

**Supplementary Information for:**

**A tumorigenic index for quantitative analysis of liver cancer initiation and progression**

Gaowei Wang<sup>a</sup>, Xiaolin Luo<sup>a</sup>, Yan Liang<sup>a</sup>, Kota Kaneko<sup>a</sup>, Hairi Li<sup>b</sup>, Xiang-Dong Fu<sup>b</sup>, and Gen-Sheng Feng<sup>a,1</sup>

<sup>a</sup> Department of Pathology, Division of Biological Sciences, Moores Cancer Center, University of California San Diego, La Jolla, CA 92093, USA

<sup>b</sup> Department of Cellular and Molecular Medicine and Institute for Genomic Medicine, University of California, San Diego, La Jolla, CA, USA.

<sup>1</sup> To whom correspondence may be addressed. email: [gfeng@ucsd.edu](mailto:gfeng@ucsd.edu)

**This PDF file includes:**

Supplementary text  
Figures S1 to S9  
SI References

**Other supplementary materials for this manuscript include the following:**

Dataset S1 to S11

## Supplementary Information Text

### Materials and Methods

**Animal protocols.** Hepatocyte-specific *Shp2* KO mice (SKO), *Pten* KO mice (PKO) and *Shp2* and *Pten* double-knockout (DKO) mice were generated and characterized as described previously<sup>1-3</sup>. All animal experimental protocols (S09108) have been approved by the Institutional Animal Care and Use Committee (IACUC) of the University of California, San Diego, following NIH guidelines.

**RNA sequencing and data analysis.** Total RNAs were extracted from liver tissues using QIAGEN RNeasy columns, and RNA-sequencing (RNA-seq) was performed using the multiplex analysis of polyA-linked sequence and the Illumina HiSeq2000 machine. Raw reads generated by RNA-seq experiments were mapped to the mm9 mouse reference genome using Star (2.3.0). The expression level of each gene under different conditions was obtained using cuffdiff. Quality control of RNA-seq data was performed by calculating Pearson's correlations between samples and analyzing expression level changes of *Shp2/Ptpn11* and *Pten* in WT and mutant livers. We retained 16230 genes that had a non-zero expression level in at least 12 of all the samples for the following analyses. Differentially expressed genes between mutant and WT livers at different time points were identified using statistical test (t-test, q values<0.05). In order to identify when and how these differentially expressed genes were changed during liver tumorigenesis in SKO, PKO and DKO mice, we selected genes that have differential expression in at least two mutant livers. We also performed targeted analysis of the expression of ligands/receptors and epigenetic regulators during HCC development. Ligand and receptor genes were collected from Database of Ligand-Receptor Partners (DLRP), and epigenetic genes were collected from EpiFactors database. Changes in gene expression levels were visualized using heatmaps, which were generated using the heatmap package in R. Based on differentially expressed genes, top-enriched biological processes and pathways of mutant livers at different ages were identified by performing gene set enrichment analysis (GSEA). In order to identify when and how these biological processes and pathways were changed during the course of tumorigenesis in PKO

liver, we selected biological processes and pathways that have significant changes in at least two time point at PKO liver.

### ***Significantly changed TF clusters between WT livers and tumors***

Transcription factors (TF) and their downstream target genes often co-express to control specific cell activities. We downloaded and integrated gene-gene correlative relationships from Cellnet database (<http://cellnet.hms.harvard.edu/downloads/>) and geneFriend database (<http://www.genefriends.org/RNAseq/about/>). Correlation threshold was set as 0.5, genes that have high positive correlation with each TF were identified and defined as a TF cluster. In total, 1568 TF clusters have been collected and defined. We further assessed the ability of these genes in TF cluster to capture key features of a given TF. We tested TF clusters by comparing their well-documented functions and the top-enriched biological processes of the co-expressed genes. For example, consistent to the known function of E2F1 in control of cell cycle progression, gene ontology analysis of the E2F1 cluster members identified cell cycle as the top-enriched biological process. In a similar way, we tested other TFs by comparing their well-documented functions and the top-enriched biological processes of the co-expressed genes (Table S6).

R package GSA has improved GSEA based on Kolmogorov-Smirnov-like statistics to identify differentially expressed genes from the transcriptomes. Using the gene list of each TF cluster and gene expression profiles of WT livers and tumors as input, GSA evaluated whether each TF cluster was differentially expressed between the two groups of samples by calculating a set score and p value for each TF cluster. We set the set score  $>0.8$  and p-value  $<0.05$  to identify significantly up-regulated TF clusters and set score  $<-0.8$  and p-value  $<0.05$  to identify significantly down-regulated TF clusters. In total, 36 up-regulated and 25 down-regulated TF clusters from WT livers to tumors were identified. Regression analysis method LASSO (least absolute shrinkage and selection operator) was performed using Glmnet package in R, and Random Forests was performed using Random Forests package in R.

### ***Human data and prognosis analysis***

Gene expression profiles, tumor stages, survival information, gender, genetic alterations, age, gender information of samples are available in 371 HCC samples in the TCGA dataset.

Transcriptomes and clinical information of human HCC patients in TCGA were downloaded from The Cancer Genome Atlas (TCGA) website (<https://portal.gdc.cancer.gov/>).

Gene expression profiles, tumor stages, survival information, gender, age, cirrhosis status, one or more nodules, tumor sizes and AFP levels of samples are available in 247 HCC samples in the GSE14520 dataset. Transcriptome and clinical information of HCC patients in GSE14520 dataset were downloaded from National Center for Biotechnology Information (NCBI) using GEO accession ID. Gene expression profiles, tumor stages, survival information, gender, age, cirrhosis status, one or more nodules, tumor sizes and AFP levels of samples are available for 100 HCC samples in GSE16757 dataset. Transcriptomes and clinical information of human HCC patients in GSE16757 dataset was kindly provided by J. Lee (MD Anderson). Kaplan-Meier plots and log-rank test were used to determine the significant difference between survival curves.

### ***Inferring correlations between TF clusters***

First, we calculated Pearson's correlations and p-values between TFs using expression levels of TFs in all samples. As each TF cluster included a list of genes, we further calculated correlations between genes in each TF cluster. The medium of these correlations was used to define the relationship between TF clusters. We set a threshold to obtain significant negative correlation (Pearson's correlations <-0.5 and p-value <0.05, and medium correlation <0) and positive correlation (Pearson's correlations >0.5 and p-value <0.05 and medium correlation >0) between TF clusters. Correlations between TF clusters were visualized using network.

### ***Quantitative description and analysis of a coarse-grained model***

Quantitative description of the coarse-graining model in Fig. 4C consists of a set of differential equations.

$$\tau_{TF\_WT} \frac{d[TF\_WT]}{dt} = B_{TF\_WT} + V_{TF\_WT} * \frac{[TF\_WT]^n}{k_{TF\_WT} + [TF\_WT]^n + [TF\_HCC]^n} - [TF\_WT]$$

$$\tau_{TF\_HCC} \frac{d[TF\_HCC]}{dt} = B_{TF\_HCC} + V_{TF\_HCC} * \frac{[TF\_HCC]^n}{k_{TF\_HCC} + [TF\_WT]^n + [TF\_HCC]^n} - [TF\_HCC]$$

Where  $[TF\_WT]$  and  $[TF\_HCC]$  denote expression levels of activated TF clusters in WT liver and HCC, respectively (Fig. 2C). The right-hand side of the equation is given by the sum of basal production rate, integrated nonlinear production rate and linear natural degradation rate.

Integrated nonlinear production rate was determined by the combined effects of regulatory interactions. Each parameter in the equations has a specific biological meaning.  $B_i$  denotes the basal production rate of cluster  $i$ . Activations and inhibitions are modeled using sigmoidal form functions in a standard way.  $V_i$  denotes the maximal production of cluster  $i$ ,  $k_i$  and the expression level of cluster  $i$  at which  $V = V_i/2$ ,  $n$  denotes the cooperative coefficients.  $\tau_i$  denotes the time constant of cluster  $i$  normalized to the degradation rate. We assumed that the TF expression levels were regulated at similar time scale. As it is difficult to obtain accurate values of these parameters, we used a dimensionless modeling framework in which the activation level of each protein was normalized to a range from 0 to 1, 0 denotes minimal activation, 1 means full activation. Therefore  $V_{TF\_WT} = 1, V_{TF\_HCC} = 1$  by definition. Simulation result in Fig. 2D was obtained using time constant  $\tau_{TF\_WT}$  and  $\tau_{TF\_HCC}$  as 1, cooperative coefficients as 3 (as most cooperative coefficients in biochemical reactions ranged from 2 to 4), and  $k_i$  as 1/8. This modeling framework provided a coarse-graining description of network topology, and the microscopic detail and parameter values are not required to be overly precise. Given parameters in the model, calculations of attractors from equations and simulations were conducted as described previously<sup>4</sup>.

### ***A multi-layer computational framework to calculate index from the transcriptomes***

Establishing a multi-layer model consisting of model training and testing, we first outline these steps and then provide detailed description of each step.

Model training:

1. Collect transcriptomic data of adult livers and HCCs for model training.
2. Calculate significantly up/down-regulated TF clusters in HCCs compared with the WT adult livers.
3. Calculate optimized weights of TF cluster's target genes to determine the changes of TF clusters between adult livers and HCCs.
4. Calculate the activities of TF clusters in each training WT adult liver or HCC sample.
5. Calculate the tumor-promoting and -inhibiting strengths in each training WT liver or HCC sample using the averaged activities of up- and down-regulated TF clusters of each sample.
6. Calculate the weight of tumor-promoting and -inhibiting strengths to define a TI score.

### ***Model testing:***

7. Using the trained model and the RNA-seq data in query, we calculated the activity of each TF cluster, averaged activity of up/down-regulated TF cluster, and TI score.

**Step 1, Data collection for model training.** We divided the collected samples into four groups (Fig. 1A). The transcriptomic data of WT adult livers and HCC samples were used for model training.

**Step 2, Identification of significantly changed TF clusters.** Based on the whole transcriptomes of training samples and the gene list of each TF cluster, we identified 36 up-regulated (p-value <0.05 and set score >0.8) and 25 down-regulated TF clusters (p-value <0.05 and set score <-0.8) in HCC. These 61 significantly changed TF clusters were used as nodes in the second layer.

**Step 3, Calculation of an optimized gene weight to define the changes of TF clusters.** For each significantly changed TF cluster, we used its target genes as predictor variables and calculated the optimized weight of each target genes, to define the change of a TF cluster between WT adult liver and HCC using Random Forests. We used  $w_{i,j}$  to denote the weight of gene  $j$  to distinguish TF cluster  $i$ .

**Step 4, Quantification of TF clusters' activities in each WT adult liver or tumor tissue.** We identified significantly up/down-regulated TF clusters in HCC samples compared to WT adult livers. To quantify the activity of these TF clusters in each sample, we reasoned that the gene expression profile of a WT adult liver should fall within a range determined by the expression profiles of all WT adult liver tissues, and the expression profile of a HCC sample should fall within a range determined by the expression profiles of all HCC tissues. We formulated a notion called gene expression deviation/difference to quantify the differences between each and all WT adult liver samples or HCC samples. First, we normalized the expression level of each gene in the query sample using expression distribution of each gene in all WT adult liver samples or HCC samples as reference/background. The normalized expression level of each gene was calculated using Z score. The difference/deviation of TF cluster between gene expression patterns in query and WT liver/HCC expression pattern can be defined by adding contribution, according to weight and expression deviation, for each gene. Therefore, the deviation of TF cluster  $i$  in a query  $q$ , compared with that in WT liver samples can be quantified using equation:

$$TF_i^{WT}(q) = \sum_{j=1}^{n_i} (|z_j^{WT}(q)| * w_{i,j})$$

Where  $n_i$  denotes the number of target genes in a TF cluster  $i$  and  $z_{i,j}^{WT}(q)$  denotes normalized expression level of gene  $j$  in sample  $q$  with gene expression level distribution of gene  $j$  in the complete adult WT liver samples as reference/background.  $w_{i,j}$  denotes weight of gene  $j$  to distinguish TF cluster  $i$ , which is derived in Step 3 using Random Forests. Therefore, the extent to which TF cluster  $i$  in gene expression of sample  $q$  is deviated from adult WT liver is quantified. We transformed and normalized the raw deviation score using equation:

$$NTF_i^{WT}(q) = \frac{TF_i^{WT}(q)}{(\sum_{s_k \text{ is WT sample}} TF_i^{WT}(q)) / \text{WT sample number}}$$

Where  $(\sum_{s_k \text{ is WT sample}} TF_i^{WT}(q)) / \text{WT sample number}$  is the averaged deviation of samples in the complete training data of adult WT liver. In this way, we quantified the deviation of each TF signature in query sample compared to WT liver.

Using this approach, we calculated the deviation of each TF cluster in each sample compared to WT liver sample and tumors in PKO and DKO mice, respectively. Based on these deviations of each TF cluster in a query expression pattern, we quantified its status using equation:

$$TF(i, q) = \frac{NTF_i^{WT}(q)}{NTF_i^{WT}(q) + NTF_i^{HCC}(q)} * TF_{i_{HCC}} - \frac{NTF_i^{HCC}(q)}{NTF_i^{WT}(q) + NTF_i^{HCC}(q)} * TF_{i_{HCC}}$$

Where  $TF(i, q)$  denotes the change of TF cluster  $i$  in sample  $q$ . By definition, the  $TF(i, q)$  range from -1 to 1, where 1 denotes this TF cluster was fully up-regulated compared to WT liver, 0 no significant change, and -1 down-regulation. In this way, we quantified changes of each TF cluster using gene expression profile of WT liver and HCC.

**Step 5, Calculation of tumor-promoting and -inhibiting strengths of each sample in the third layer.** In the third layer, these significantly changed TF clusters were divided into two groups according to their up- or down-regulated expression in tumors. These 36 TF clusters up-regulated in tumors were assumed to be pro-tumorigenic, with the 25 down-regulated TF clusters being anti-tumorigenic. Activities of these TF clusters in each sample were obtained in step 4. For each sample  $q$ , the average activity of these 36 up-regulated TF clusters was used to quantify

its tumor-promoting strength, and the average activity of these 25 down-regulated TF clusters was used to quantify its tumor-inhibiting strength.

$$G_{promoting}(q) = \frac{1}{N_{pro}} \sum_1^{N_{pro}} TF(i_{pro}, q)$$

$$G_{suppressing}(q) = \frac{1}{N_{sup}} \sum_1^{N_{sup}} TF(i_{sup}, q)$$

Where  $G_{promoting}(q)$  and  $G_{suppressing}(q)$  denote the average activity of pro-tumorigenic strength and average activity of anti-tumorigenic strength, respectively.  $N_{pro}$  and  $N_{sup}$  denote the number of up-regulated TF cluster and down-regulated TF cluster, respectively.

**Step 6, Calculation of the weight of tumor-promoting and -inhibiting strengths in defining the TI.** The tumor-promoting strength and tumor-inhibiting strength of adult WT liver samples and HCC samples were calculated in step 5, which were used to define HCC index using the following equation.

$$Index(q) = c_1 * G_{promoting}(q) + c_2 * G_{suppressing}(q) + c_0$$

Where  $c_1$  and  $c_2$  denote weights of tumor-promoting strength and tumor-inhibiting strength in defining TI. We used our training adult WT liver sample and HCC samples to infer the optimized values of these weights. We defined the index range from -1 to 1, -1 denotes adult WT liver and 1 denotes HCC. By definition, the index of adult WT liver samples in training datasets is -1 and HCC samples is 1. We then performed regression analysis to calculate optimized parameter values using LASSO.

**Step 7, Model testing.** For each query transcriptome, we calculated activities of TF clusters, tumor-promoting and -inhibiting strengths in the sample, and finally a TI score using the trained multi-layer computational model. In order to make data from different backgrounds and platforms comparable, we performed quantile normalization to diminish batch effects using R package preprocess Core. According to step 4, we calculated the activity of each TF cluster query gene expression pattern. According to step 5, we calculated tumor-promoting and -inhibiting strengths in the samples. At last, we calculated its index according to tumor-promoting and -inhibiting strengths according to step 6. According to the definition, if a TI score  $<0$ , it indicates this sample is more like normal liver or liver with chronic liver diseases before passing the transition point. However, if a TI score  $>0$ , it indicates this sample has passed the critical tumorigenic switch point more like a tumor tissue. The tumorigenic index scores of WT, SKO,



PKO and DKO livers were calculated according to their transcriptomes. The TIs of human liver samples were also calculated from their transcriptomic data deposited in the public datasets.

Supplementary Figures

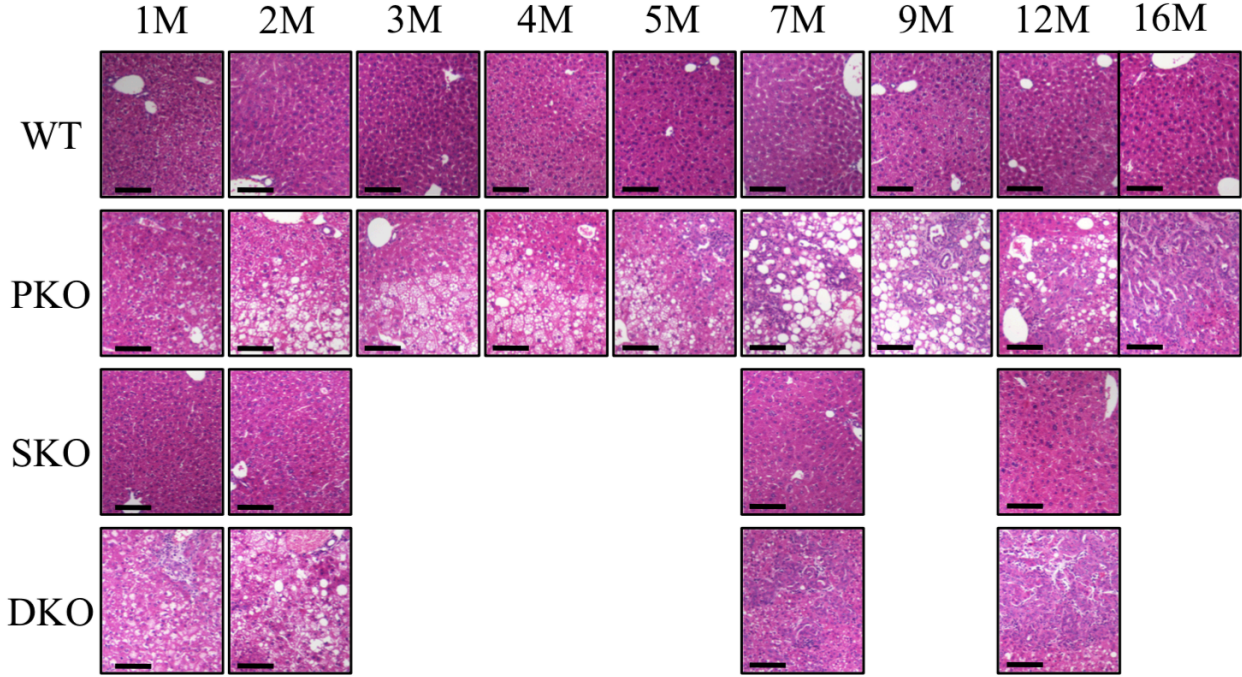
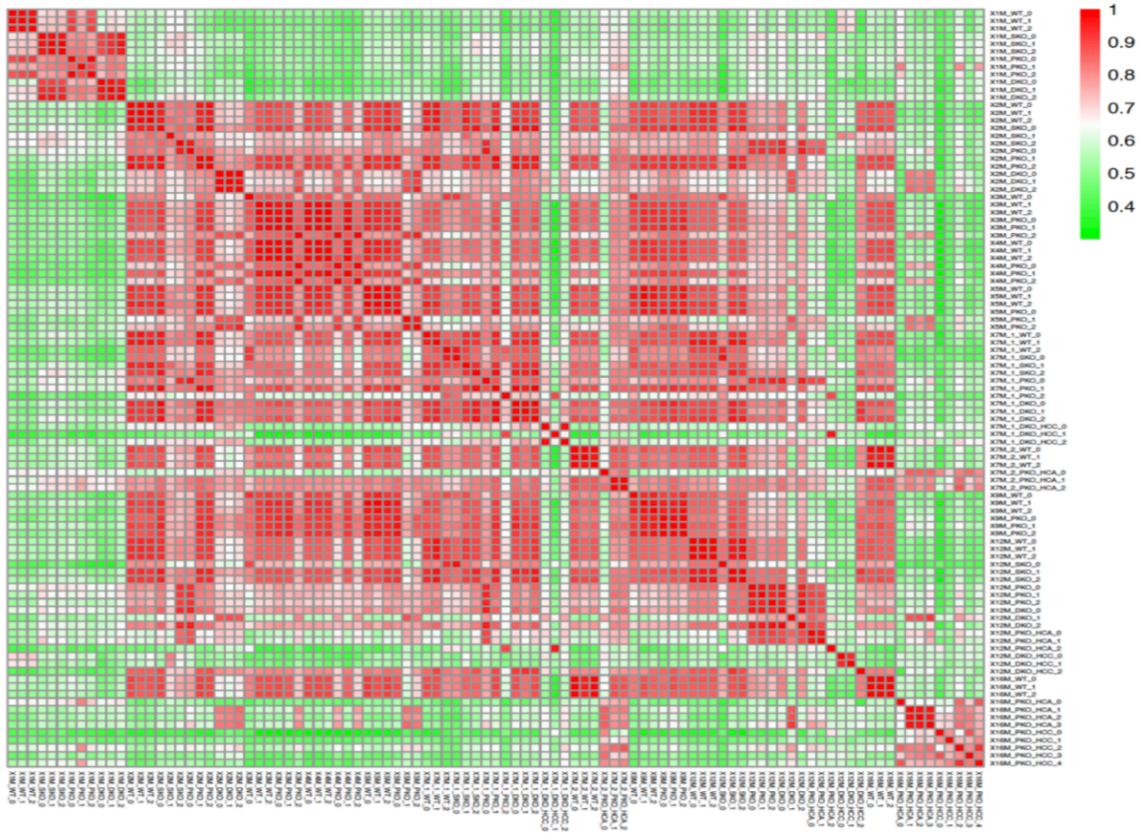


Figure S1

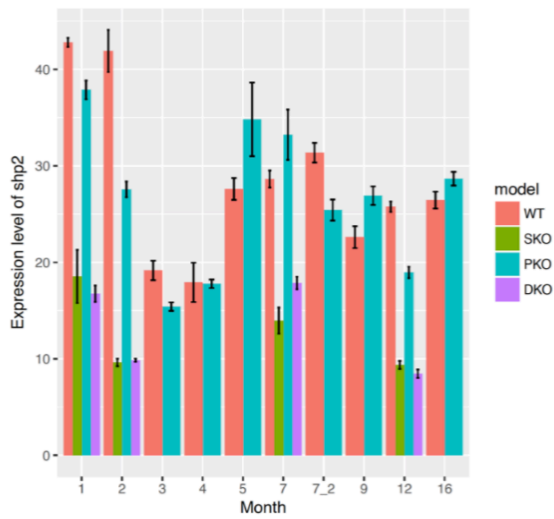
**Figure S1. Histopathological characterization of liver samples**

Representative liver sections with H&E staining are shown for the liver samples of four genotypes collected at various time points as shown in Fig.1A. Scale bars: 100  $\mu$ M.

A



B



C

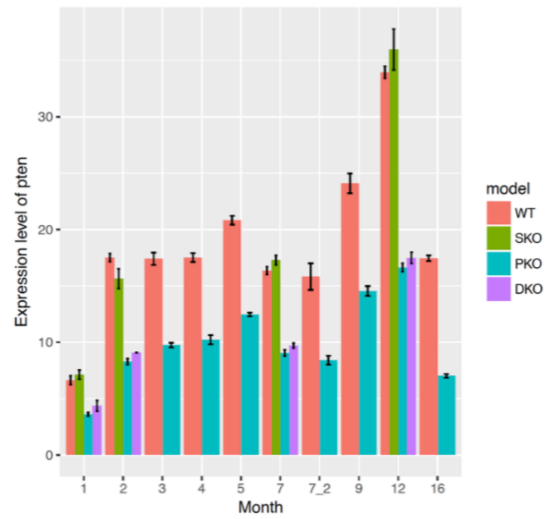


Figure S2

Figure S2. Quality control of RNA-seq data

(A). Pearson's correlations between samples were calculated and visualized using heatmap.

(B). Expression levels of *Shp2* in *SKO* livers (hepatocyte-specific deletion of *Shp2*) and *DKO* livers (hepatocyte-specific deletion of *Shp2* and *Pten*) were significantly lower than *WT* livers.

(C). Expression levels of *Pten* in *PKO* livers (hepatocyte-specific deletion of *Pten*) and *DKO* livers (hepatocyte-specific deletion of *Shp2* and *Pten*) were significantly lower than *WT* livers.

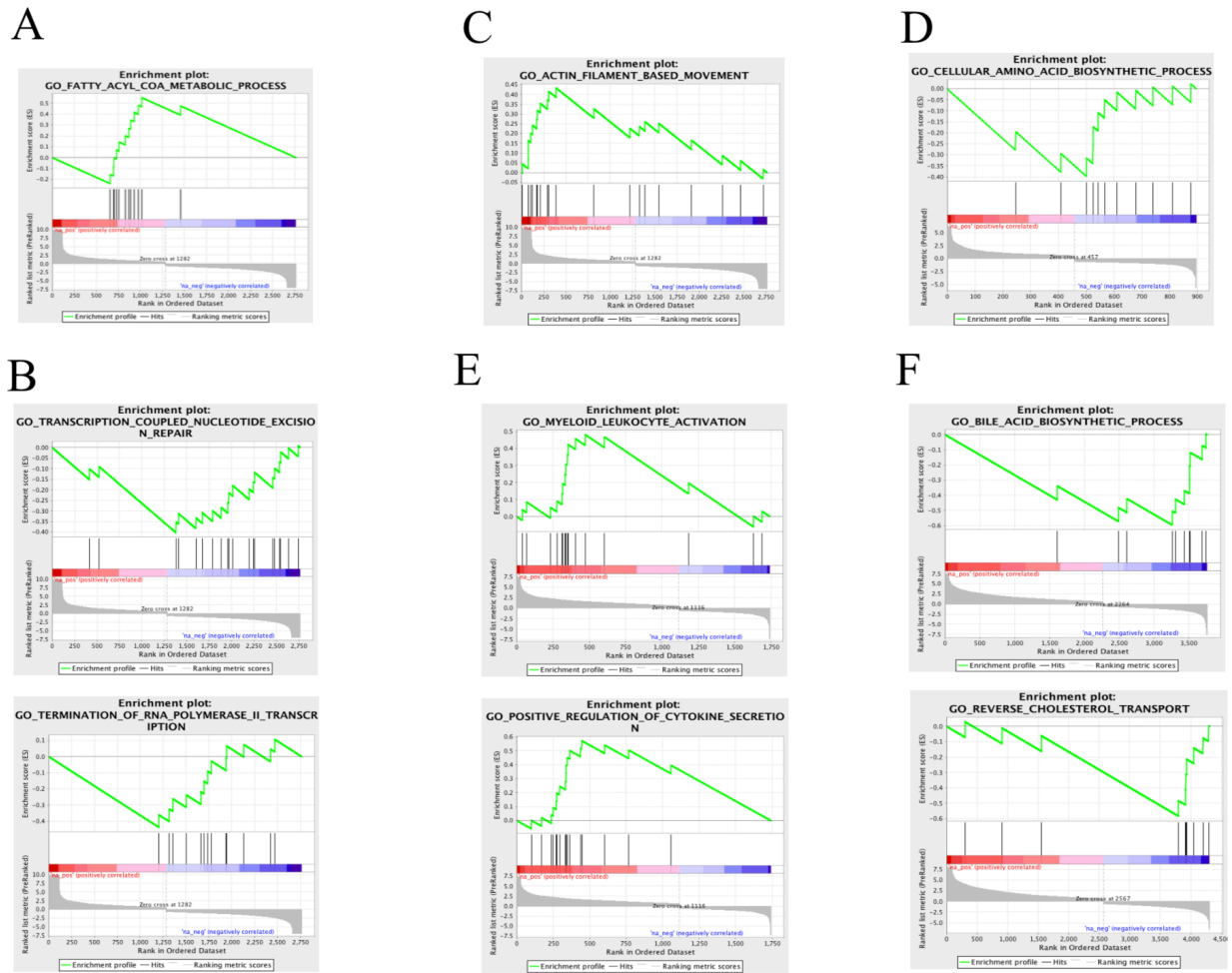


Figure S3

**Figure S3. Significantly changed biological processes during tumorigenesis in PKO livers**

- (A). GSEA showed significant changes in fatty acid and lipid metabolic processes in PKO livers starting from 1-month (1M).
- (B). Expression of epigenetic and DNA repair components started to change significantly at 1M.
- (C). Expression of extracellular structural elements changed significantly starting from 1M.
- (D). Expression of amino acid metabolic processes started to change from 2M.
- (E). Expression of several immunological and inflammatory processes started to change from 5M.
- (F). Expression of bile acid and cholesterol metabolic processes had significant changes at 12M.

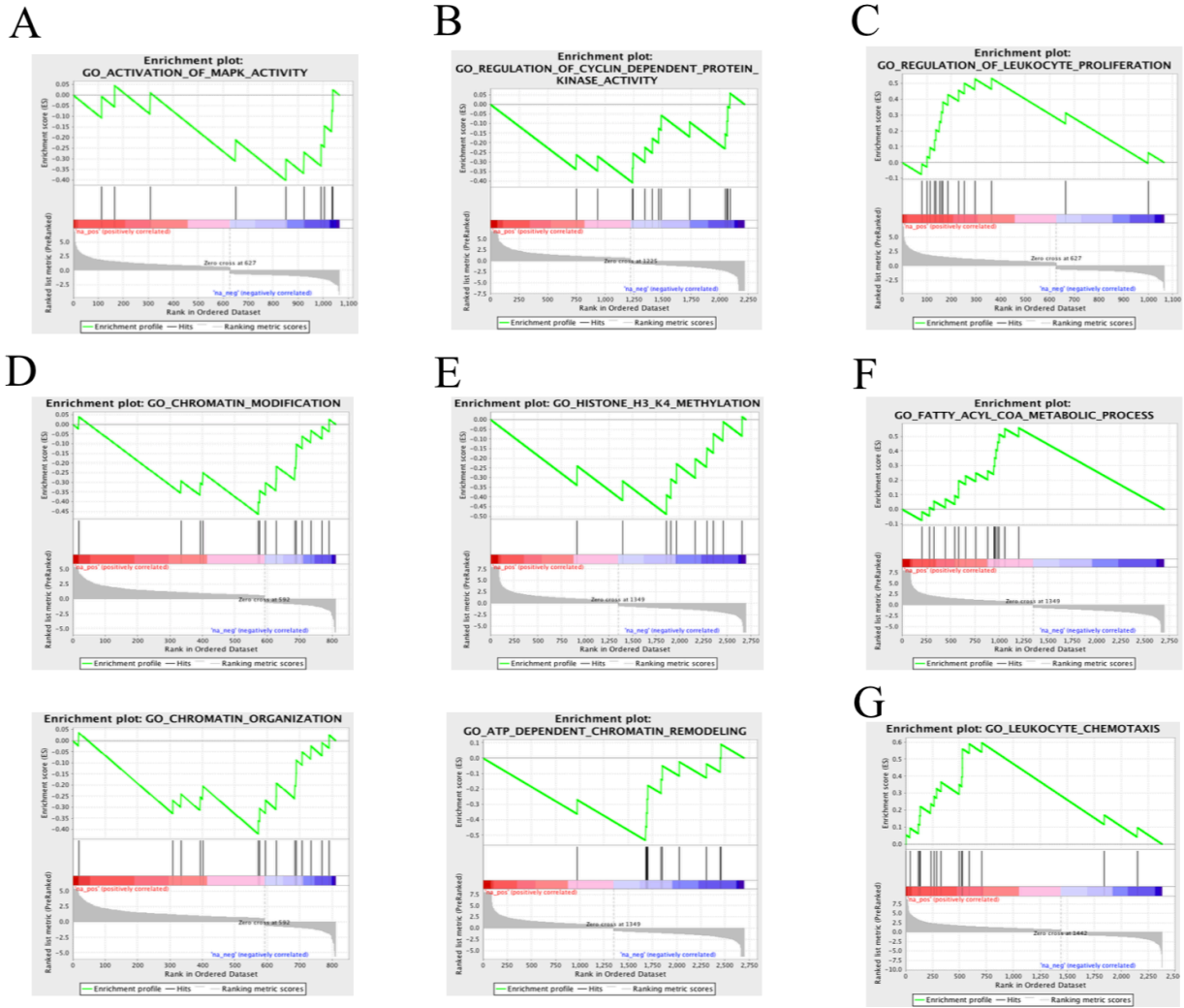
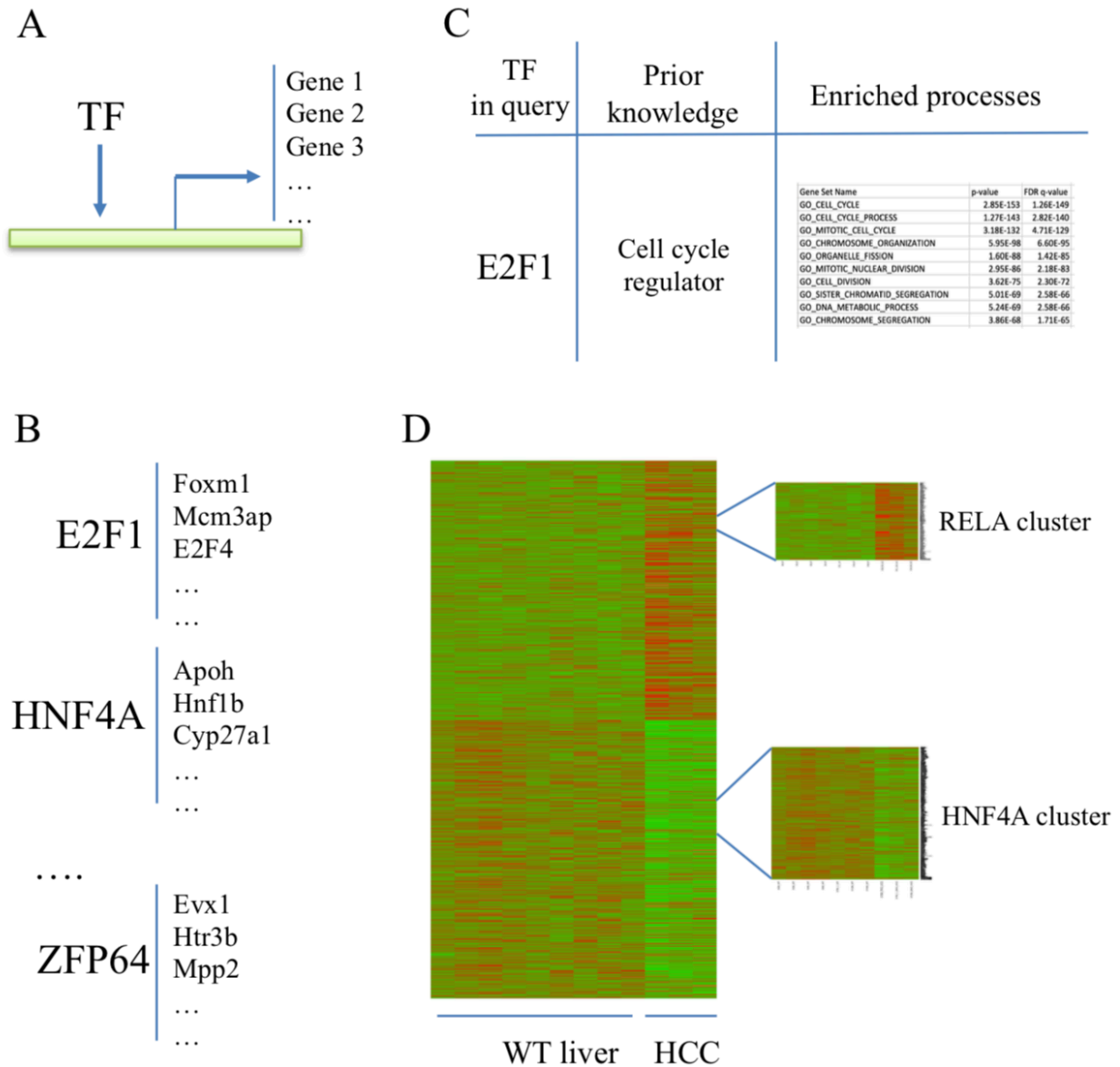


Figure S4

**Figure S4. Significantly changed biological processes during tumorigenesis in SKO and DKO livers**

- (A). GSEA showed significant changes in MAPK activity in SKO livers starting from 1-month (1M).
- (B). Cell cycle processes started to change from 1M in SKO liver.
- (C). Immunological and inflammatory process started to change from 1M in SKO liver.
- (D). Epigenetic machinery changed significantly at 12M in SKO liver.
- (E). Epigenetic machinery started to change from 1M in DKO liver.
- (F). fatty acid metabolism started to change from 1M in DKO liver.
- (G). Immunological and inflammatory process started to change from 1M in DKO liver.



**Figure S5**

**Figure S5. Significantly changed TF clusters in tumors relative to WT adult livers**

(A). Schematic illustration of TF downstream genes that co-express with a TF to control specific cell activities. A TF and its co-expressed genes are combined into a TF cluster for further analysis.

(B). Co-expressed genes of each TF were obtained from TF-gene correlation in public databases Cellnet and geneFriend. Genes that have high positive correlation with each TF (Pearson's

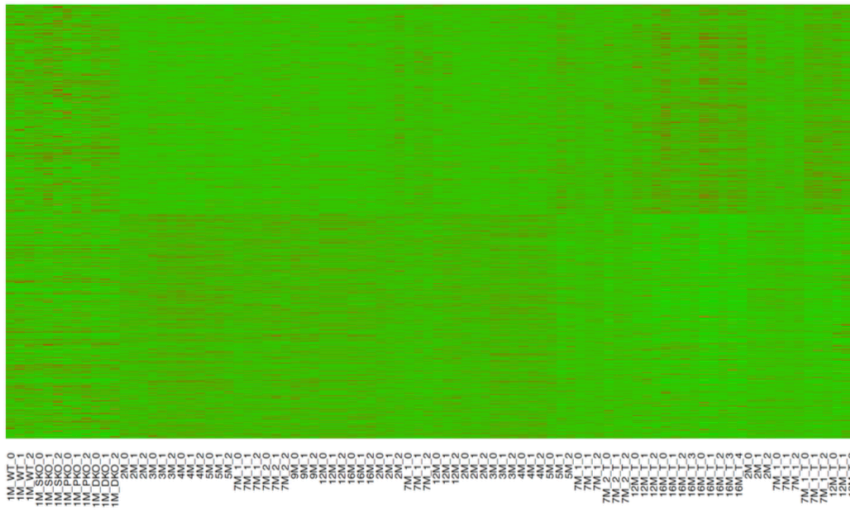
correlation>0.5) were identified and grouped into a TF cluster, with a total of 1568 TF clusters established.

(C). Functional assessment of co-expressed genes in a TF cluster to capture its key biological features using E2F1 as an example. Consistent to the known function of E2F1 in control of cell cycle progression, gene ontology analysis of the E2F1 cluster members identified cell cycle as the top-enriched biological process. In a similar way, we tested other TF clusters by comparing their well-documented functions and the top-enriched biological processes of the co-expressed genes (Table S6).

(D). A total of 36 up-regulated and 25 down-regulated TF clusters were identified in tumors (PKO at 16 months, DKO at 7 and 12 months), as compared to WT livers (WT adults at 2, 3, 4, 5, 7, 9, 12, and 16 months) using gene set analysis. Expression of the up- and down-regulated TF clusters in WT livers or tumors was visualized using heatmap. For example, the RelA and HNF4 $\alpha$  clusters were significantly up- and down-regulated in tumors, respectively.

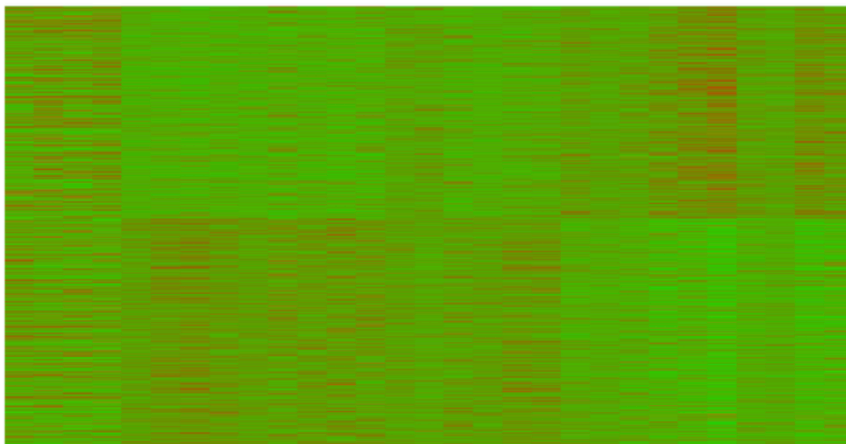


A



Youth WT SKO PKO DKO

B



1M\_WT  
1M\_SKO  
1M\_PKO  
1M\_DKO  
2M  
3M  
4M  
5M  
7M  
7M  
9M  
12M  
16M  
2M  
7M  
12M  
2M  
3M  
4M  
5M  
7M  
7M\_T  
12M\_T  
16M\_T  
16M\_T  
2M  
7M  
7M\_T  
12M\_T

Youth WT SKO PKO DKO

Figure S6

**Figure S6. Temporal gene expression patterns in WT, SKO, PKO and DKO livers**

(A). Heatmap for temporal gene expression profiles of TF clusters including all replicated samples collected from WT, SKO, PKO and DKO livers at various time point.

(B). Heatmap of the averaged gene expression levels in all replicates of samples in (A).

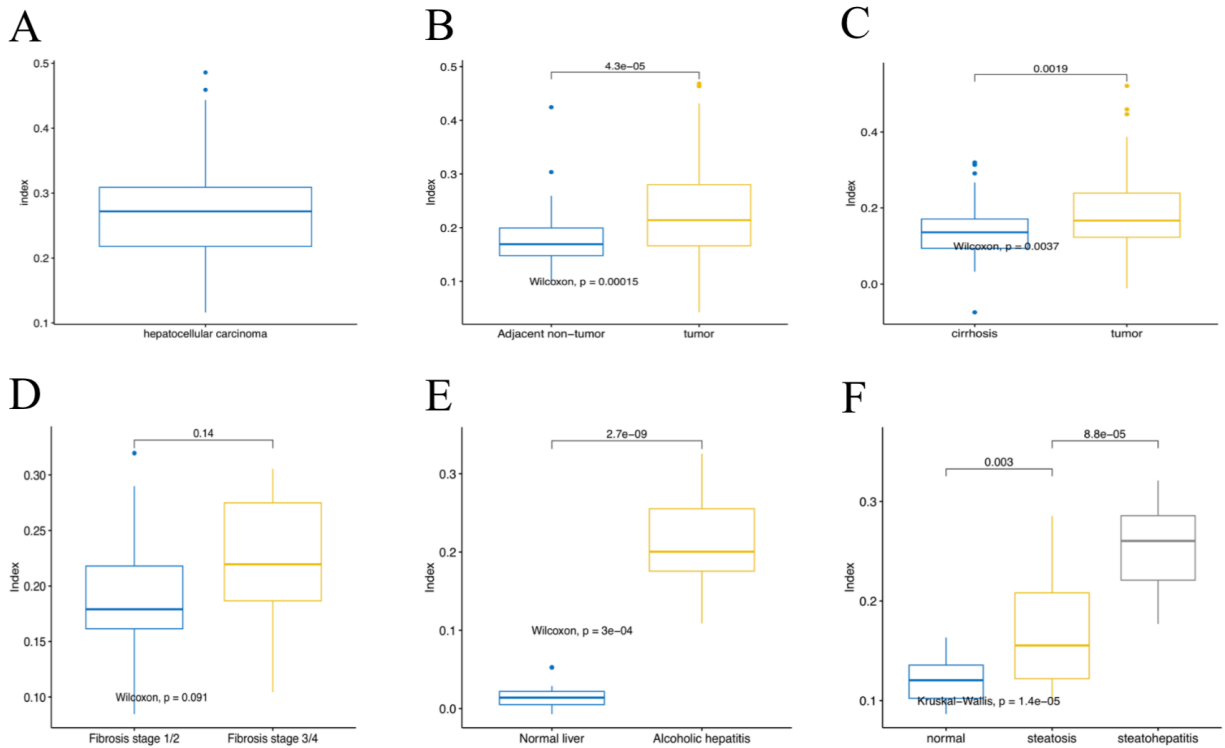


Figure S7

**Figure S7. TI values of human non-tumor, pre-cancer and cancer samples with different backgrounds**

(A). The TI values of 371 HCC samples in GSE19977 dataset were obtained at the time of surgical resection.

(B). The TIs of 52 adjacent non-tumor liver samples and 115 liver tumor samples in GSE76427.

(C). The TIs of 79 cirrhotic liver samples and 61 tumor samples in GSE54236.

(D). The TIs of 36 HCV-related fibrotic liver samples (fibrosis stages 1/2:  $n=18$ ; fibrosis stages 3/4:  $n=18$ ) in GSE33258. The averaged TI value of fibrosis stages 3/4 was higher than that of fibrosis 1/2.

(E). The TIs of 7 control liver samples and 15 alcoholic hepatitis samples in GSE28619. The TI values increased significantly from control liver to alcoholic hepatitis.

(F). The TIs of 13 non-tumor liver samples, 19 steatosis liver samples and 12 steatohepatitis liver samples in GSE33814. The TI values increased gradually from non-tumor liver, steatosis liver to steatohepatitis liver.

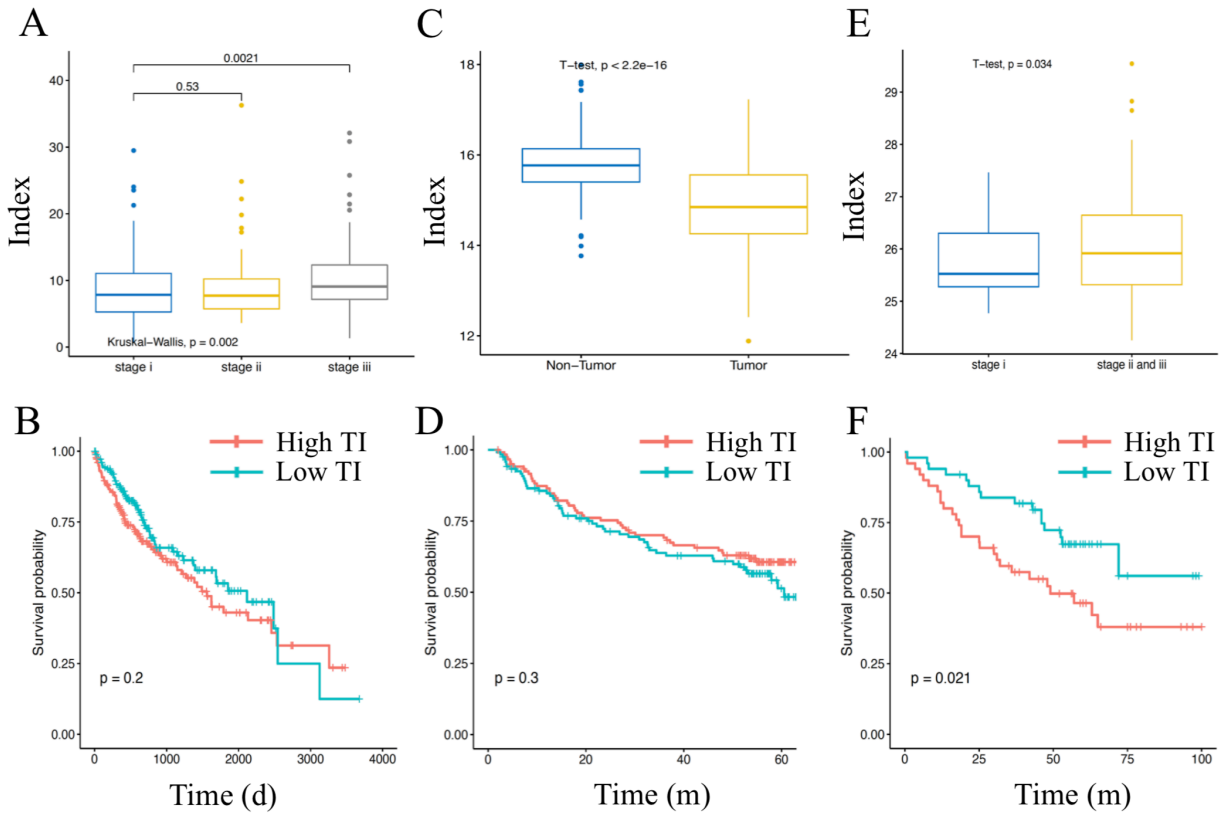


Figure S8

**Figure S8. The TI derived using LASSO was not a good predictor of clinical outcomes**

(A). TIs were calculated using LASSO based on the same data in TCGA dataset as in Fig. 5G. The TI values failed to distinguish stage I and II liver tumors.

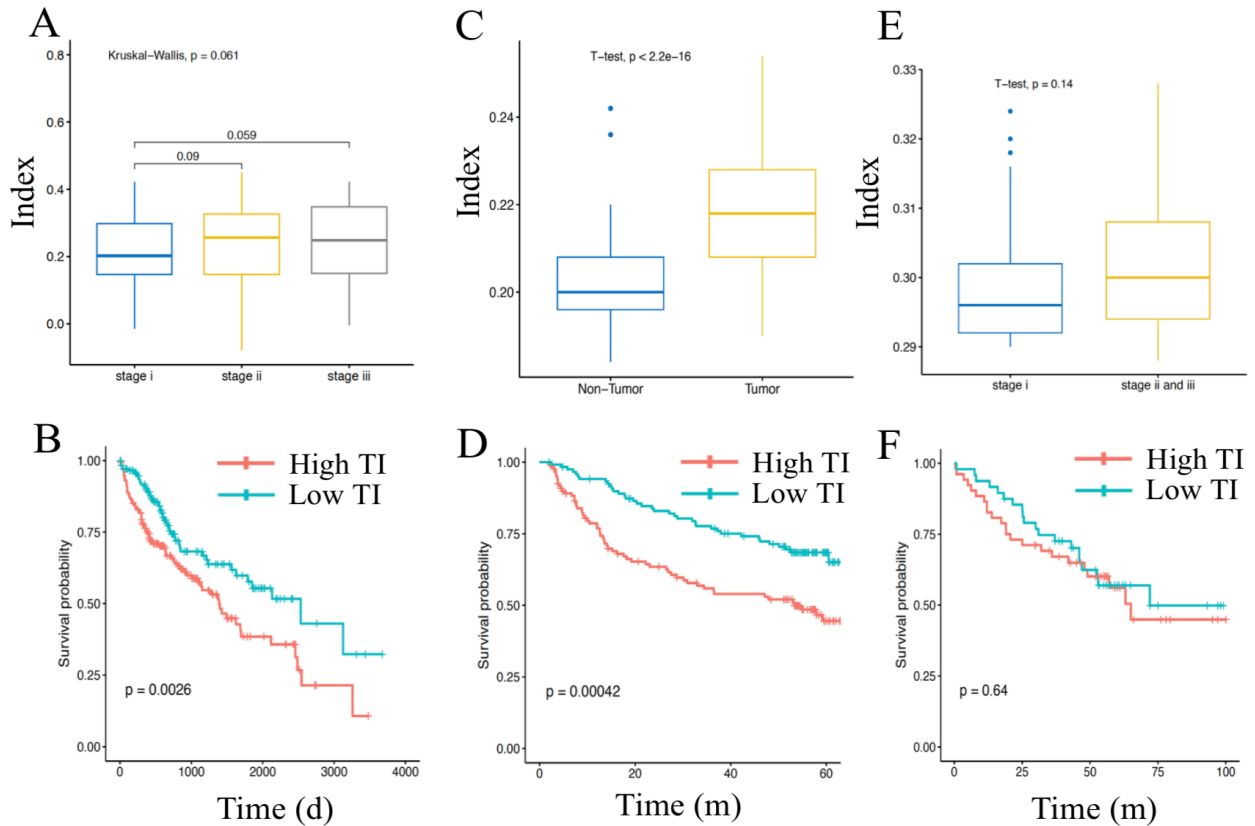
(B). The 371 liver tumor samples in TCGA dataset were divided into low and high TI groups at the median TI (high TI: n = 186; low TI: n= 185), Kaplan-Meier analysis showed no significant difference of survival between the two subgroups (log-rank test, p=0.2).

(C). TIs were calculated using LASSO based on the same microarray data in Fig. 5I of 445 non-tumor or tumor samples in GSE14520 (tumor: n=221; non-tumor: n=224). The TI values failed to predict tumor stages, the average TI of non-tumor samples was even higher than the tumor samples.

(D). The 221 tumor samples in GSE14520 were divided into low and high TI groups at the median TI (high TI: n = 111; low TI: n = 110), Kaplan-Meier survival analysis showed no significant difference between the two subgroups (log-rank test, p=0.3).

(E). TIs were calculated using LASSO based on the same microarray data in Fig. 5K of 100 liver tumor samples in GSE16757, with well-documented clinical data of tumor stages (stage I: n=35; stage II and III: n=65). The TI values correlated with the tumor stages in GSE16757, with advanced HCCs having significantly higher TI.

(F). The 100 tumor samples in GSE16757 were divided into low and high TI groups by the median TI (high TI: n = 50; low TI: n = 50). Kaplan-Meier survival analysis showed shorter survival and poor prognosis for patients with high TI, compared to the low TI group (log-rank test,  $p=0.021$ ).



**Figure S9**

**Figure S9. The TI derived using Random Forest was less effective in predicting clinical outcomes**

(A). TIs were calculated using Random Forest based on the same data of 371 liver tumor samples in TCGA dataset in Fig. 5G. A total of 347 tumor samples were deposited with well-documented tumor stages (stage I:  $n=171$ ; stage II:  $n=86$ ; stage III and more advanced:  $n=90$ ). The TI values calculated using Random Forest failed to distinguish the tumor stages.

(B). The 371 liver tumor samples in TCGA dataset were divided into low and high TI groups by the median TI (high TI:  $n=186$ ; low TI:  $n=185$ ). Kaplan-Meier survival analysis showed short survival and poor prognosis for patients with high TIs (red), compared to the low TI group (blue) (log-rank test,  $p < 0.003$ ).

(C). TIs were calculated using Random Forest based on the same data in Fig. 5I of 445 non-tumor or tumor liver samples in GSE14520 (tumor:  $n=221$ ; non-tumor:  $n=224$ ).

(D). The 221 tumor samples in GSE14520 were divided into low and high TI groups at the median TI (high TI:  $n=111$ ; low TI:  $n=110$ ). Kaplan-Meier survival analysis showed shorter

survival and poorer prognosis for patients with high than the low TI group (log-rank test,  $p < 0.0005$ ).

(E). TIs were derived using Random Forest from the same data in Fig. 5K of 100 liver tumor samples in GSE16757 dataset, with well-documented tumor stages (stage I:  $n=35$ ; stage II and III:  $n=65$ ). The TI values failed to distinguish tumor stages.

(F). The 100 tumor samples in GSE16757 dataset were divided into low and high TI groups at the median TI (high TI:  $n=50$ ; low TI:  $n=50$ ). Kaplan-Meier survival analysis showed no significant difference of survival between the two subgroups (log-rank test,  $p=0.64$ ).

**Dataset S1 (separate file)**, differentially expressed genes in mutant livers at different months

**Dataset S2 (separate file)**, significantly changed biological processes

**Dataset S3 (separate file)**, significantly changed ligands and receptors

**Dataset S4 (separate file)**, significantly changed epigenetic regulators

**Dataset S5 (separate file)**, significantly changed pathways

**Dataset S6 (separate file)**, correlated genes of each TF cluster

**Dataset S7 (separate file)**, test of TF clusters

**Dataset S8 (separate file)**, significantly changed TF clusters in tumors compared to WT liver

**Dataset S9 (separate file)**, correlation between the TF clusters inferred from RNA-seq data

**Dataset S10 (separate file)**, Mouse liver disease datasets

**Dataset S11 (separate file)**, Human liver disease datasets

## References

- 1 Bard-Chapeau, E. A. *et al.* Concerted functions of Gab1 and Shp2 in liver regeneration and hepatoprotection. *Mol Cell Biol* **26**, 4664-4674, 2006.
- 2 Bard-Chapeau, E. A. *et al.* Ptpn11/Shp2 acts as a tumor suppressor in hepatocellular carcinogenesis. *Cancer Cell* **19**, 629-639, 2011.
- 3 Luo, X. *et al.* Dual Shp2 and Pten Deficiencies Promote Non-alcoholic Steatohepatitis and Genesis of Liver Tumor-Initiating Cells. *Cell Rep* **17**, 2979-2993, 2016.
- 4 Wang, G., Zhu, X., Gu, J. & Ao, P. Quantitative implementation of the endogenous molecular–cellular network hypothesis in hepatocellular carcinoma. *Interface Focus* **4**, 20130064, 2014.