

**The intratumoral bacterial metataxonomic signature predicts hepatocellular
carcinoma**

Supplementary Material

Supplementary methods

Fresh tissue collection, preparation of paraffin sections, and hematoxylin-eosin (H&E) staining

Fresh cancer and peri-tumor tissues were collected during surgery with caution to avoid contamination. The fresh tissues were stored in sterile tubes cooled on ice immediately after being isolated from the patient's body. All subsequent steps were performed on ice. All items were wiped clean with alcohol and autoclaved, including slides, knives, tubes, forceps, clippers, and containers. Each tissue was divided into three parts. One piece was processed immediately and used for tissue culture. Another part was fixed with paraformaldehyde (PFA) and used for H&E staining, Gram staining, and fluorescence in situ hybridization (FISH) analyses. The remaining part was stored at -80°C and used for PCR or other analysis. All these steps should be completed within one hour after the tissues were isolated from the body. To prepare paraffin sections, fresh surgery tissues were fixed overnight with 4% PFA. The tissues were trimmed and consecutively dehydrated using the following protocol: 4 h in 75% ethanol; 2 h in 85% ethanol; 2 h in 90% ethanol; 2 h in 95% ethanol; 1 h in 100% ethanol; 30 min in 100% ethanol; 30 min in ethanol/xylene (1:1); 5 – 10 min in 100% xylene; 5 – 10 min in 100% xylene. Next, the sections were soaked in paraffin using the following steps: 5 – 10 min in melted paraffin at 65°C; 1 h in melted paraffin at 65°C; 1 h in melted paraffin at 65°C. The slides were embedded and allowed to cool at -20°C. The paraffin-embedded tissues were cut into 4- μ m sections. The sections were incubated overnight in a 56°C oven.

After deparaffinization, the sections were soaked in 100% xylene for 20 min (two changes) and then dehydrated in 100% ethanol for 5 min (two changes) and 75% ethanol for 5 min. The sections were washed with running water. Next, the sections were stained with 2.5% hematoxylin solution for 3~5 min and then rinsed with running water. The sections were treated with 1% HCl in ethanol and washed with running water. The slides were washed with 1% ammonia solution for several seconds and were then rinsed using running water. Next, the sections were continuously dehydrated with 85% and 95% ethanol for 5min, respectively, staining and then stained with 1% eosin for 5 min. The sections were then consecutively dehydrated with 100% ethanol for 5 min (three changes) and 100% xylene for 5 min (two changes). Finally, the sections were mounted with neutral balsam.

Gram staining of HCC or peri-tumor tissues

Gram staining of tissues was performed with the Gram Staining kit (Wuhan Servicebio Technology Co., Ltd. Wuhan, Hu Bei, China) according to the manufacturer's instructions. Briefly, paraffin sections were deparaffinized using the procedure: 20 min in xylene (two changes), 5 min in 100% ethanol (two changes), 5 min in 75% alcohol, and washed with tap water. sections were stained with Gram Staining Solution A (2.5% crystal violet/ammonium oxalate in ethanol) for 10 s - 30 s, washing, and then washed with distilled water. slides were stained with Gram Staining Solution B (1% iodine) for 1 min and then washed with distilled water. The Gram Staining Solution C (95% ethanol) was used to rinse the slides from the frosted surface side of the slides for differentiation until the solution running off the section was colorless. The sections were immersed in Gram Staining Solution D (2.5% safranin) for one second and then washed with distilled water. The sections were dried in

an oven at 60°C. The dried sections were quickly immersed in 100% ethanol (three changes of 1s, 3s, and 5s in sequence) for dehydration. The sections were immersed in clean xylene for 5 min to increase transparency and then sealed with neutral balsam. The slides were scanned using the NanoZoomer S60 Digital Slide Scanner (C13210, Hamamatsu Photonics K.K. Co., Ltd., Japan). The images were viewed with NDP.view version 2.7.43 (Hamamatsu Photonics K.K. Co., Ltd., Japan).

16S rRNA fluorescence in situ hybridization (FISH)

Formalin-fixed paraffin-embedded (FFPE) tissue slides were deparaffinized and rehydrated using the procedure: 15 min in 100% xylene (two changes), 5 min in 100% ethanol (two changes), 5 min in 85% ethanol, and 5 min in 75% ethanol. The sections were washed with diethylpyrocarbonate (DEPC)-treated water. Slides were incubated for 10~15 min in boiled citrate-EDTA antigen retrieval solution and the solution was let cool naturally. Sections were incubated in prehybridization solution at 37°C for one hour. Hybridization was performed overnight at 37°C using Cy3-labelled probes EUB338 5'-GCTGCCTCCCGTAGGAGT-3' and Cy5-labelled non-specific complement probe 5'-CGACGGAGGGCATCCTCA-3' (5 ng/μL) as reported by Nejman D., *et al* (1). Sections were washed with 2 × SSC at 37°C for 10 min, 1 × SSC at 37°C for 5 min (two changes), and 0.5 × SSC at room temperature for 10 min. The nucleus was counterstained with DAPI (2 μg/μL) for 8 min in a dark place. Sections were mounted with antifade Mountant. The signal was captured using the PANNORAMIC MIDI digital slide scanner (3DHISTECH Ltd. Budapest, Hungary). The images were viewed with CaseViewer version 2.4 (3DHISTECH Ltd. Budapest, Hungary).

Aerobic and anaerobic cultures of cancer or peri-tumor tissues

Fresh cancer and peri-tumor tissues were immersed in normal saline in a sterile culturing dish. The surgical blade was wiped with ethanol and passed through a flame. After cooling down, the blade was used to cut the tissues into two pieces, which were then put into two 1.5-mL Eppendorf tubes. One part of the tissues was subjected to aerobic culture using the Brain Heart Infusion (BHI) Broth (HB8297-5, Hopebio Biotechnology Co., LTD., Qingdao, Shandong, China). Another part of the tissues was used for anaerobic culture using Gifu Anaerobic Medium (GRM) (HB8518-1, Hopebio Biotechnology Co., LTD., Qingdao, Shandong, China) (supplemented with 0.0005% vitamin k1 and 25 µg/mL hemin). The tissue in each tube was ground for ~250 seconds in a freezing grinder (4°C, 2500 rpm) using a pre-autoclaved stainless-steel bead (6.2 mm in diameter) (Shanghai Jingxin Industrial Development Co. LTD, Shanghai, China). During the grinding process, the tubes were examined several times to ensure that the tissue was completely homogenized. A 150-µL homogenate of each sample was pipetted to an agar plate (1.5% agar in GAM or BHI medium) and spread on the surface with a sterilized triangle end glass spreader. For aerobic culturing, the plates were incubated in an incubator and cultured at 37°C for 3~5 days. For anaerobic culture, the plates were placed in a 7.0 L sealed culture tank (Mitsubishi Gas Chemical Company, INC. Japan) was equipped with three 2.5 L anaerobic gas-producing bags (HBYY001, Hopebio Biotechnology Co., LTD., Qingdao, Shandong, China) and oxygen indicator (HBYY004, Hopebio Biotechnology Co., LTD., Qingdao, Shandong, China). The sealed culture tank was placed in an incubator at 37°C for 3~6 days. The negative controls for aerobic and anaerobic cultures were blank BHI and GAM

media, respectively. To prepare positive controls for aerobic culture, we collected the environmental microbiota by wiping a lab door handle with a sterile cotton swab, and the cotton swab was immersed in the BHI medium used for aerobic culture. The positive control for anaerobic culturing was done in the same way except for the use of the GAM medium.

Taxonomic identification of bacteria colony isolated from the tissue culturing cultures

The colonies from the aerobic and anaerobic cultures of cancer or peri-tumor tissues were identified by 16S rRNA gene sequencing anaerobic using the 2 × EasyTaq® PCR SuperMix (TransGen Biotech Co., LTD., Beijing, China). Aerobic colonies were picked and amplified by agitation overnight at 37°C in BHI media, while anaerobic colonies were picked and cultured overnight at 37°C in GAM media without agitation. 10 µL of the culture was diluted with 90 µL of H₂O. The PCR reaction mix contained 2 µL of diluted DNA template, 1 µL of each primer (forward primer 27F, 5'-AGAGTTTGATCCTGGCTCAG-3' and reverse primer 1492R, 5'-GGTTACCTTGTTACGACTT-3'), 10 µL of 2 × EasyTaq PCR SuperMIX and 6 µL of H₂O. The PCR amplification cycle consisted of a pre-denaturing step at 94°C for 5 min, followed by 30 cycles of 94°C for 30 s, 53°C for 30 s, 72°C for 2 min, and a final extension step at 72°C for 2 min. The PCR products were analyzed using agarose electrophoresis followed by Sanger DNA sequencing using the same primer pairs. The taxonomic identities of the sequences were determined by blasting the NT database using the BLASTN program in GenBank (<http://www.ncbi.nlm.nih.gov>). The PCR sequences, taxonomic identities, and the BLAST results were deposited in **Supplementary Table S2**.

Quantitative PCR (qPCR) validation of 16S rRNA gene of selected taxa

Total cellular DNA was extracted from liver tissues as described above. qPCR was performed with BeyoFast Probe qPCR Mix (2X, Low ROX) (Beyotime, Shanghai, China) using the Roche LightCycler® 480 II instrument (Roche Diagnostics Corporation, Indianapolis, USA) according to the manufacturer's instructions. The primers Gamma395f (5'-CMATGCCGCGTGTGTGAA-3') and Gamma871r (5'-ACTCCCCAGGCGGTCDACTTA-3') for Gammaproteobacteria were used in the qPCR analysis, as described by Pfeiffer S, et al (2). The TestPrime 1.0 program (<https://www.arb-silva.de/search/testprime/>) was used to evaluate the specificity of this primer pair on the SILVA database. The results indicate that this primer pair targets 28.4% of Gammaproteobacteria in the SILVA database. A qPCR analysis of human albumin exon 12 was performed in parallel using the primer pair: forward-TGTTGCATGAGAAAACGCCA and reverse-GTCGCCTGTTCAACCAAGGAT (3). The PCR reaction mix contains 0.5 µL diluted DNA template (100 ng/µL), 1 µL of each primer (final concentration 1 nmol), 5 µL 2 × BeyoFast Probe qPCR Mix (2X, Low ROX), and 2.5 µL H₂O. The three-step qPCR amplification cycle consists of a pre-denaturing step at 95°C for 10 min, followed by 42 cycles of 95°C for 5 s, 55°C for 10 s, and a final extension step at 72°C for 30 s. The melting curve was generated using the same conditions as above. The analysis was performed in a 384-well plate and 10 negative controls using blank reagents (no DNA template) were analyzed in parallel on each plate. The fluorescence cycle threshold (Ct) values were recorded for each amplicon. Delta Ct (Δ Ct) method was used to normalize the expression of 16S rRNA gene (Ct(16S rRNA) – Ct

(albumin)). The ΔC_t values were log₂-transformed and the p-value was calculated using a two-sided, unpaired Student's *t*-test.

Random-forest machine learning

The random-forest machine learning model was built based on OTU or Class. A training cohort containing 75% of total cases was built with samples randomly selected from the normal and HCC samples using the R “sample” function. The remaining 25% of samples were used as the validation cohort. Random-forest machine learning analysis was performed using R package randomForest v4.6.14 and 500 trees were built with 10-fold cross-validation. The normalized additive predicting probability was calculated and used as the final predicting probability. The higher probability of the binary classification was used as the predictive label. The features were ranked from high to low Mean Decrease Accuracy value and selected based on the cross-validation curve. The top-5 features were selected to build the model. The prediction power of the selected features was validated in the validation cohort.

Chi-square test

The association between taxon abundance and the clinicopathological parameters of HCC patients was analyzed using the Chi-square test. Taxon abundance was determined based on the median read count of each taxon. Taxon with an abundance lower than the median was classified as “low”, while taxon with an abundance higher than the median was classified as “high”. A two-sided Chi-square test was performed using GraphPad Prism 7.04. $p < 0.05$ was considered as statistically significant.

Statistics

Statistical significance was calculated using a two-sided Student's test. P-value or corrected p-value < 0.05 is considered as statistically significant. All measured values were expressed as mean \pm SE.

Figure S1

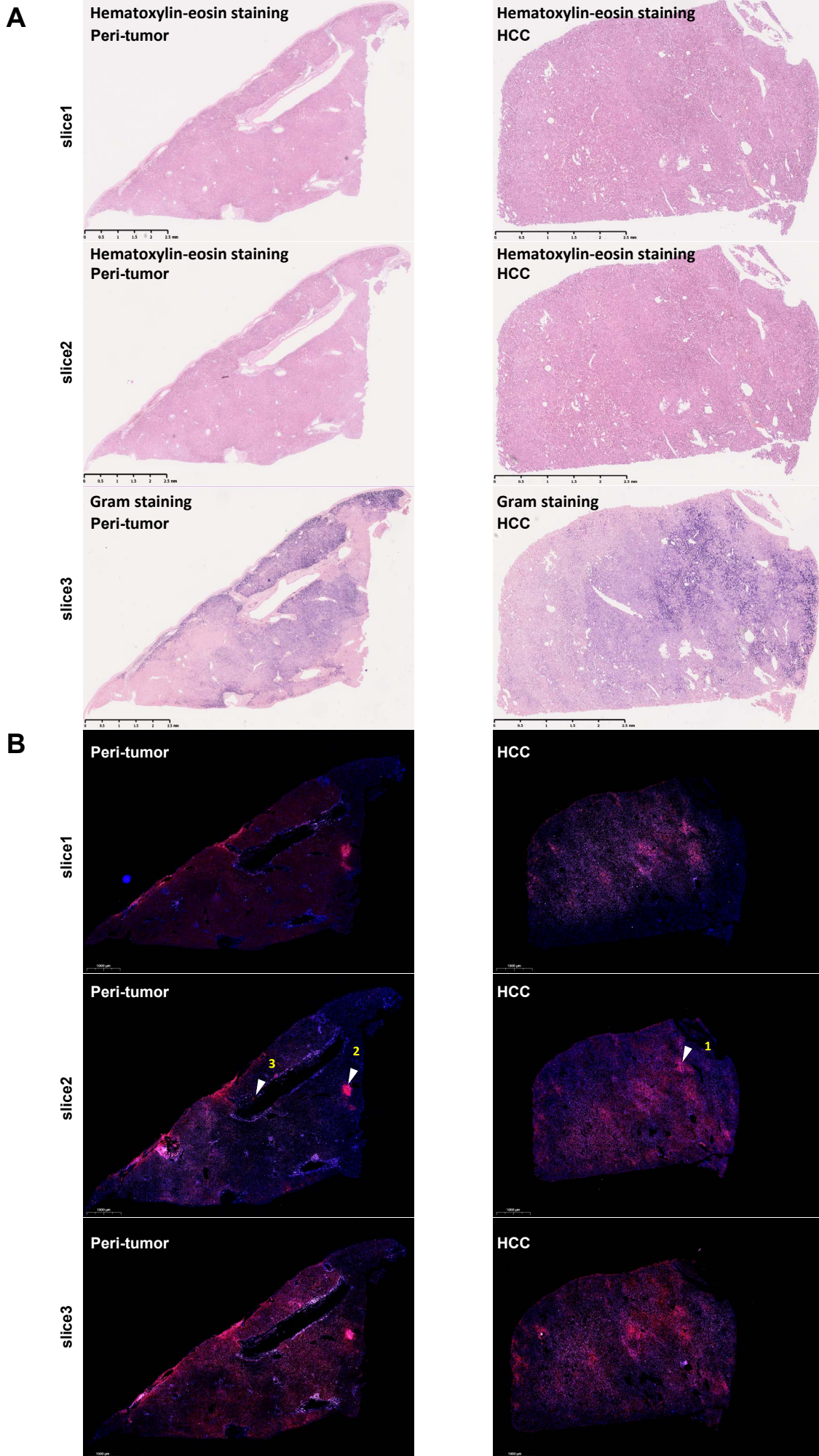


Figure S1. 16S rRNA fluorescence in situ hybridization (FISH) analysis of an HCC tissue and its associated peri-tumor tissue. (A) Three pieces of sequential formalin-fixed paraffin-embedded (FFPE) sections of peri-tumor (left) or liver cancer (right) were deparaffinized, rehydrated, and stained with hematoxylin-eosin (H&E, panel 1 and panel 2) or Gram staining (panel 3). (B) Three pieces of sequential FFPE liver cancer sections were deparaffinized, rehydrated, and probed with Cy3-labelled probes EUB338 (5'-GCTGCCTCCCGTAGGAGT-3') (red color) or Cy5-labelled non-specific complement probe (5'-CGACGGAGGGCATCCTCA-3') (pink color) Both probes have been recently used to analyze the microbiota of human cancer tissues (1). The nuclei were counterstained with diamidino-phenyl-indole (DAPI). The white arrowheads indicate the regions highlighted in Fig. 1 and Figure S2, Figure S3 (region 2), and Figure S4 (region 3).

Figure S2

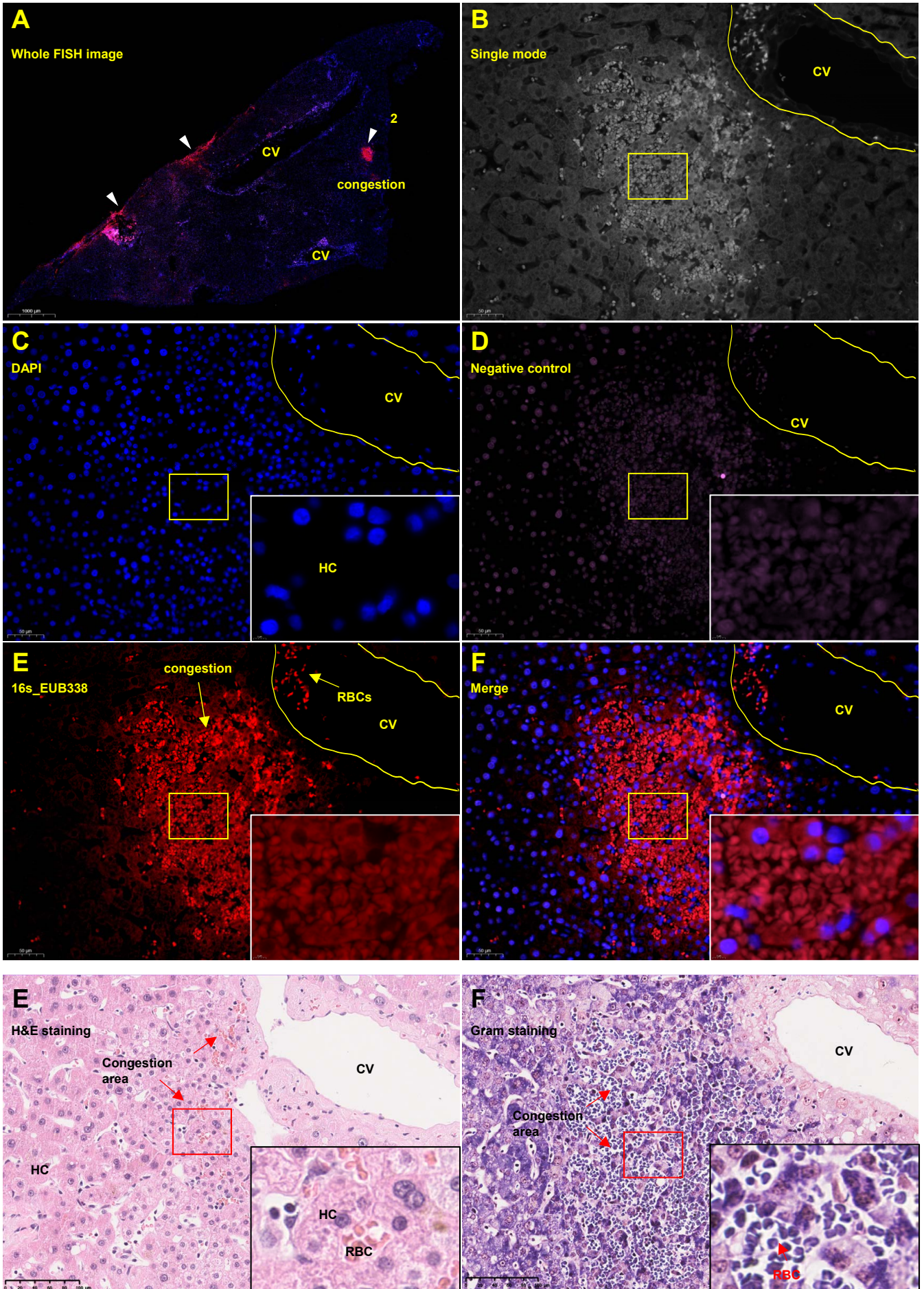


Figure S2. FISH analysis of a peri-tumor tissue revealing accumulation of RBCs in the congestion area. (A) The overview of the analyzed region. The white arrowheads indicate the congestion area. CV, central vein. (B) The single-mode of region 2 is labeled in (A). The CV area was illustrated with the yellow curve. The yellow box indicated the region that was enlarged. (C) DAPI staining. The region of the yellow box was enlarged in the lower right corner. HC, hepatocyte. (D) Staining of the complement probe. (E) EUB338 staining. The extensive staining was concentrated in the congestion region. The yellow arrows indicate the clusters of RBCs in the congestion area. (F) The merged image of different staining. (G) The H&E staining image. The red arrows indicate the RBCs in the congestion area. The red boxes indicate the same region as highlighted in the lower right corner. (H) The Gram staining image of a sequential tissue section of the corresponding region. Since the sections used for FISH and H&E staining were sequential sections, they could not be 100% merged to show the overlapping due to small changes between the sections.

Figure S3

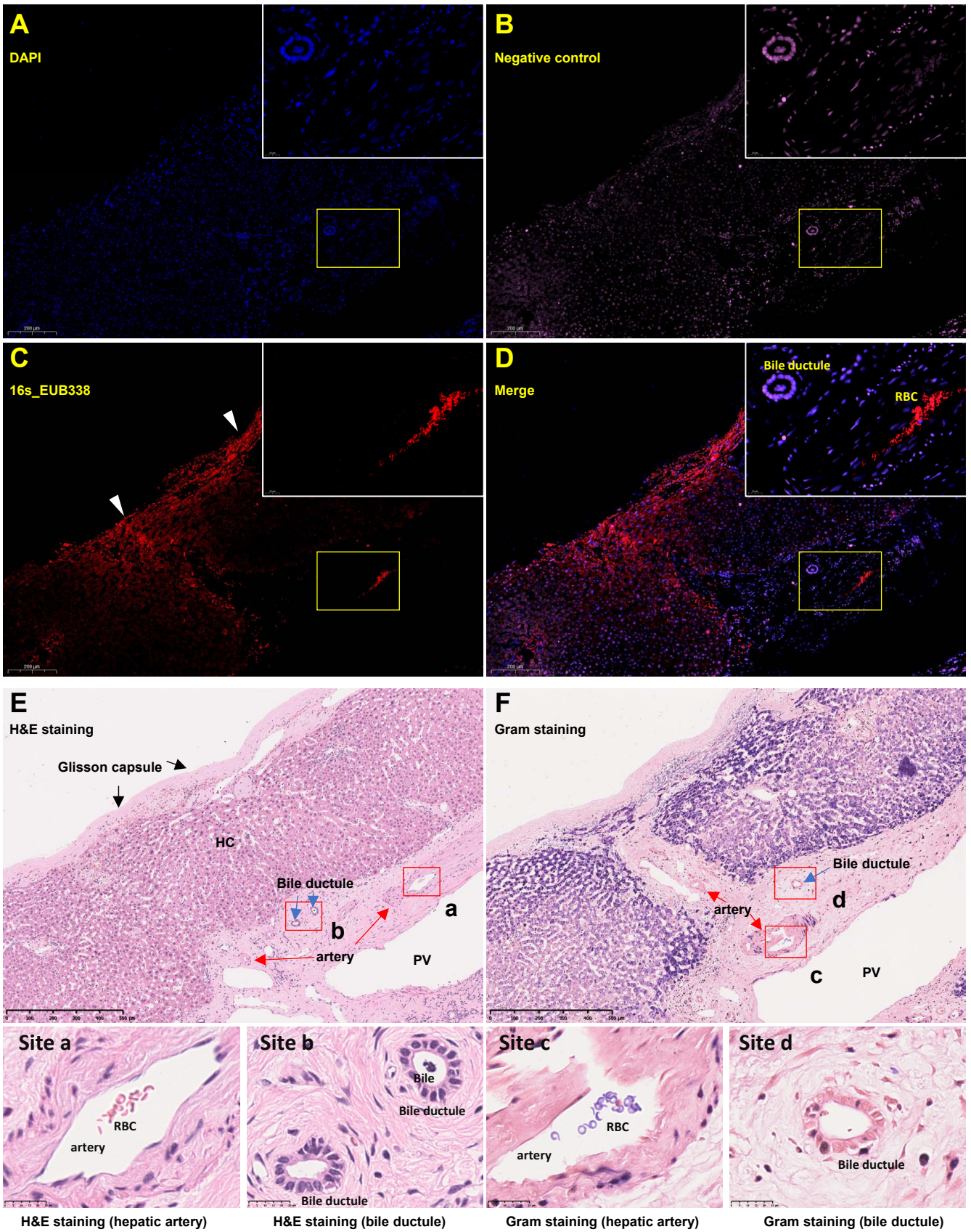
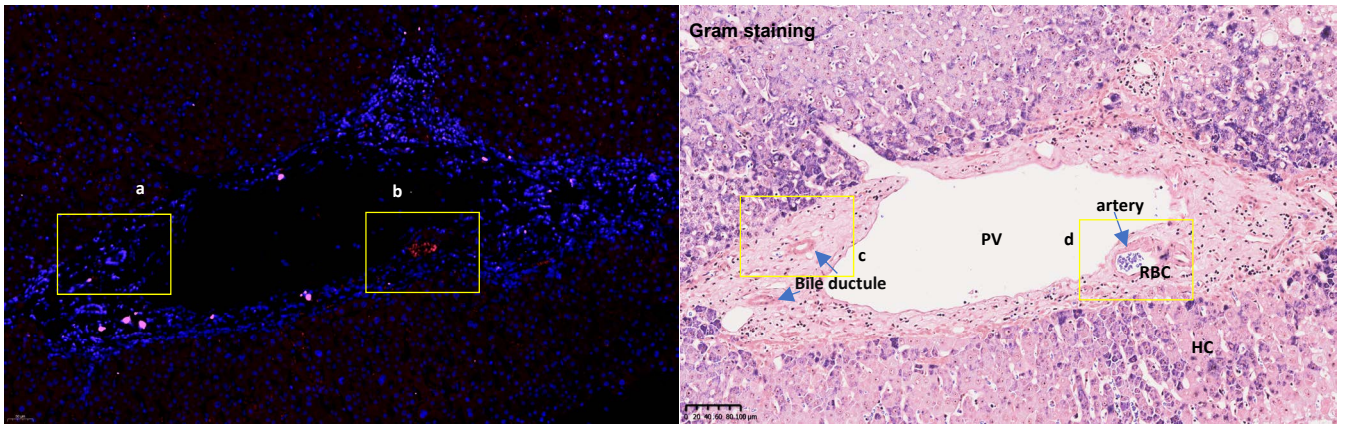


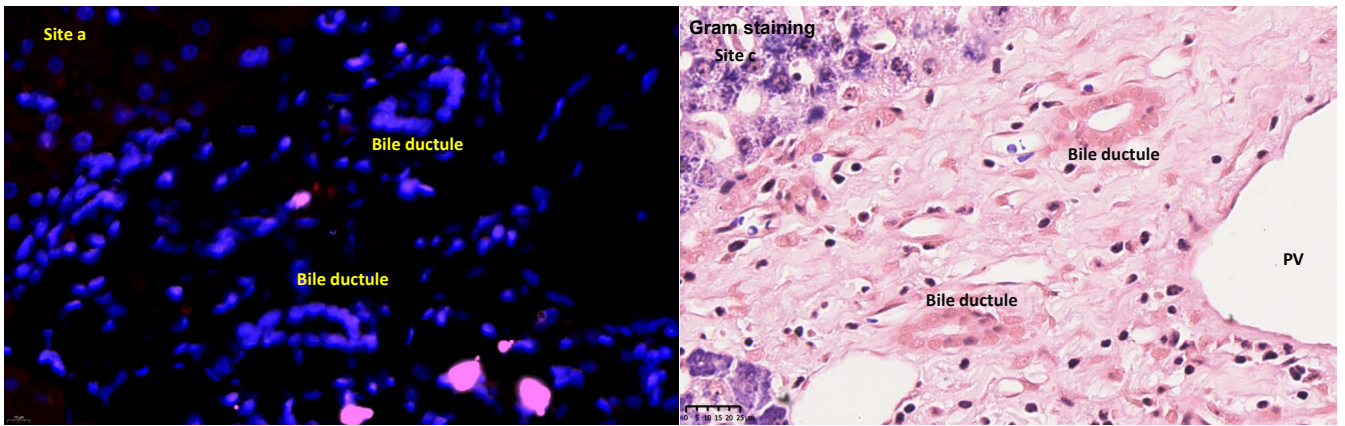
Figure S3. FISH staining of the peri-tumor tissue shows negative staining of the bile duct. (A) DAPI staining. The yellow boxes indicated the region enlarged in the top right corner. (B) Staining of the complement probe. (C) The EUB338 staining. The white arrowheads indicated the region of congestion. (D) The merged image of different staining. RBC, red blood cell. The magnified FISH image showed strong staining of RBCs in the artery and negative staining of bile ductule cells. (E) H&E staining. The black arrows indicate the Glisson capsule. The blue arrows indicate the bile ductule. The red arrows indicate the artery. The red boxes indicated the sites (a, b,) enlarged in the bottom panel. The magnified H&E staining images showed the artery (bottom left) or bile ductule (bottom right). HC, hepatocyte. PV, portal vein. (F) Gram staining image of a sequential tissue section of the corresponding region. The bile and the bile duct cells were negatively stained, compared with the intensive staining of RBCs in the artery lumen. The bottom panel shows the magnified Gram staining images of the artery and bile ductule (sites c and d). The RBCs in the artery lumen were stained blue-violet.

Figure S4

A



B



C

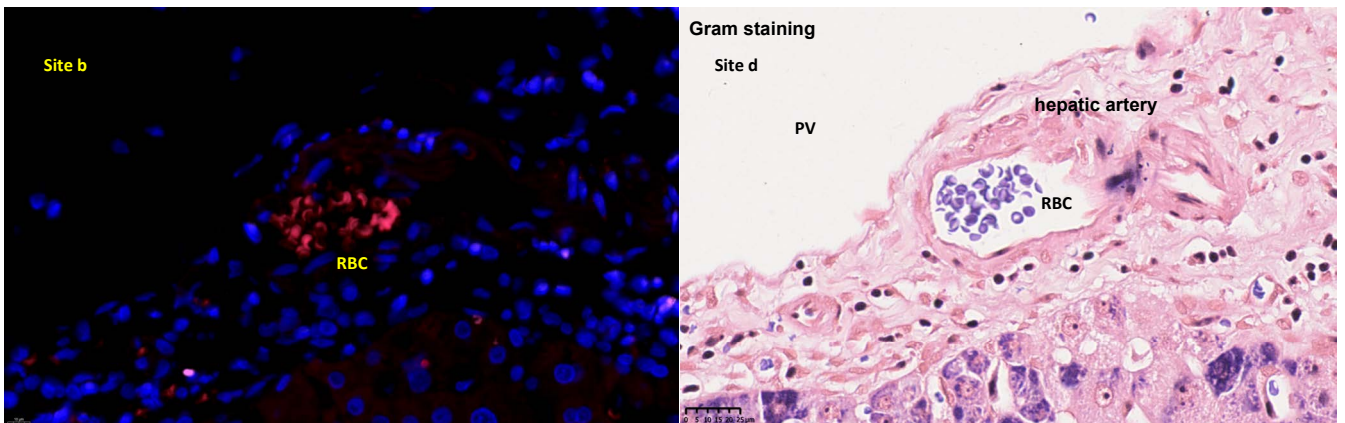


Figure S4. RBCs in the artery were positively stained, while bile ductules were negatively stained by a 16S rRNA FISH probe. (A) The merged image of EUB338, complement probe, and DAPI staining (left) and the image of Gram staining of the same region. Yellow boxes a and b were enlarged in (B) and (C), respectively. RBC, red blood cell. PV, portal vein. HC, hepatocytes. (B) The regions highlighted by yellow box a were enlarged. (C) The region highlighted by yellow box b was enlarged.

Figure S5

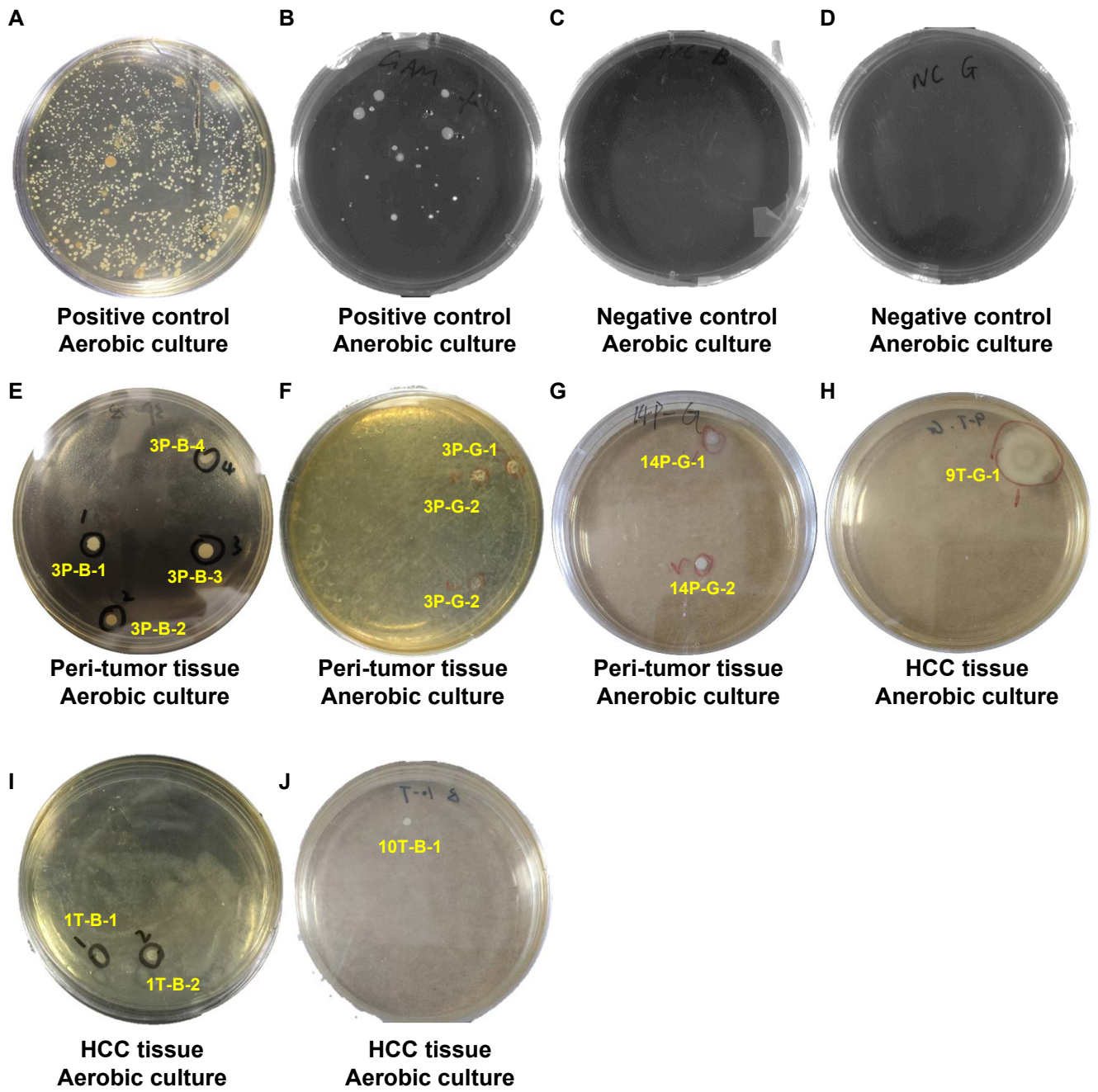


Figure S5. Viable bacteria were cultured from freshly prepared homogenates of liver cancer or peri-tumor tissues. (A) Positive control for aerobic culturing. The environmental microbiota was collected from the door handle using a dry cotton swab. The dish was incubated at 37°C. The picture was captured regularly on the desktop. (B) Positive control of anaerobic culturing. The dish was incubated at 37°C in a container used for anaerobic culturing. The picture was captured using the Fujifilm LAS3000 system. (C) and (D) The negative controls of aerobic or anaerobic culture using blank BHI or GAM medium. (E) and (F) Aerobic or anaerobic culturing of a peri-tumor tissue (3P). The labeled colonies were subjected to sub-culturing, DNA extraction, and DNA sequencing using the primers targeting the V3-V4 variable region of bacteria 16S rRNA gene. The tissue debris was left on the plate. The bacteria colonies were picked up with sterilized toothpicks for sub-culturing. (G) Anaerobic culturing of the peri-tumor tissue (14P). (H)-(J) Aerobic culturing of HCC tissues 9T, 1T, and 10T. The tissue debris was spread on the plate.

Figure S6

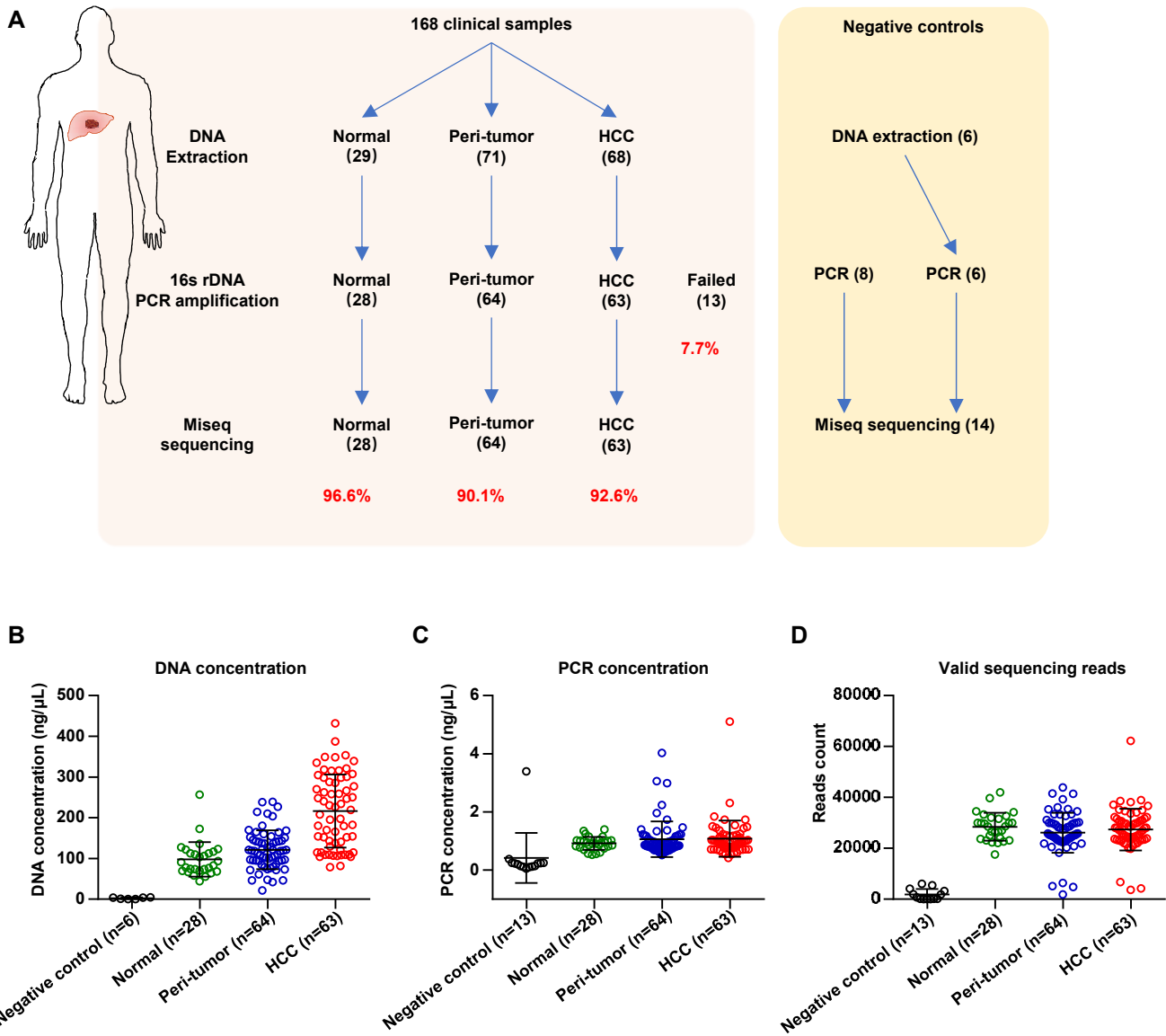


Figure S6. 16S rRNA sequencing of 168 normal, peri-tumor, and HCC tissues. (A)

Flowchart of the metataxonomic analysis of 168 liver tissues. Thirteen samples, including one normal, 7 peri-tumor, and 5 HCC tissues, failed in 16S rRNA gene PCR amplification. Negative control analyses of blank reagents (DNA extraction and PCR experiments) were also performed in parallel. Fourteen samples were subjected to Miseq sequencing. **(B)** DNA concentration ($\mu\text{g}/\mu\text{L}$) of the liver samples. **(C)** PCR products were quantified before being used for Miseq sequencing. **(D)** Valid reads of liver samples after filtering.

Figure S7

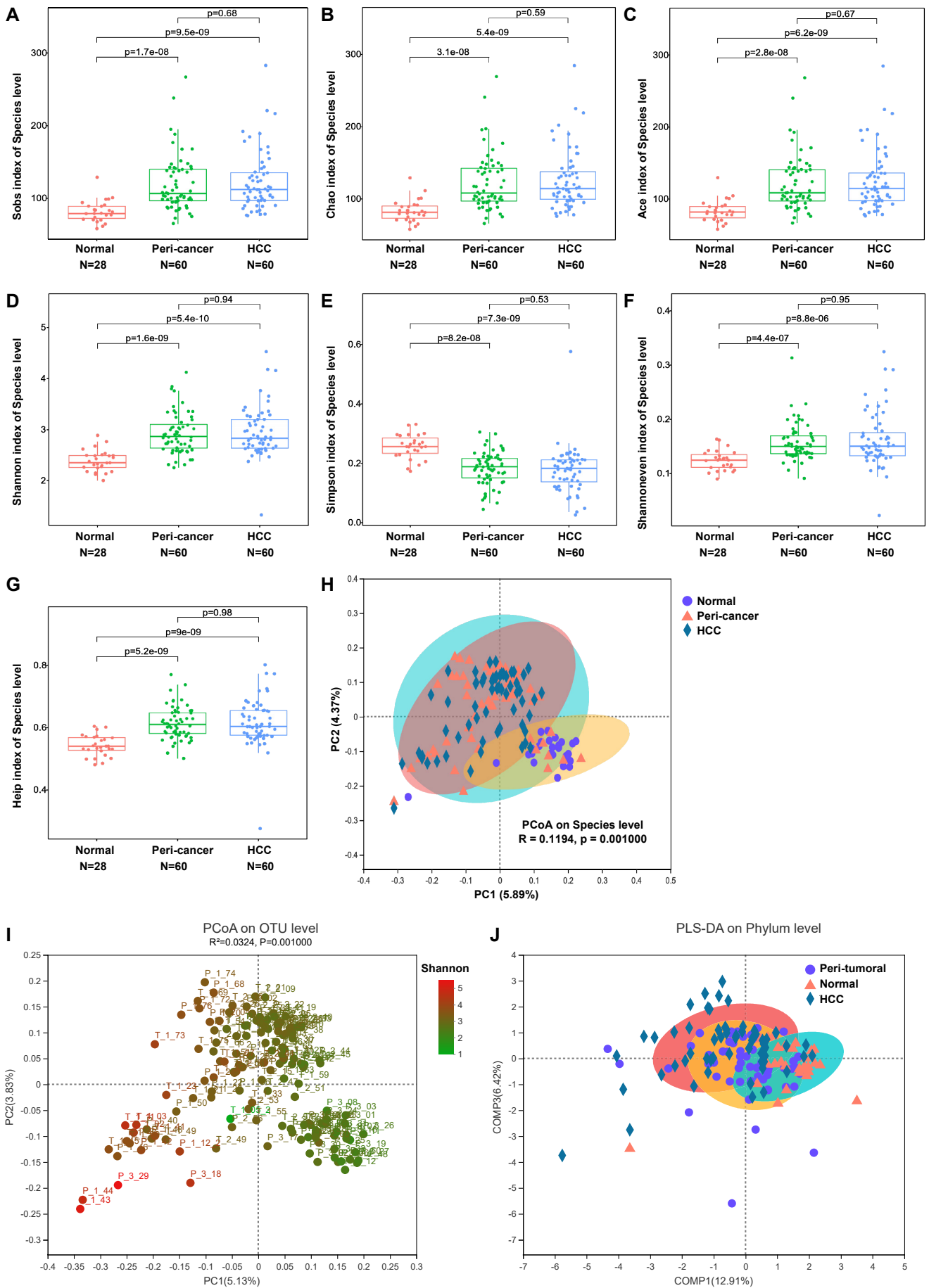


Figure S7. Alpha diversity of different liver microbiota at the species level. (A)-(C)

Richness indexes Sobs, Chao, and Ace were significantly increased in peri-tumor and HCC microhabitats. N indicates the sample number. (D) and (E) Diversity index Shannon was significantly increased while Simpson was significantly decreased in peri-tumor and HCC microhabitats compared with normal microhabitats. (F) and (G) Evenness indexes Shannoneven and Heip were significantly increased in peri-tumor and HCC microhabitats compared with normal microhabitats. (H) PCoA analysis on Species-level of normal, peri-tumor, and HCC subjects. The difference between the group was tested using the Adonis algorithm and the number of replacements was 999. (I) PCoA analysis on OTU level of normal, peri-tumor, and HCC microbiota. The samples were colored based on the Shannon index. The difference between the group was tested using the Adonis algorithm and the number of replacements was 999. (J) PLS-DA analysis revealed specific bacterial taxa which contribute to the major differences between different microhabitats. The ellipses indicated the confidence interval.

Figure S8



Figure S8. The structure and abundance of the microbiota in normal, peri-tumor, and HCC microhabitats. The structure and relative abundance of the normal, peri-tumor, and HCC microbiota at Order (**A**), Family (**B**), Genus (**C**), and Species (**D**) levels using OTU data.

Figure S9

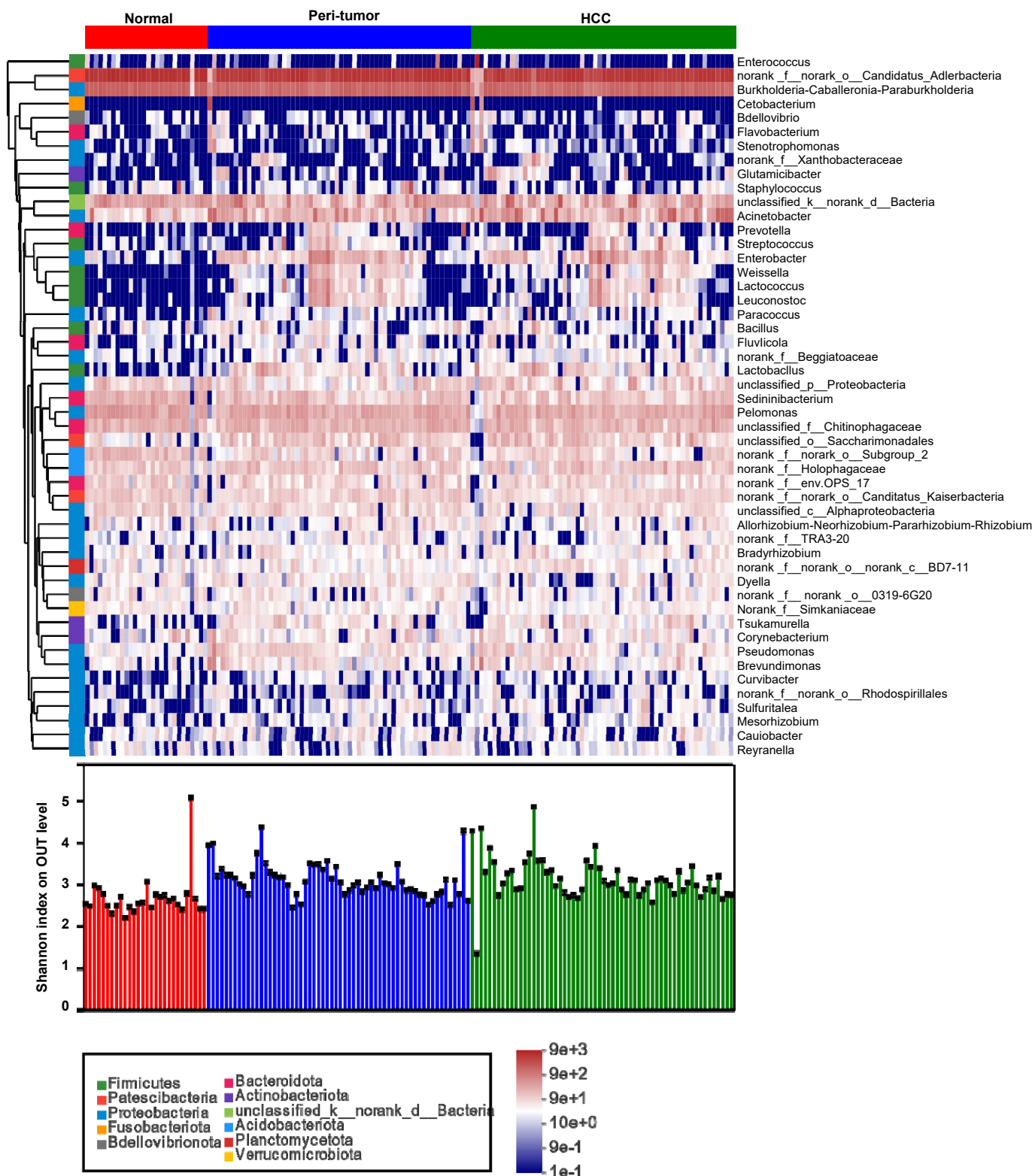


Figure S9. Hierarchical clustering of taxa at Genus level across different liver microbiota. Heatmap analysis revealed the relationship between identified Genera and the abundance across different liver microbiota. The Phyla of each Genus name was illustrated on the right side of the heatmap. The Shannon index on the OTU level was displayed at the bottom of the heatmap. The heatmap was generated using the Vegan package of R.

Figure S10

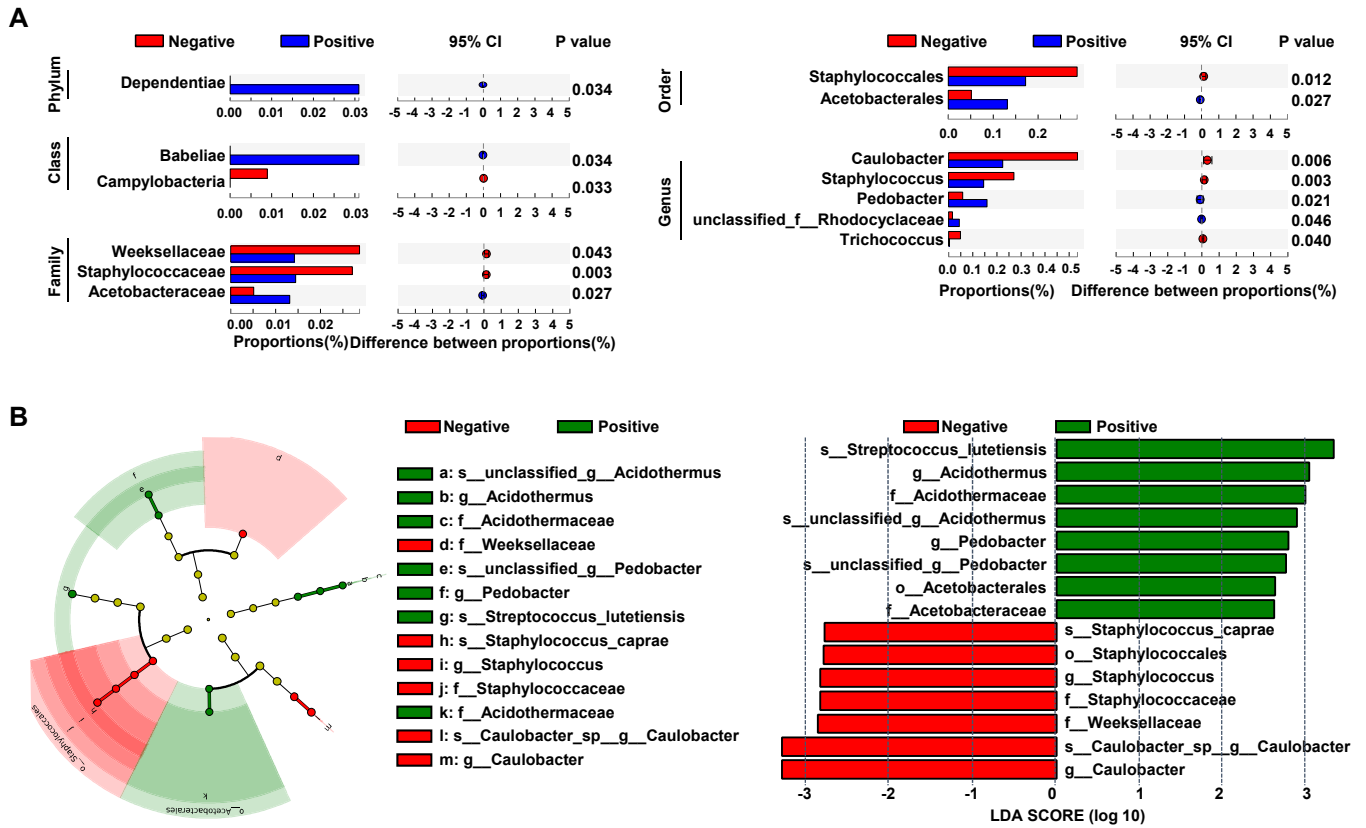


Figure S10. Bacterial taxa enriched in HBV-negative HCCs. A. The bacterial taxa showed a significant difference between HBV-negative (n=22) HccM and HBV-positive (n=49) HccM. The two-tailed p-value was calculated using Wilcoxon sum-rank test, and the p-value was adjusted using the FDR method. The 0.95% confidence intervals (CIs) were calculated using the bootstrap method. **B.** The cladogram and bar plot of LEfSe analysis showed the major taxa with the most abundance and a significant LDA score >2.5 in HBV-negative (n=22) HccM and HBV-positive (n=49) HccM.

Figure S11

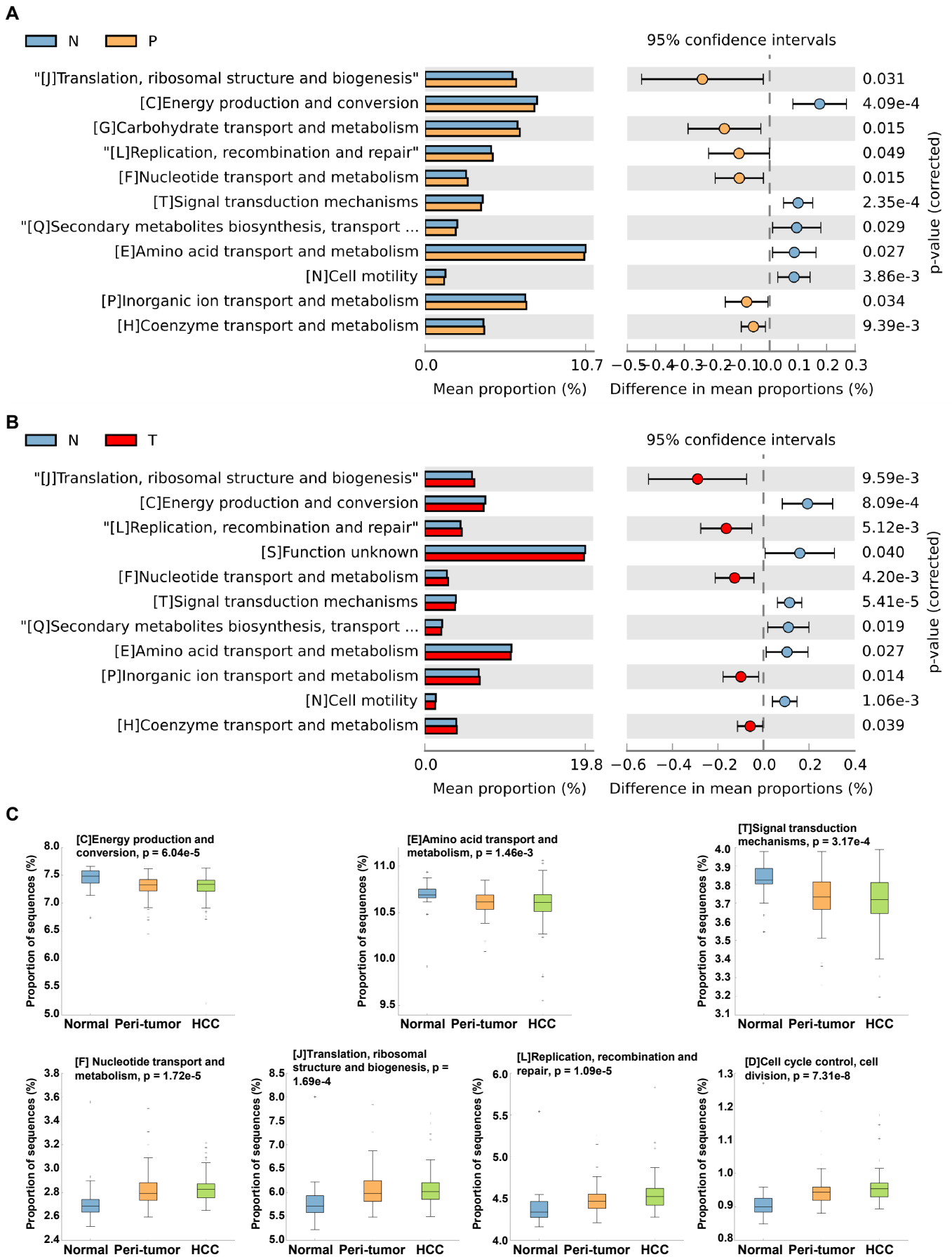


Figure S11. Inferred COG function by PICRUSt across different liver microbiota. (A) and **(B)** Comparison of the abundance of COG categories between the normal and peri-tumor microbiota or between the normal and HCC microbiota. The significance was calculated using a two-sided Welch's t-test and the 95% confidence interval of the difference in mean proportion was calculated using Welch's inverted method. **(C)**. The comparison of the abundance of selected COG categories in the normal, peri-tumor, and HCC microbiota. The comparison of function categories among the three groups was evaluated using the Kruskal–Wallis test followed by a post hoc Tukey-Kramer test (threshold = 0.95).

Figure S12

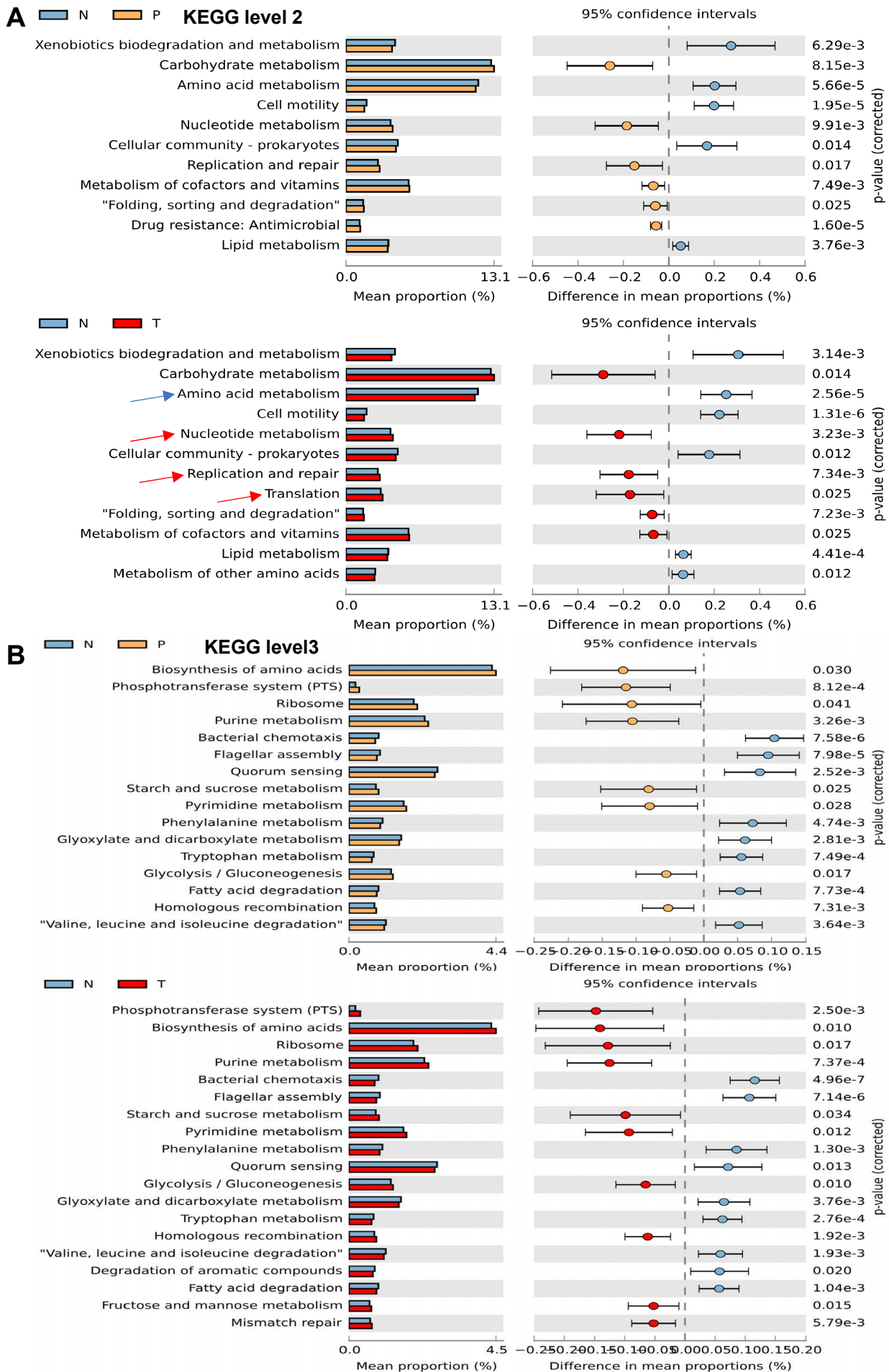


Figure S12. Inferred KEGG function by PICRUSt across different liver microbiota. (A)

Comparisons of the abundance of KEGG level2 categories between the normal and peri-tumor microbiota or between the normal and HCC microbiota. **(B)** Comparisons of the abundance of KEGG level3 categories between the normal and peri-tumor microbiota or between the normal and HCC microbiota.

Figure S13

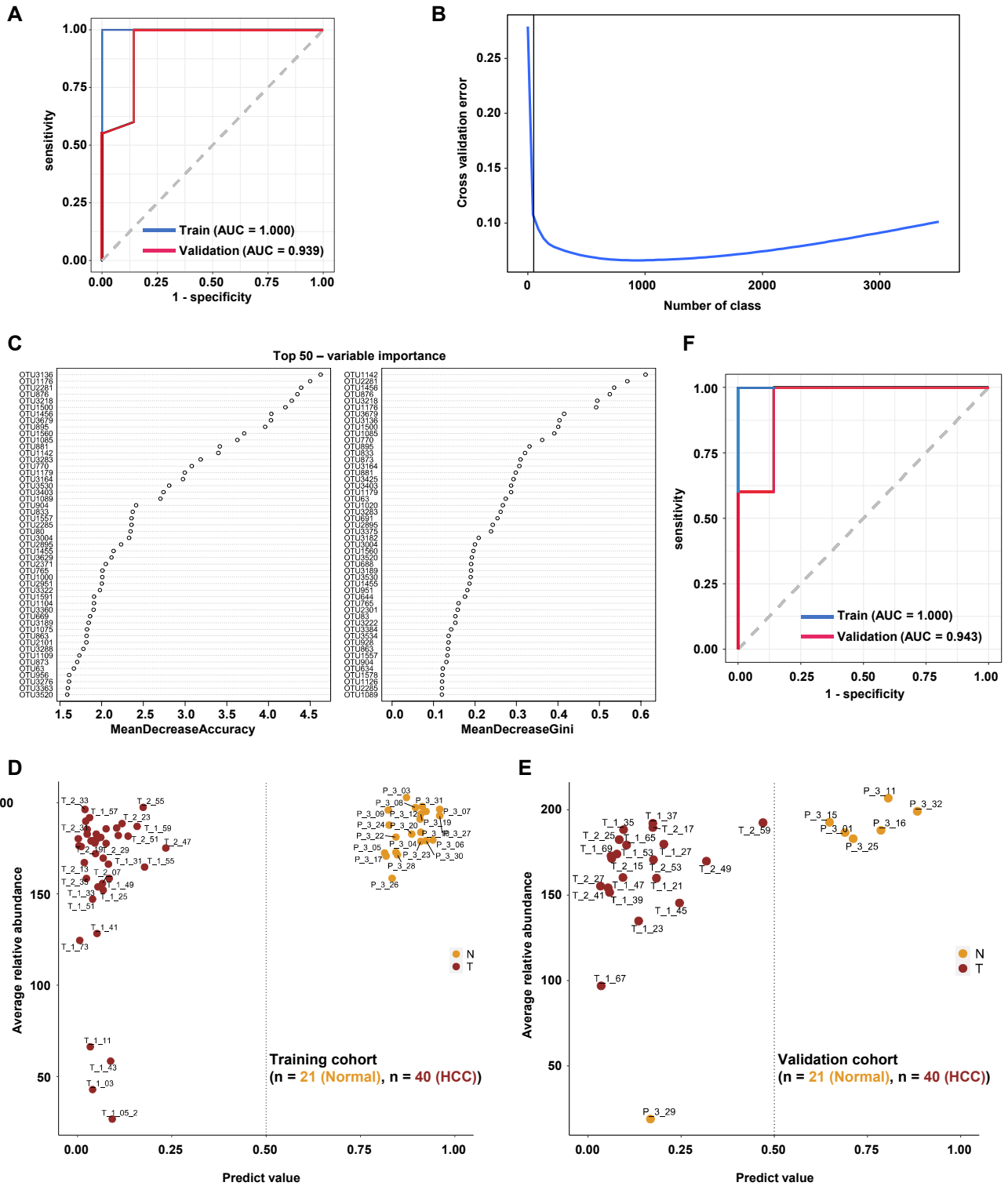


Figure S13. Microbial OTU-based diagnostic biomarker for HCC. (A) Receiver operating characteristic (ROC) curves of prediction performance for the OTU-based predictors of the RF model built on all OTUs. The blue and red curves indicated the performance of the model in the training cohort (HCC patient = 40, healthy = 20) and validation cohort (HCC patient = 20, healthy = 7), respectively. (B) The cross-validation error curve shows the 10-cross validation method to prioritize the top features used to build a simplified model. (C) The top 50 OTUs are prioritized by the 10-cross validation method. (D) Performance of the simplified model using the top 50 OTUs in the training cohort containing 40 HCC patients and 20 healthy individuals. (E) Performance of the simplified model using the top 50 OTUs in the validation cohort containing 20 HCC patients and 7 healthy individuals. (F) ROC curves of prediction performance for the OTU-based predictors of the simplified RF model built on top-50 OTUs. The blue and red curves indicated the performance of the model in the training cohort (HCC patient = 40, healthy = 20) and validation cohort (HCC patient = 20, healthy = 7), respectively.

References

1. Nejman D, Livyatan I, Fuks G, Gavert N, Zwang Y, Geller LT, Rotter-Maskowitz A, Weiser R, Mallel G, Gigi E, Meltser A, Douglas GM, Kamer I, Gopalakrishnan V, Dadosh T, Levin-Zaidman S, Avnet S, Atlan T, Cooper ZA, Arora R, Cogdill AP, Khan MAW, Ologun G, Bussi Y, Weinberger A, Lotan-Pompan M, Golani O, Perry G, Rokah M, Bahar-Shany K, Rozeman EA, Blank CU, Ronai A, Shaoul R, Amit A, Dorfman T, Kremer R, Cohen ZR, Harnof S, Siegal T, Yehuda-Shnaidman E, Gal-Yam EN, Shapira H, Baldini N, Langille MGI, Ben-Nun A, Kaufman B, Nissan A, Golan T, Dadiani M, et al. 2020. The human tumor microbiome is composed of tumor type-specific intracellular bacteria. *Science* 368:973-980.
2. Pfeiffer S, Pastar M, Mitter B, Lippert K, Hackl E, Lojan P, Oswald A, Sessitsch A. 2014. Improved group-specific primers based on the full SILVA 16S rRNA gene reference database. *Environ Microbiol* 16:2389-407.
3. Sookoian S, Salatino A, Castano GO, Landa MS, Fijalkowky C, Garaycochea M, Pirola CJ. 2020. Intrahepatic bacterial metataxonomic signature in non-alcoholic fatty liver disease. *Gut* 69:1483-1491.