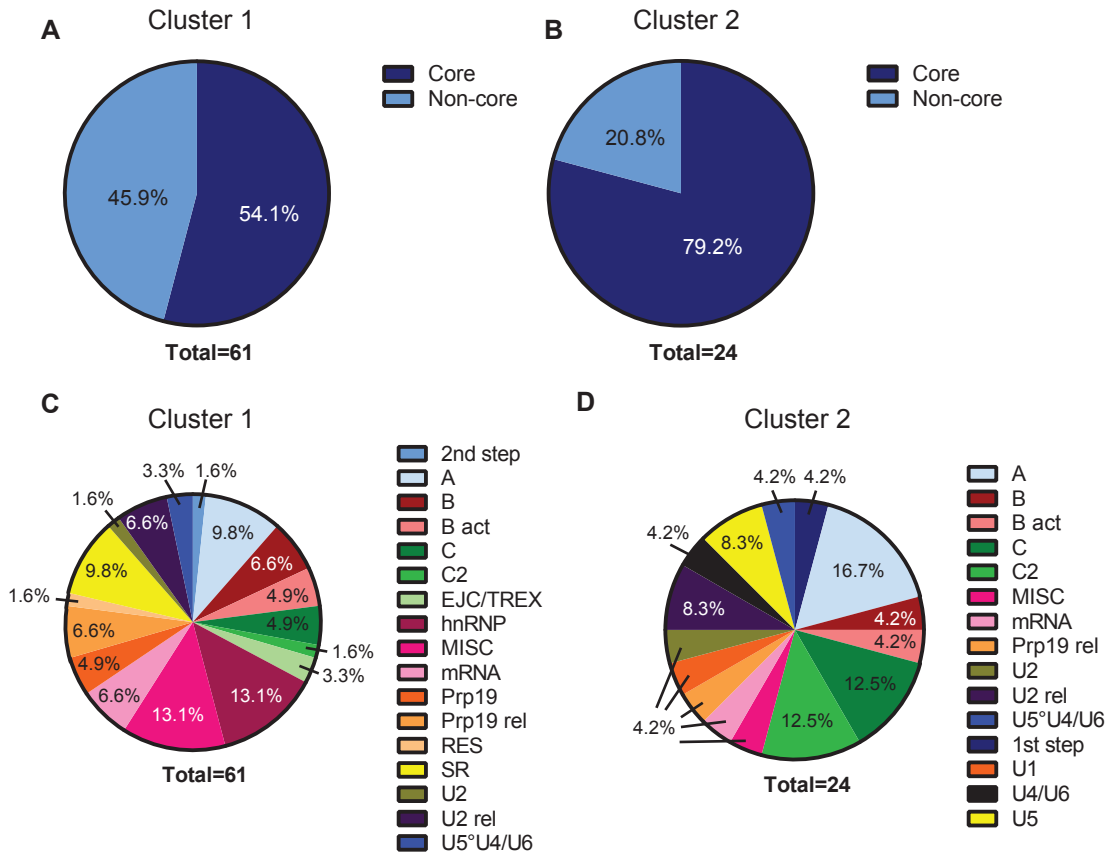
Supplementary Figure 3



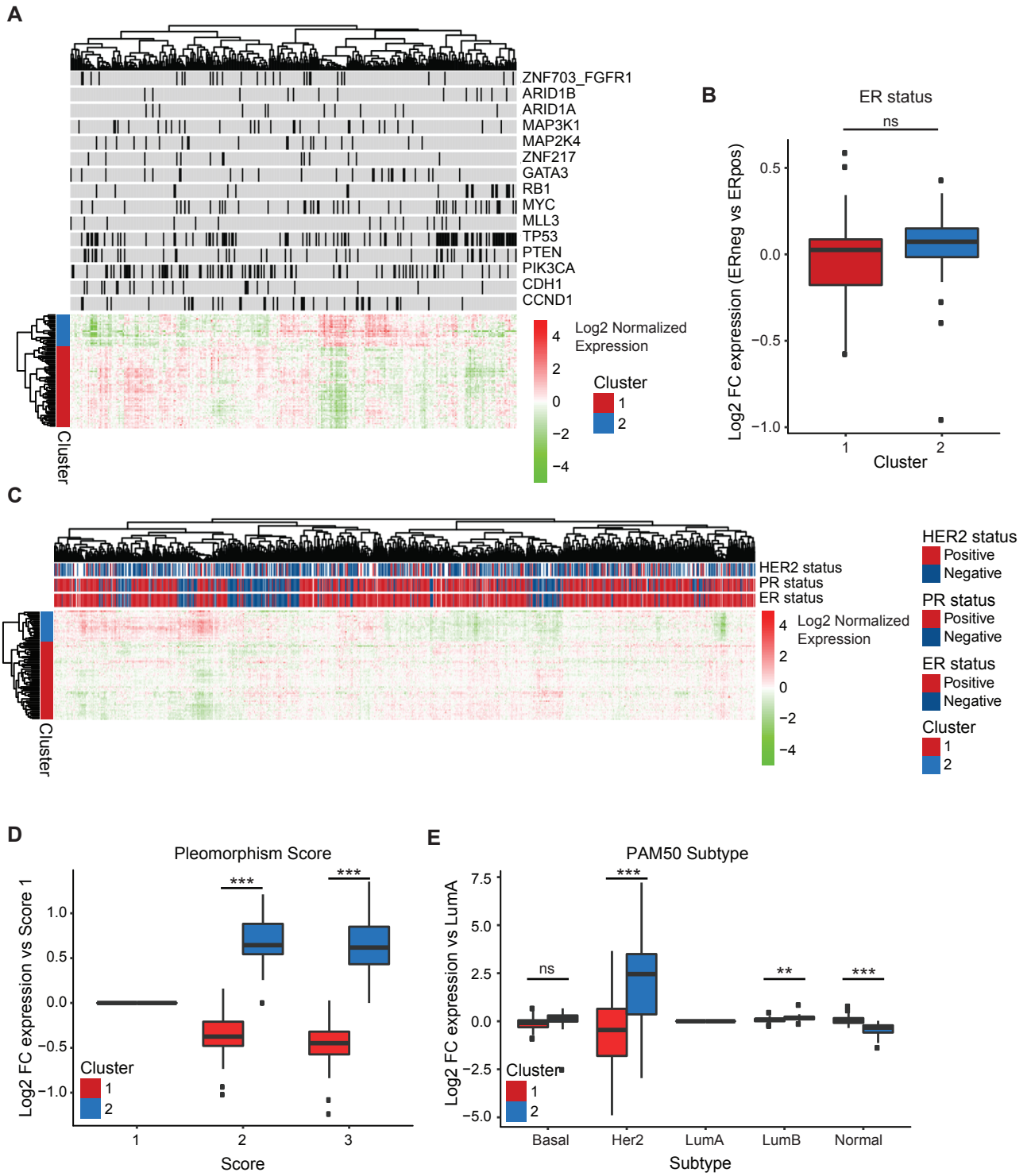
Supplementary Figure 4



Supplementary Figure 5

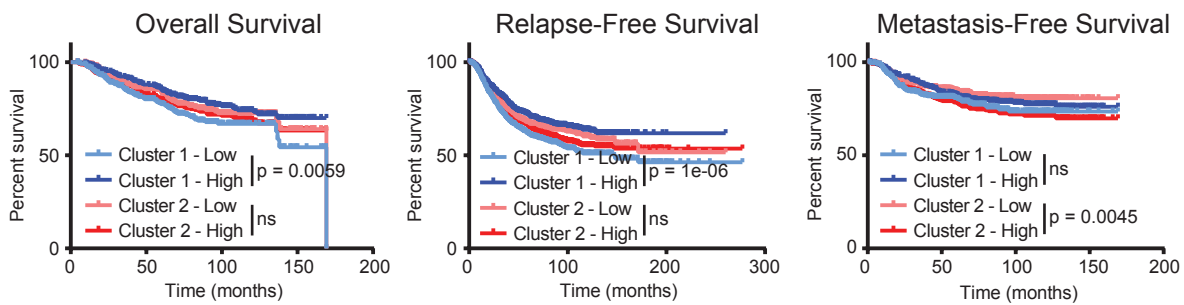


Supplementary Figure 6

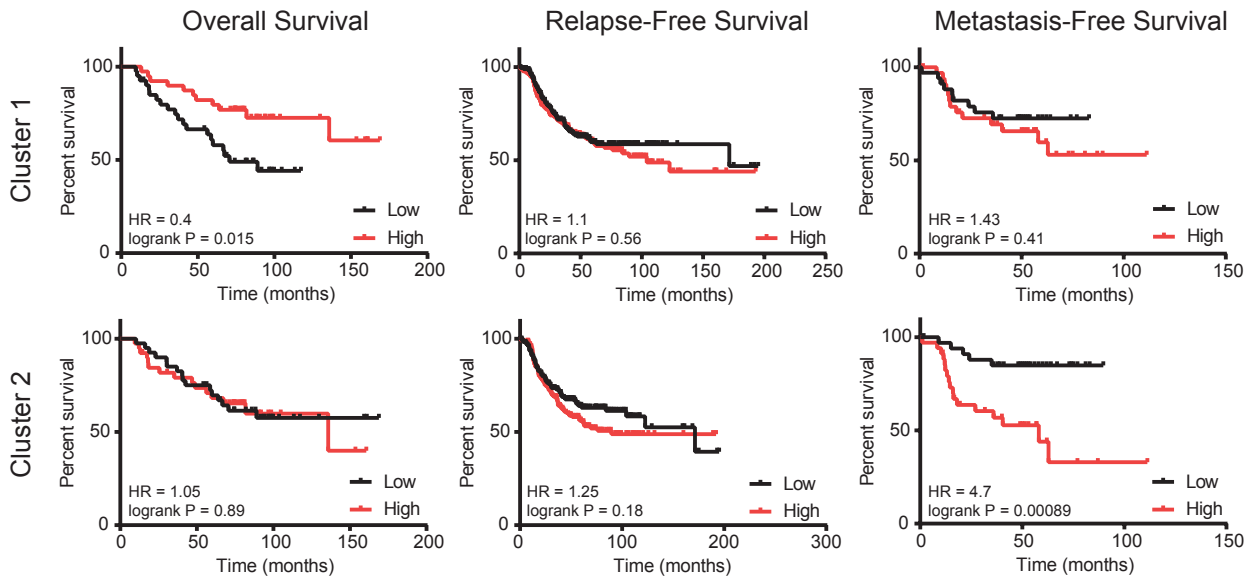


Supplementary Figure 8

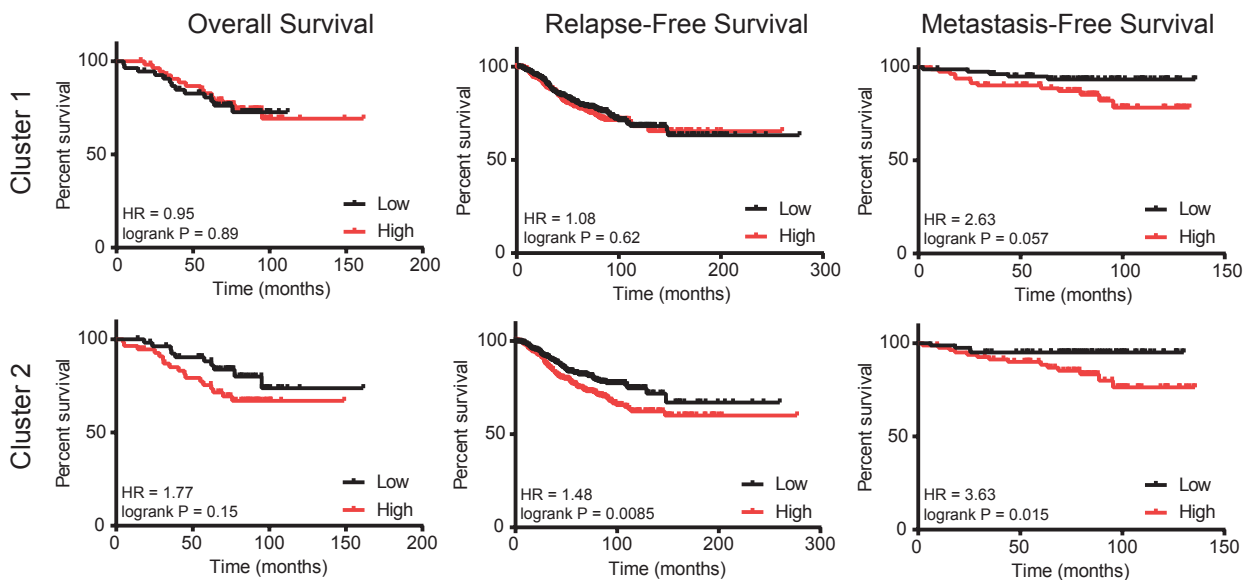
A



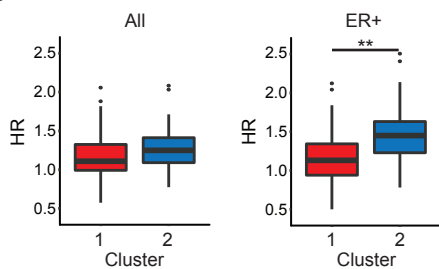
B



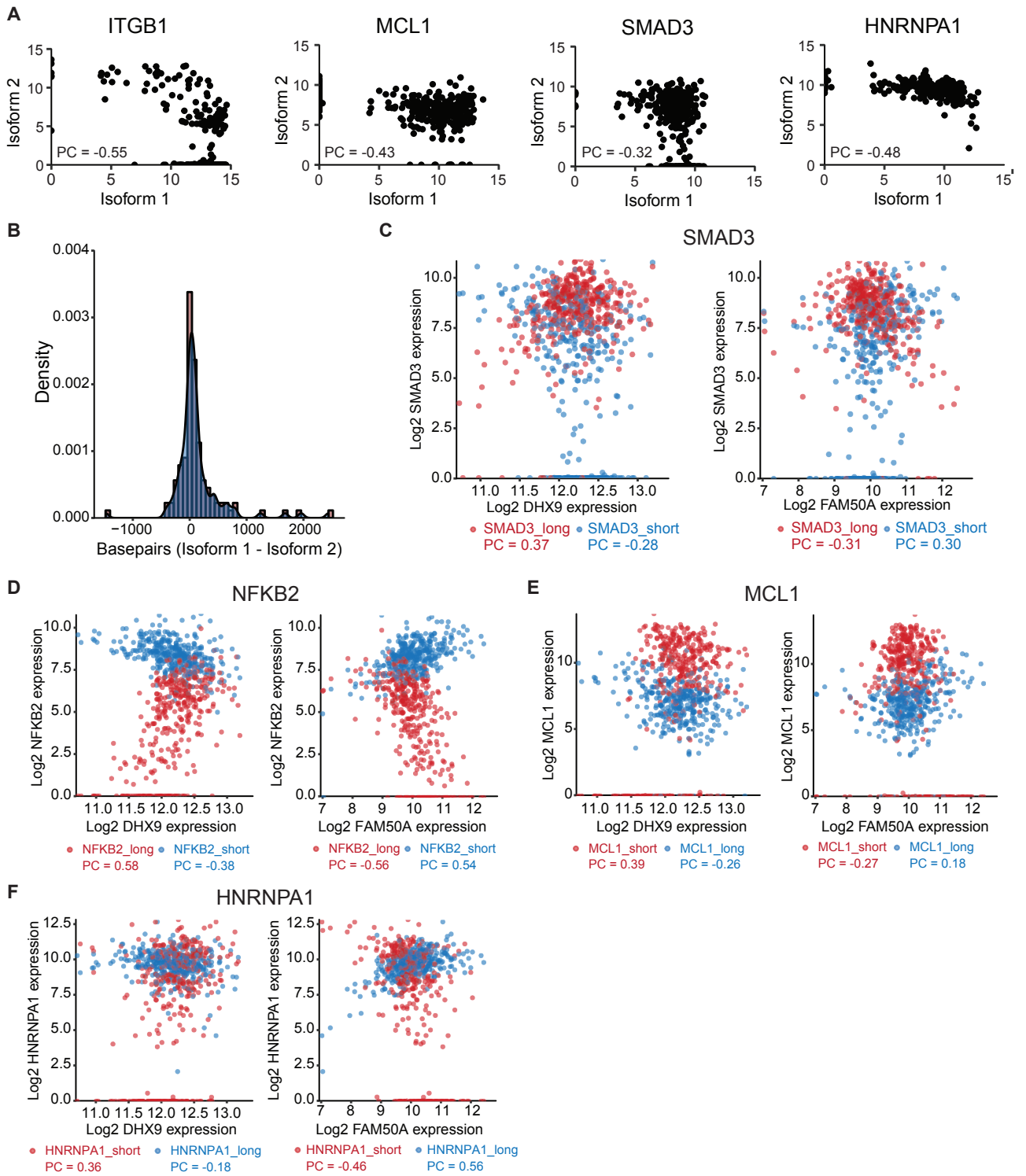
C



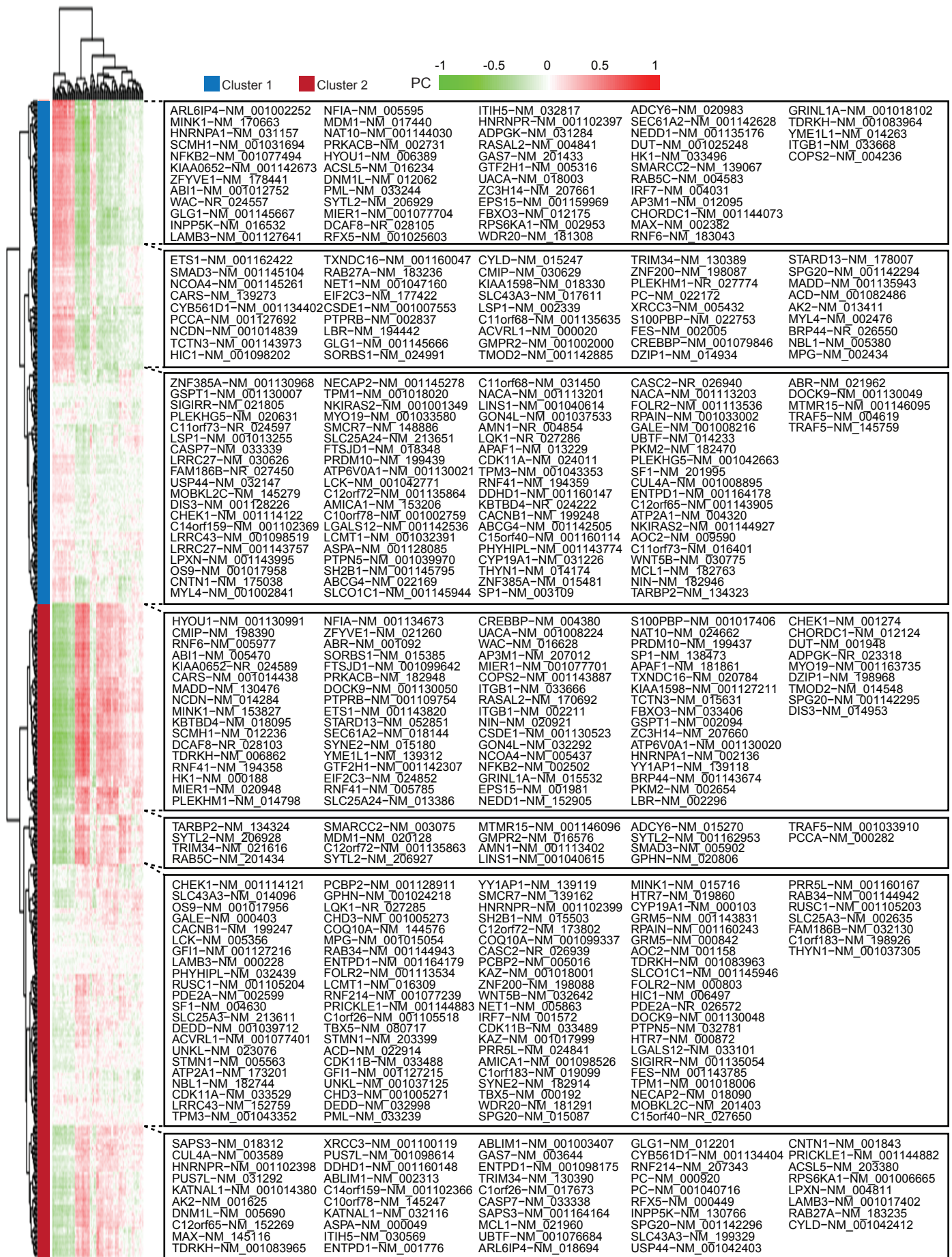
D



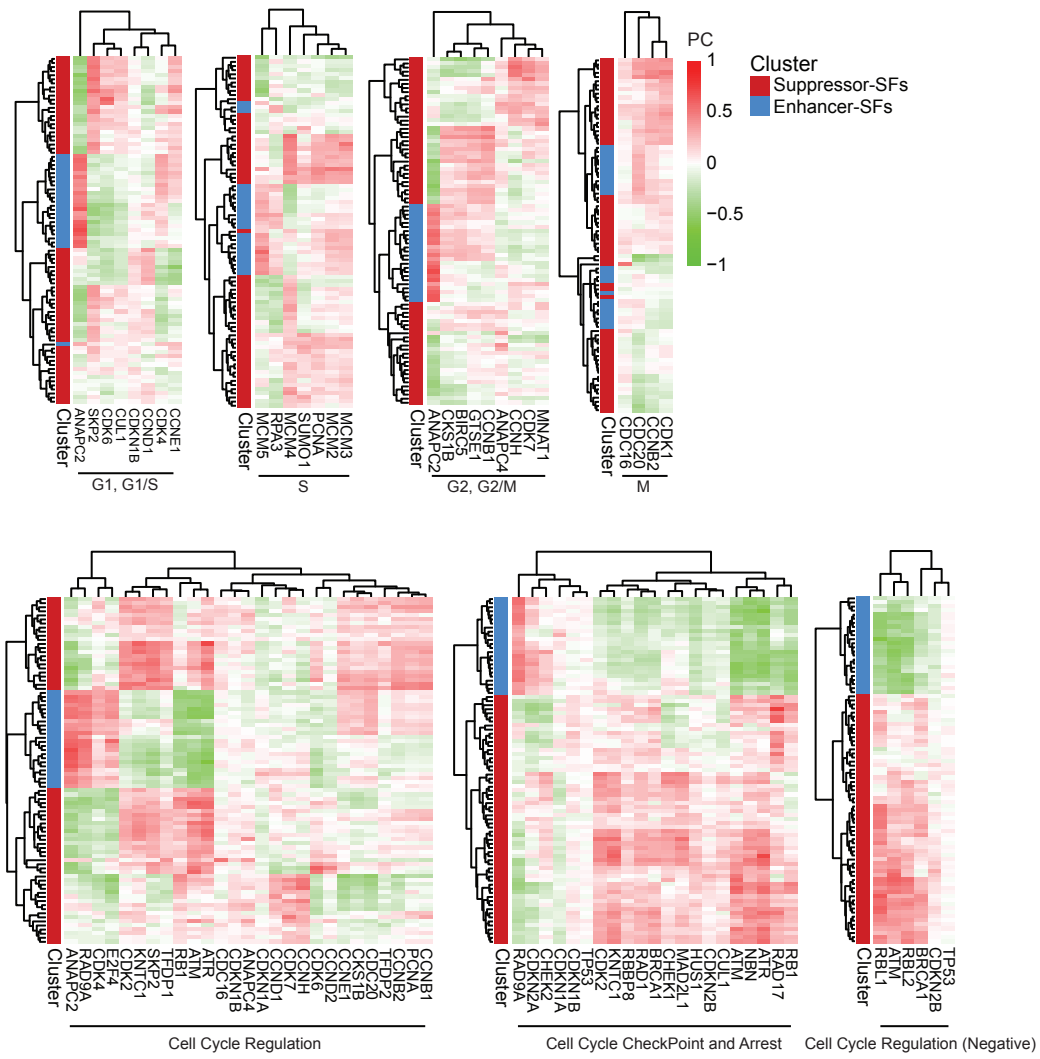
Supplementary Figure 9



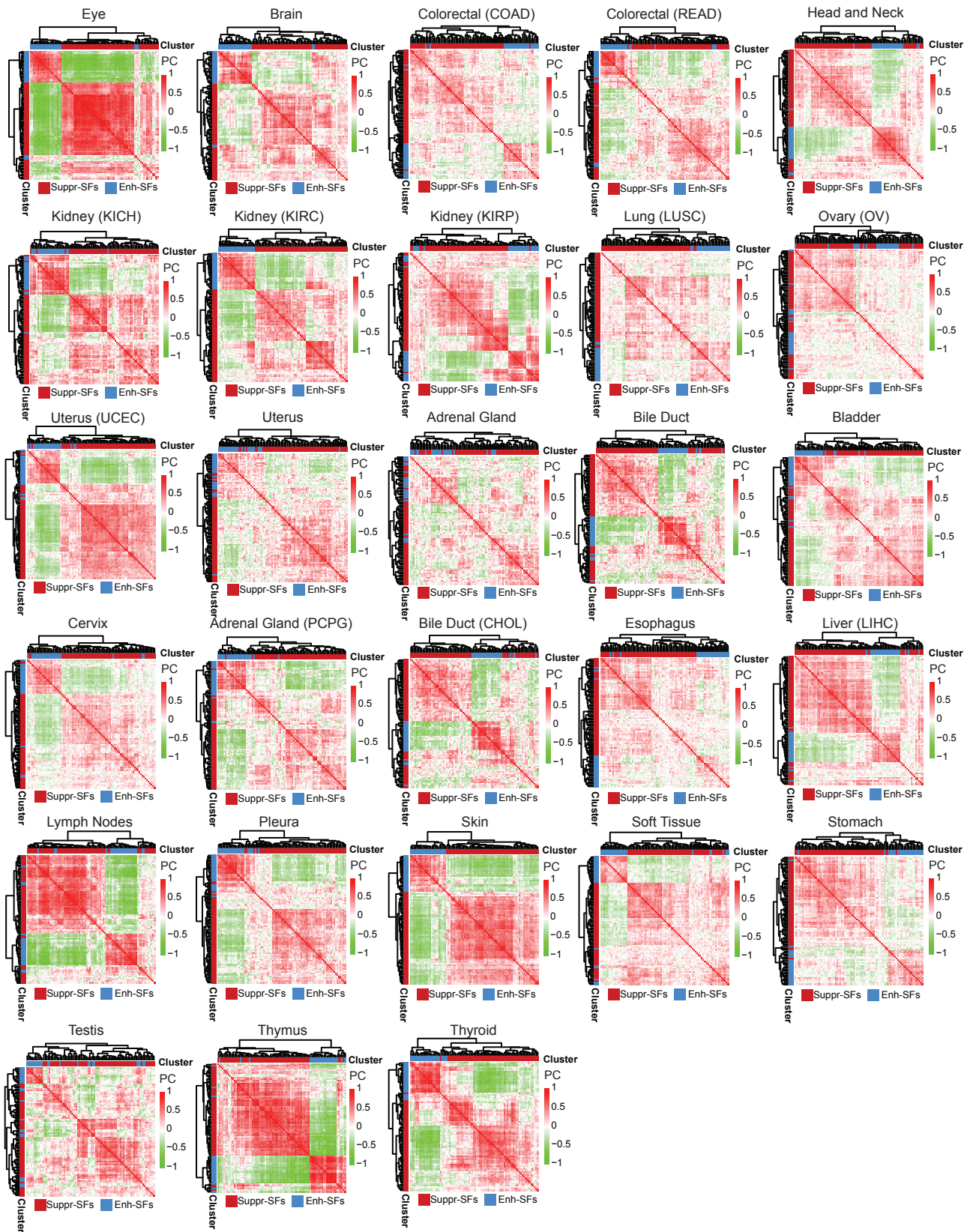
Supplementary Figure 10



Supplementary Figure 11

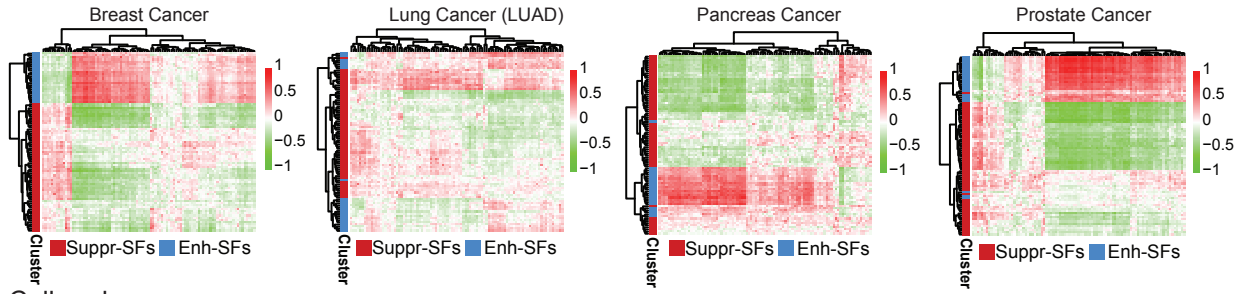


Supplementary Figure 13

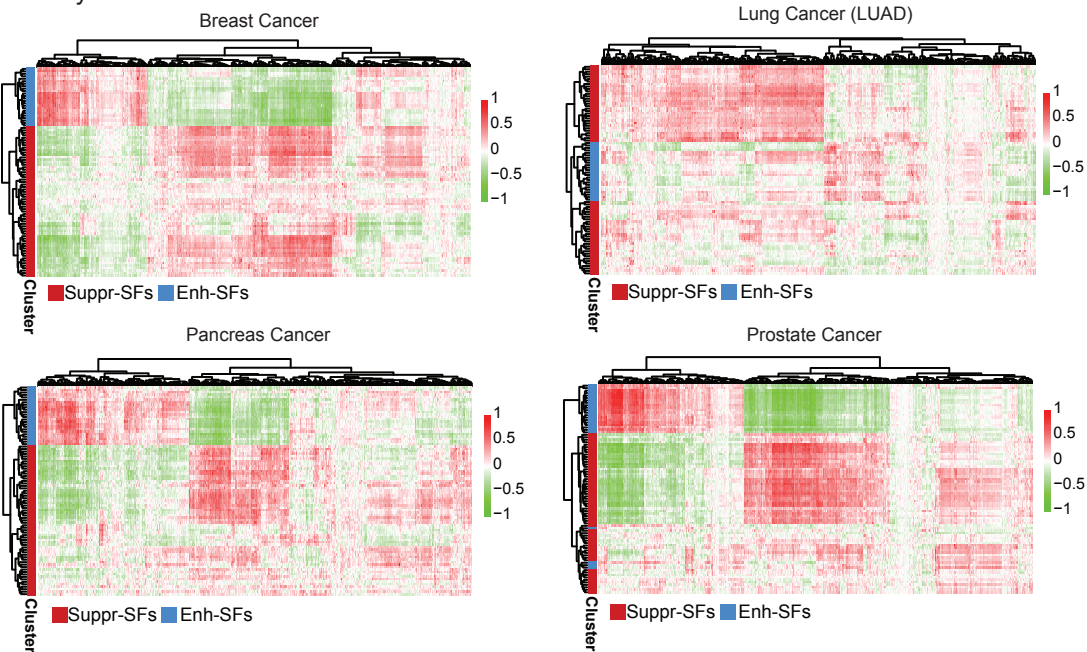


Supplementary Figure 14

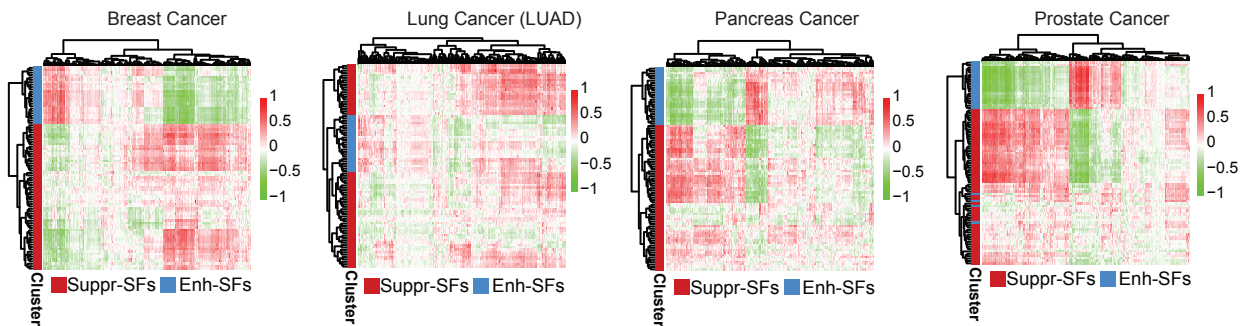
Mitochondrial translation



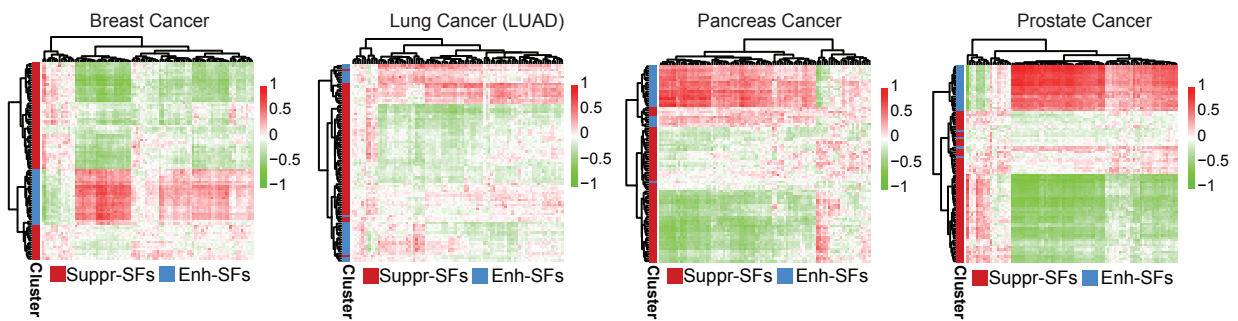
Cell cycle



M phase

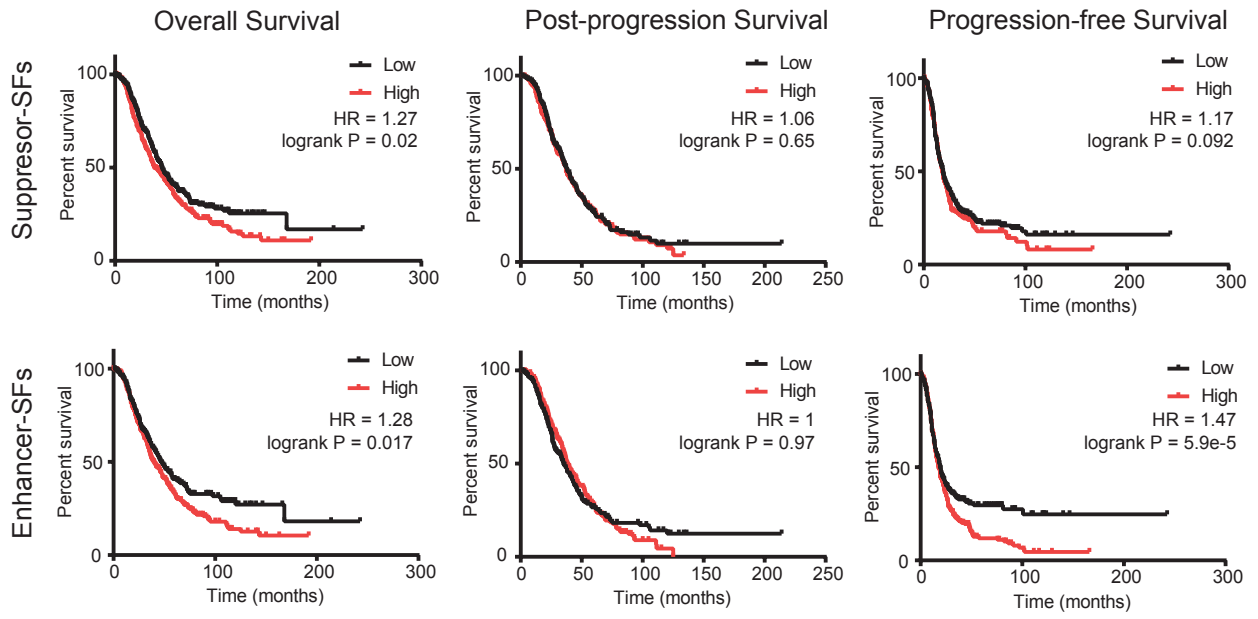


Respiratory electron transport

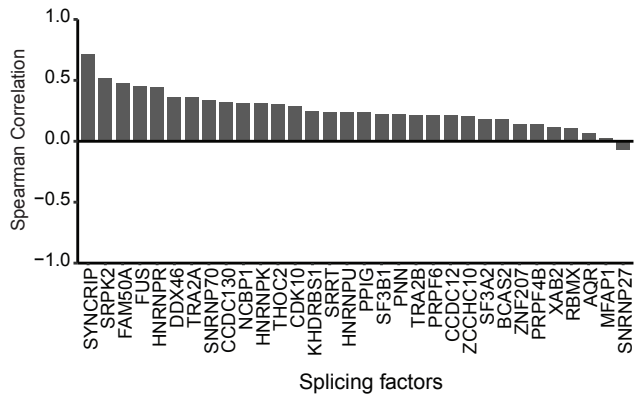


Supplementary Figure 15

Ovarian Cancer



Supplementary Figure 16



Supplementary Figure 17

