

# 1 uniLIVER: a Human Liver Cell Atlas for Data-Driven 2 Cellular State Mapping

3

4 Yanhong Wu<sup>1,#</sup>, Yuhan Fan<sup>1,#</sup>, Yuxin Miao<sup>1,#</sup>, Yuman Li<sup>1</sup>, Guifang Du<sup>2,3</sup>, Zeyu Chen<sup>1</sup>,  
5 Jinmei Diao<sup>2,3</sup>, Yu-Ann Chen<sup>2,3</sup>, Mingli Ye<sup>4</sup>, Renke You<sup>4</sup>, Amin Chen<sup>4</sup>, Yixin Chen<sup>1</sup>,  
6 Wenrui Li<sup>1</sup>, Wenbo Guo<sup>1</sup>, Jiahong Dong<sup>2,3</sup>, Xuegong Zhang<sup>1,5</sup>, Yunfang Wang<sup>2,3,\*</sup>, Jin  
7 Gu<sup>1,\*</sup>

8

9 <sup>1</sup>. MOE Key Lab of Bioinformatics, BNRIST Bioinformatics Division, Department of  
10 Automation, Tsinghua University, Beijing, China

11 <sup>2</sup>. Hepato-Pancreato-Biliary Center, Beijing Tsinghua Changgung Hospital, Tsinghua  
12 University, Beijing, China

13 <sup>3</sup>. Clinical Translational Science Center, Beijing Tsinghua Changgung Hospital,  
14 Tsinghua University, Beijing, China

15 <sup>4</sup>. Fuzhou Institute of Data Technology, Fuzhou, China

16 <sup>5</sup>. Center for Synthetic and Systems Biology, School of Life Sciences and School of  
17 Medicine, Tsinghua University, Beijing, China

18

19 <sup>#</sup> These authors are equally contributed to this work.

20 <sup>\*</sup> Correspondences should be addressed to Jin Gu ([jgu@tsinghua.edu.cn](mailto:jgu@tsinghua.edu.cn)) and Yunfang  
21 Wang ([wyfa02717@btch.edu.cn](mailto:wyfa02717@btch.edu.cn)).

22

23

# Abstract

The liver performs several vital functions such as metabolism, toxin removal and glucose storage through the coordination of various cell types. The cell type compositions and cellular states undergo significant changes in abnormal conditions such as fatty liver, cirrhosis and liver cancer. As the recent breakthrough of the single-cell/single-nucleus RNA-seq (sc/snRNA-seq) techniques, there is a great opportunity to establish a reference cell map of liver at single cell resolution with transcriptome-wise features. In this study, we build a unified liver cell atlas uniLIVER by integrative analyzing a large-scale sc/snRNA-seq data collection of normal human liver with 331,125 cells and 79 samples from 6 datasets. Besides the hierarchical cell type annotations, uniLIVER also proposed a novel data-driven strategy to map any query dataset to the normal reference map by developing a machine learning based framework named LiverCT. Applying LiverCT on the datasets from multiple abnormal conditions (1,867,641 cells and 439 samples from 12 datasets), the alterations of cell type compositions and cellular states were systematically investigated in liver cancer.

# 1 Main

2 The liver is a major metabolic organ, which performs many essential physiological  
3 functions, including toxin removing, albumin and bile production, glucose and amino  
4 acid processing, and vitamin storage, etc. For humans, hepatocytes occupy about 80%  
5 liver volume and cholangiocytes, immune cells and stromal cells consist of the  
6 majority of the remaining part. These cells are organized into hexagonal hepatic  
7 lobules as the basic functional units of liver. As the recent breakthrough of the single-  
8 cell/single-nucleus RNA-seq (sc/snRNA-seq) techniques<sup>1,2</sup>, there is a great  
9 opportunity to establish a reference cell map of liver at single cell resolution with  
10 transcriptome-wise features. Besides, the reference map is very useful for studying the  
11 altered cell type compositions and cellular states under diverse physiological and  
12 pathological conditions in liver, such as acute injury, virus infection, cirrhosis and  
13 cancer.

14 In this study, we collected 79 normal human liver samples from 6 datasets<sup>3-7</sup>, and  
15 439 abnormal or disease samples from 12 datasets<sup>4,8-18</sup>. Based on this collection, we  
16 hierarchically annotated 63 cell types/subtypes and the hepatocytes in 4 different  
17 lobular regions/zones of normal liver, and then constructed an integrated and data-  
18 driven human liver cell atlas uniLIVER. Beyond the traditional cell atlases mainly  
19 providing comprehensive cell type annotations, uniLIVER also aims at establishing a  
20 novel data-driven strategy to map any query dataset to the normal reference map.

21 Analogy to the genome sequence mapping, we proposed a concept for cell type or  
22 cellular state mapping: the query cells are computationally “mapped” to the reference  
23 map based on gene expression features. Those cells whose gene expressions are  
24 dissimilar to any annotated cell subtype in the reference are defined as “variant” state  
25 cells (analogy to nucleotide variants in genome sequence analysis). The “variant”  
26 states are broadly categorized into two types, the “deviated” state and the  
27 “intermediate” state: the “deviated” state means that the gene expressions of the query

1 cells are shifted from a single cell type, and the “intermediate” state means that the  
2 gene expressions of the query cells located between any two reference cell types. We  
3 developed LiverCT, a machine learning based liver Cell-Type mapping method using  
4 the annotated reference data, to achieve the task of cellular state mapping in liver.

5 We applied LiverCT on the collected abnormal liver datasets. Results show that  
6 almost all types of cells are strongly “deviated” from their normal states in  
7 hepatocellular carcinoma (HCC) and the deviated scores of T cells are positively  
8 correlated with the stress response pathway signatures. Interestingly, the results also  
9 show that the hepatic stellate cells and granulocytes (mainly neutrophils) are highly  
10 deviated in adjacent non-tumor tissues. For the “intermediate” state analysis, it was  
11 found that the cancer cells with high intermediate scores have strongly up-regulated  
12 glycolysis and hypoxia pathways. Also, the up-regulated genes of those cells are  
13 significantly overlapped with poor prognosis genes. Another interesting task is to  
14 analyze the zonation tendency of the HCC tumor cells by mapping them to the  
15 hepatocyte states in different lobular zones. We found that the zonation tendency is  
16 highly associated with the expression of multiple malignant signatures and the  
17 composition of multiple immune and stromal cell types.

18 uniLIVER tends to establish a new framework of cell atlas by introducing  
19 machine learning into the traditional data portal only design. Both the reference cell  
20 map (as three portraits similar to hECA<sup>19</sup>) and the computationally cellular state  
21 mapping tool LiverCT are freely available via a web-based database.

22

23

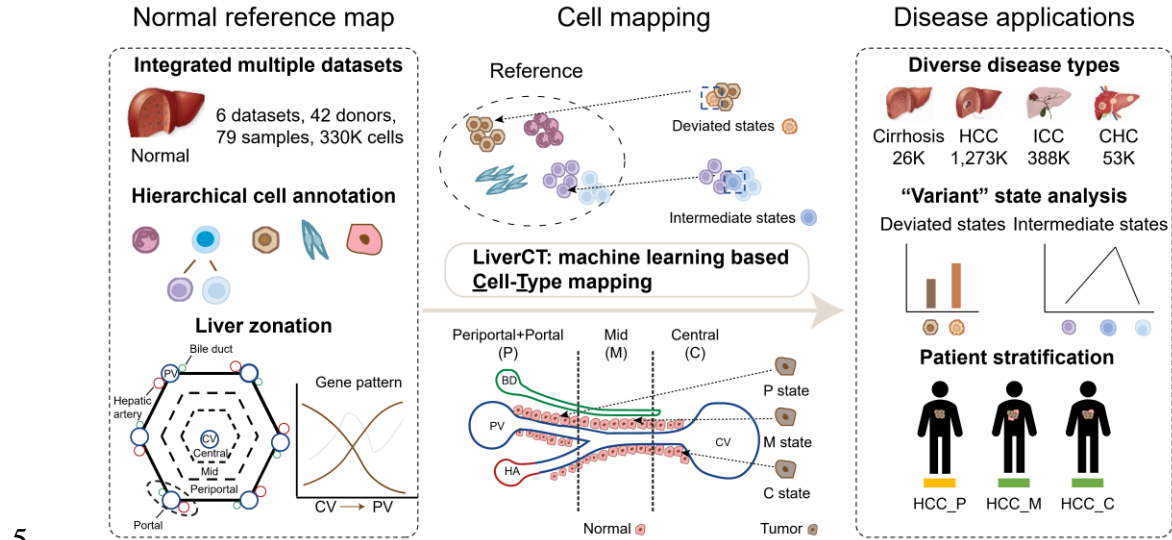
# 1 Results

## 2 Overview of the human liver cell atlas uniLIVER

3 We collected scRNA-seq data from 6 datasets, which include 79 samples from 42  
4 donors (Table 1). After stringent quality control, 331,125 cells were integrated and  
5 annotated hierarchically to form the normal reference map. Population and clinical  
6 information, such as age and gender, were also collected if available (Extended Data  
7 Fig. 1a, Supplementary Table 1). Besides, we also curated 1,867,641 cells from  
8 diverse abnormal or disease samples including cirrhosis, hepatocellular carcinoma  
9 (HCC), intrahepatic cholangiocarcinoma (ICC), combined hepatocellular and  
10 intrahepatic cholangiocarcinoma (CHC) and some liver metastases (Table 1).

11 Then, the atlas tends to “map” the cells from disease samples to the normal  
12 reference map based on gene expression similarities. To find the “variant” state cells  
13 in disease samples, we developed LiverCT, a machine learning based Cell-Type  
14 mapping method that can distinguish “deviated” states and “intermediate” states and  
15 can also classify the abnormal hepatocytes in HCC into P-like (P: Portal/Periportal)  
16 state, M-like (M: Mid) state and C-like (C: Central) state (Fig. 1). One unique aspect  
17 of LiverCT is that it adopts the concept of genomic variant analysis to identify and  
18 analyze the “variant” states of cells based on their expression patterns. Firstly,  
19 LiverCT embedded the cells in the normal reference map into a latent space. In the  
20 latent space, we trained a hierarchical classifier using ensemble learning to predict the  
21 cells’ type labels. This was followed by another one-class classifier to identify the  
22 margin of a normal cell type and a one-vs-one classifier to identify the boundary  
23 between any two cell types. Further, if a cell was predicted as hepatocyte, we utilized  
24 an additional classifier to distinguish its lobular zonation tendency. When applied to  
25 query datasets, LiverCT embedded the query cells into the pre-calculated latent space  
26 of the normal reference by scArches<sup>20</sup> to remove the batch effects. After the cell type

1 prediction, a “deviated” score was calculated to measure the degree of deviation from  
2 the corresponding normal cell type and an “intermediate” score was calculated to  
3 measure the degree to which the cell was in the middle of top two predicted cell types.  
4 (Methods, Extended Data Fig. 1b).



5  
6 **Fig. 1| uniLIVER overview.** The normal reference by hierarchically annotating the cell types and  
7 zonation tendency via an integrative analysis of multiple sc/snRNA-seq datasets (left). The  
8 machine learning based Cell-Type mapping method named LiverCT can identified the cells with  
9 deviated or intermediate states and calculate the hepatocyte zonation tendency of any query  
10 dataset (middle). LiverCT was applied for liver cancer and other abnormal conditions (rights).

11

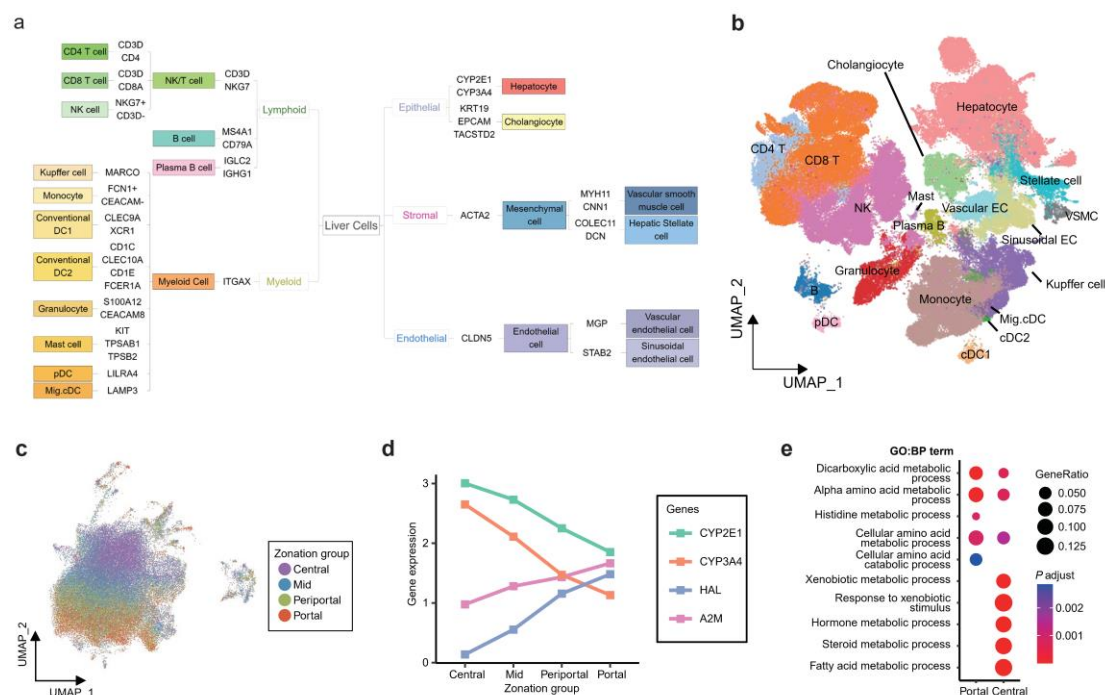
# 1 Hierarchical annotations of the normal liver cells

2 A unified normal reference map needs in-depth and harmonized annotations.  
 3 Currently, there exists some inconsistency in cell type definition across different  
 4 studies. To harmonize the cell type labels from different datasets, firstly we built a  
 5 unified hierarchical annotation framework (uHAF<sup>19</sup>) for liver (Fig. 2a). The major cell  
 6 type annotations provided by the original references were harmonized into 8 major  
 7 cell types (Level 1) in uniLIVER (Supplementary Table 2). In cases where the  
 8 annotations were not provided, the datasets were annotated manually. Using the Level  
 9 1 annotations as prior knowledge, we fine-tuned scANVI<sup>21</sup>, a tool that proved to be  
 10 one of the top-performing integration methods<sup>22</sup>, to remove the batch effects between  
 11 different studies (Fig. 2b). Under each of these major cell types, we further performed  
 12 un-supervised graph-based clustering for in-depth annotations, which generated 17  
 13 stable cell types (Level 2) (Fig.2c).

14 The hepatic lobule is the basic unit for liver function, with a central vein (CV) in  
 15 the middle and portal vein in the six corners (PV). The hepatocytes have different  
 16 states and functions along the CV to PV axis. The lobule can be roughly divided in  
 17 four regions: central, middle, portal and peri-portal zones<sup>5</sup> (in a few studies the portal  
 18 and peri-portal zones are combined as a single zone<sup>23,24</sup>). We annotated hepatocyte  
 19 zonation utilizing the gene signatures obtained from spatial transcriptomes (Methods,  
 20 Fig. 2d, Extended Data Fig. 2a-c, Supplementary Table 3). The annotated zonation  
 21 patterns can be validated by the expressions of multiple canonical marker genes:  
 22 *CYP2E1* and *CYP3A4* gradually decreased along the CV-PV axis while *HAL* and *A2M*  
 23 exhibited an opposite trend, highly consistent with previous studies<sup>5,25</sup> (Fig. 2e).  
 24 Functional enrichment analysis shows that the up-regulated genes of the hepatocytes  
 25 annotated as in central region were enriched in xenobiotic metabolism pathways and  
 26 the result of portal region was enriched in amino acid processing (Fig. 2f, Extended  
 27 Data Fig. 2d), which was also consistent with the well-established cellular functions

1 of the hepatocytes in different zones<sup>26</sup>.

2 This large scale of data collection enables us to primarily investigate the impact  
3 of population variables (e.g. gender and age) on the cell type compositions and subtle  
4 gene expression variations in a same cell type. We found no significant difference in  
5 the proportion of cell types between genders (Extended Data Fig. 3a, b). As the age  
6 increases, the proportion of different cell types underwent complex changes  
7 (Extended Data Fig. 3c, d). The effects of these variables on gene expression  
8 variations were also explored using generalized linear mixed models<sup>27-29</sup>. We found  
9 that the biological processes, including cellular response to ions (*MT1A* and *MT1E*),  
10 material transport (for example, cholesterol and sterol transport including *APOC2* and  
11 *APOM*), homeostasis maintenance (such as cholesterol, sterol, and lipid homeostasis  
12 including *AKR1C1*, *APOC2*, and *APOM*) and immune response (*C1A*, *CFHR2*,  
13 *RARRES2*) were significantly downregulated in monocytes within the elder  
14 population (Extended Data Fig. 3e).



16 **Fig. 2 | Construction of the normal reference map.** a, The Level 1 and Level 2 uHAF tree of  
17 liver. b, The UMAP visualization of the normal reference map. c, The UMAP of hepatocytes  
18 (colors are zonation annotations). d, Gene expression of canonical markers *CYP2E1*, *CYP3A4*,



1 *HAL* and *SBDS* across the CV-PV axis. e, The significantly enriched biological process (GO:BP  
2 terms with Benjamini–Hochberg-adjusted  $P < 0.05$ ) of the genes upregulated in portal and central  
3 regions, respectively.

4

# 1 **Deviated state analysis identifies diverse disease-associated**

## 2 **cellular states**

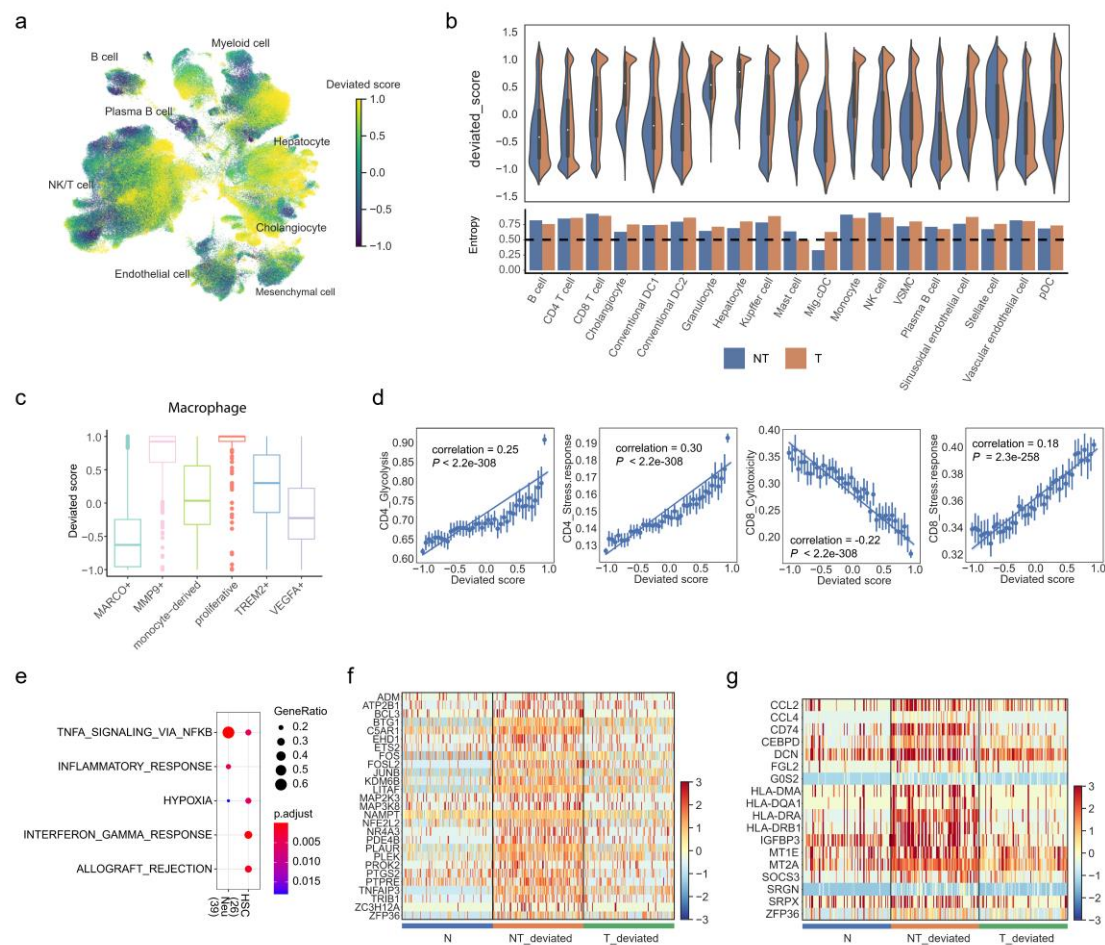
3 Compared to the normal liver, the cells in disease often exhibit certain cellular states  
 4 “deviation”, which provide potential targets for treatment. To elucidate the state  
 5 deviation extent of the cells under disease conditions, we developed LiverCT with a  
 6 supervised and hierarchical ensemble learning framework to calculate a quantitative  
 7 deviated score based on the normal reference map. The performance of LiverCT was  
 8 firstly validated on predicting the cell type labels in the normal reference map  
 9 (Methods, Extended Data Fig. 4a).

10 The collected disease datasets were mapped to the normal reference by LiverCT  
 11 (Fig. 3a). We observed higher deviated scores in the cells from tumor (T) tissues  
 12 compared to adjacent non-tumors (NT) in most cell types. In primary tumors,  
 13 hepatocytes, cholangiocyte and granulocytes changed significantly, indicating strong  
 14 deviation to the normal reference (Fig. 3b, Extended Data Fig. 4b). Hepatocytes and  
 15 cholangiocytes were parenchymal cells that were prone to oncogenic transformations,  
 16 making them distinct from normal tissues. A recent study also reported that  
 17 neutrophils (major cell group of granulocytes) in liver cancer tissues exhibited  
 18 significant gene expression changes compared to non-tumor tissues<sup>16</sup>.

19 Furthermore, we tested the correctness of the deviated states on a HCC dataset  
 20 published recently by Lu *et al.*<sup>9</sup>. They found the MMP9<sup>+</sup> macrophages to be tumor-  
 21 associated macrophages. It was observed that these MMP9<sup>+</sup> macrophages had high  
 22 deviated scores, with almost all of the cells scored positive (Fig. 3c). Also, there was  
 23 no proliferative macrophage observed in the normal reference and this group got the  
 24 highest deviated score in the HCC dataset. Besides, other known tumor-associated  
 25 cell types such as regulatory T cells, pro-metastatic hepatocytes can also be  
 26 successfully found with high deviated scores (Extended Data Fig. 4c). Taken together,  
 27 above results showed that LiverCT can accurately retrieve the deviated state cells.

1 Tumor-infiltrating T cells have paved a novel way for tumor therapy<sup>30</sup>. We  
 2 investigated the correlation between the function of T cells and their deviated scores<sup>31</sup>.  
 3 The deviated scores of CD4 T cells have a positive correlation with the signatures of  
 4 glycolysis and stress response. The scores of CD8 T cells show also a positive  
 5 correlation with stress response signature, but a negative correlation with cytotoxicity  
 6 (Supplementary Table 4, Fig. 3d). A recent pan-cancer study observed a strong  
 7 association between the stress response and immunosuppression<sup>31</sup>, suggesting that the  
 8 deviated score of T cells is a possible alternative indicator for immunotherapy  
 9 response.

10 The adjacent non-tumor tissue (NT) presents a unique state between the normal  
 11 and tumor<sup>32</sup> and may provide additional information of the oncogenic transformation  
 12 and recurrence<sup>33</sup>. Based on the cell type mapping results by LiverCT, we observed  
 13 that the granulocytes, hepatocytes, and stellate cells had the highest number of the  
 14 cells in deviated states. To unveil the NT's unique expression features, we conducted  
 15 differential expression analyses and focused on the consistently up-regulated genes in  
 16 NT's deviated cells. The up-regulated genes in both granulocytes and stellate cells  
 17 were significantly enriched in the TNF- $\alpha$  signaling pathway (Fig. 3e-g,  
 18 Supplementary Table 5), consistent with a pan-cancer study that observed NT-specific  
 19 TNF- $\alpha$  signaling pathway activation<sup>32</sup>. Collectively, the deviated state analysis  
 20 identified diverse disease-associated cellular states and illustrated the most susceptible  
 21 cell types in disease.



**Fig. 3 | Deviated state analysis reveals the changes in adjacent non-tumor and tumor samples.** a, The UMAP of the disease datasets with deviated scores. b, The deviated score distribution of each cell type at Level 2 (upper panel) and the entropy of the deviated state cells (below panel). c, Deviated scores of the macrophages with their original labels in Lu *et al.* datasets. d, Regplot of the deviated scores and T cell functional signatures for both CD4 and CD8 T cells. Pearson correlations were used to assess the associations. E, Enriched cancer hallmarks in the genes consistently upregulated in “deviated” neutrophils (Neu) and hepatic stellate cells (HSC) in adjacent non-tumor samples (adjusted P-value < 0.05). f,g, Heatmap of the expressions of the upregulated genes in the enriched hallmarks (Neu, f; HSC, g).

# 1 Intermediate state analysis reveals a population of tumor 2 cells associated with poor prognosis in HCC

3 Intermediate or transition cellular states are widely induced in tumors. For instance, in  
4 human melanoma<sup>34</sup>, a transitional CD8 state is observed, and in peripheral  
5 neuroblastic tumors, an intermediate state is observed between adrenergic and  
6 mesenchymal neuroblasts<sup>35</sup>. To find the cells with intermediate states, LiverCT can  
7 also calculated an “intermediate” score for query cells (Fig. 4a). In the collected liver  
8 cancer datasets, cells between CD4-CD8, Mono-Macro, HSC-VSMC and LSEC-VEC  
9 exhibited high intermediate states ratio (Fig. 4b).

10 Also, we observed that many tumor cells have high intermediate scores between  
11 the hepatocyte-cholangiocyte pair. Liver cancers, including HCC, ICC, and CHC,  
12 exhibited significant heterogeneity and can arise from diverse origins<sup>36</sup>. It has been  
13 reported that ICC could potentially arise from biliary-like cells that undergo trans-  
14 differentiation from hepatocytes, as well as from hepatic progenitor cells in addition  
15 to mature cholangiocytes<sup>37</sup>. So, we then tend to study the characteristics of those  
16 tumor cells with intermediate states.

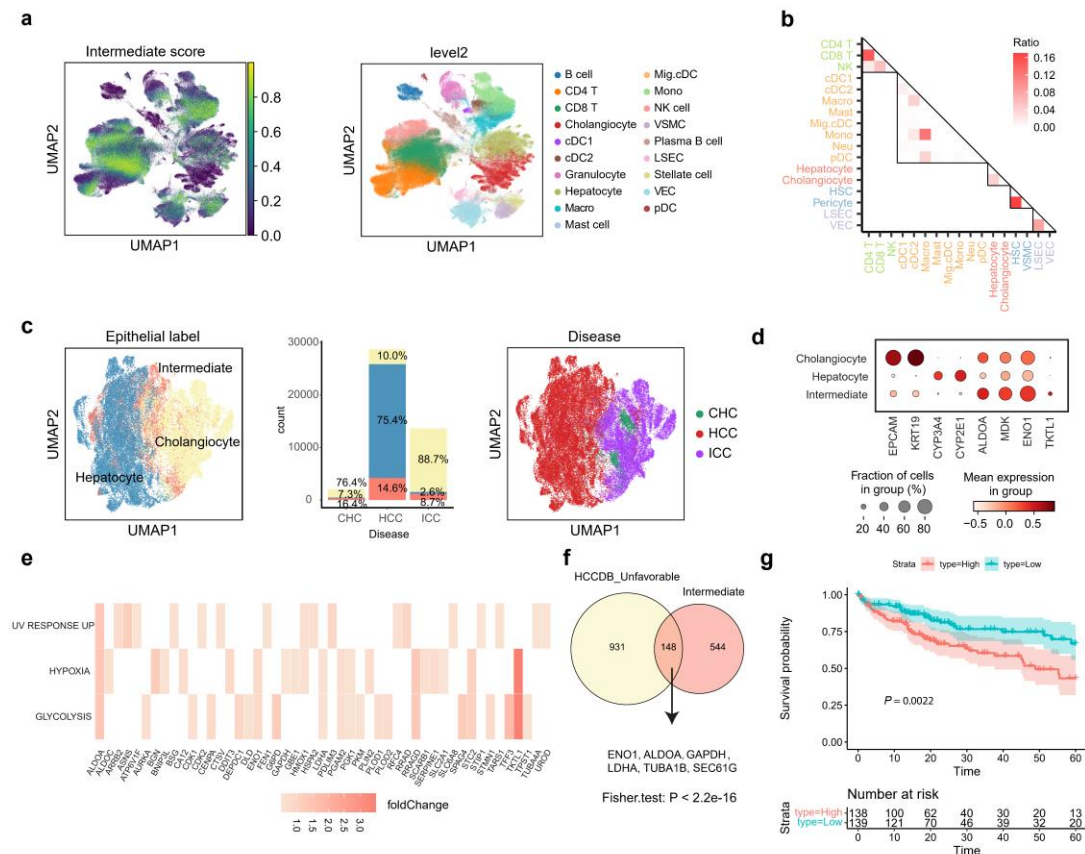
17 LiverCT classified the malignant epithelial cells from tumor samples into  
18 hepatocyte, cholangiocyte and intermediate states (Methods, Fig. 4c). Notably, of the  
19 three types of liver cancer, the ratio of intermediate states in CHC was the highest,  
20 followed by HCC. Cholangiocytes exhibited high expression of *EPCAM* and *KATI9*,  
21 while hepatocytes expressed *CYP3A4* and *CYP2E1*. The intermediate states, however,  
22 under-expressed both these cell type markers. (Fig. 4d, Extended Data Fig. 5a).

23 To investigate the cellular features of the intermediate states, we performed a  
24 differential expression analysis, comparing the gene expressions of the tumor cells  
25 with intermediate states to the other tumor cells. We found a set of up-regulated genes  
26 associated with cell growth and development (*ALDOA*, *MDK*, *ENO1*, *TKTL1*) (Fig.  
27 4e, Extended Data Fig. 5b, Supplementary Table 6). Among them, *ALDOA* was

1 proved to serve as a driver for HCC cell growth under hypoxia<sup>38</sup>. *MDK* has been  
2 proposed as a multifunctional protein in HCC development, progression, metastasis,  
3 and recurrence<sup>39</sup>. The up-regulated genes were also enriched in glycolysis, hypoxia,  
4 and UV\_response\_up besides cycling-related pathways (Fig. 4e, Extended Data Fig.  
5 5c).

6 Subsequently, we explored the association between the intermediate states and  
7 prognosis in HCC using another database HCCDB which integrates multiple large-  
8 scale clinical cohorts to examine the gene expression variations in HCC<sup>40</sup>. Using  
9 Fisher's exact test, we compared the poor prognosis-associated genes listed in  
10 HCCDB with the genes upregulated in the intermediate state tumor cells. Remarkably,  
11 the results demonstrated a significant overlap between these two gene lists, indicating  
12 a potential association between the intermediate states and the poor prognosis in HCC  
13 (Fig. 4f). We further examined the clinical relevance of the differentially expressed  
14 genes in intermediate states in the TCGA cohort<sup>41</sup> and found that higher intermediate  
15 gene signature score is significantly associated with worse overall survival (Fig. 4g).  
16 Taken together, the intermediate states analysis revealed a group of tumor cells with  
17 multiple malignant features which located between the hepatocyte and cholangiocyte  
18 states.

19 Finally, we studied the relationship between the deviated scores and the  
20 intermediate scores. It was observed that MACRO+ macrophages (kuppfer cells) had  
21 low intermediate scores and deviated scores, consistent with its preference in non-  
22 tumor tissues. And while, MMP9+ macrophages get high deviated scores and low  
23 intermediate scores, which were also consistent with that MMP9+ macrophages were  
24 in a terminal state in HCC<sup>9</sup> (Extended Data Fig. 5d).



**Fig. 4 | Intermediate state analysis enable the identification of the tumor cells associated poor prognosis.** a, The UMAP visualization of the intermediate scores in the tumor samples. b, The ratio of the intermediate state cells occupying the proportions of the two “terminal” cell types. c, The UMAP visualization of the tumor (epithelial) cells annotated as hepatocyte-like, cholangiocyte-like and intermediate state cells by LiverCT. d, The expressions of selected marker genes for cholangiocytes and hepatocytes, and also the highly expressed genes in the intermediate state tumor cells. e, Three cancer hallmarks enriched in the upregulated genes in the intermediate state tumor cells (adjusted  $P < 0.05$ ). f, Venn diagram of the genes upregulated in the intermediate state tumor cells and the unfavorable genes listed in HCCDB. g, The survival analysis of the TCGA HCC cohort based on the gene signature derived from the tumor cells with intermediate states (P value was calculated using log-rank test).



# 1 Tumor cell zonation tendency mapping defines novel HCC

## 2 subtypes

3 Hepatocytes in different lobular zones display spatial-ordered functional  
4 heterogeneities. However, in HCC tumor tissues, the lobule-like patterns are usually  
5 lost. It is an interesting question that whether the malignantly transformed hepatocytes  
6 (tumor cells) still have zonation tendency and whether the tendency patterns are  
7 associated with clinical outcome.

8 Leveraging the second module of LiverCT, the HCC tumor cells were mapped  
9 into three zones as the P (Periportal/Portal) state, the M (Mid) state or the C (Central)  
10 state (Methods). We observed different distributions of tumor cells' zonation labels in  
11 different patients. Based on the different compositions of the LiverCT annotated  
12 zonation labels of tumor cells, we divided the patients into three groups, namely  
13 HCC\_P (dominated by "Periportal+Portal" state cells), HCC\_M (dominated by "Mid"  
14 state cells), and HCC\_C (dominated by "Central" state cells) (Fig. 6a) (Methods).

15 To explore the distinctive features of the three HCC subtypes, we calculated the  
16 differentially expressed genes (DEGs) of each of the subtypes by using patient-  
17 specific pseudo-bulk data (Fig. 6b, Supplementary Table 7). Results show that in  
18 addition to several well-known zonal marker genes, different subtypes of HCC  
19 patients exhibited unique expression of many other non-zonal genes. Notably, *CD24*  
20 was among the top DEGs of the HCC\_P subtype, suggesting a more malignant  
21 phenotype<sup>42</sup>. We further performed a survival analysis in TCGA<sup>41</sup> bulk data, and  
22 found that higher HCC\_P signature scores were significantly associated with poorer  
23 clinical prognosis (Fig. 6c). Multiple independent cohorts in HCCDB<sup>40,43</sup> showed the  
24 similar observations (Extended Data Fig. 6a).

25 By scoring the curated stemness gene sets (*ANPEP*, *CD24*, *CD44*, *PROM1*,  
26 *EPCAM*)<sup>44</sup> and metastatic gene sets<sup>9</sup> on the pseudo-bulk data, we found statistically  
27 significant differences in the distribution of the stemness scores ( $P = 2e-5$ , ANOVA



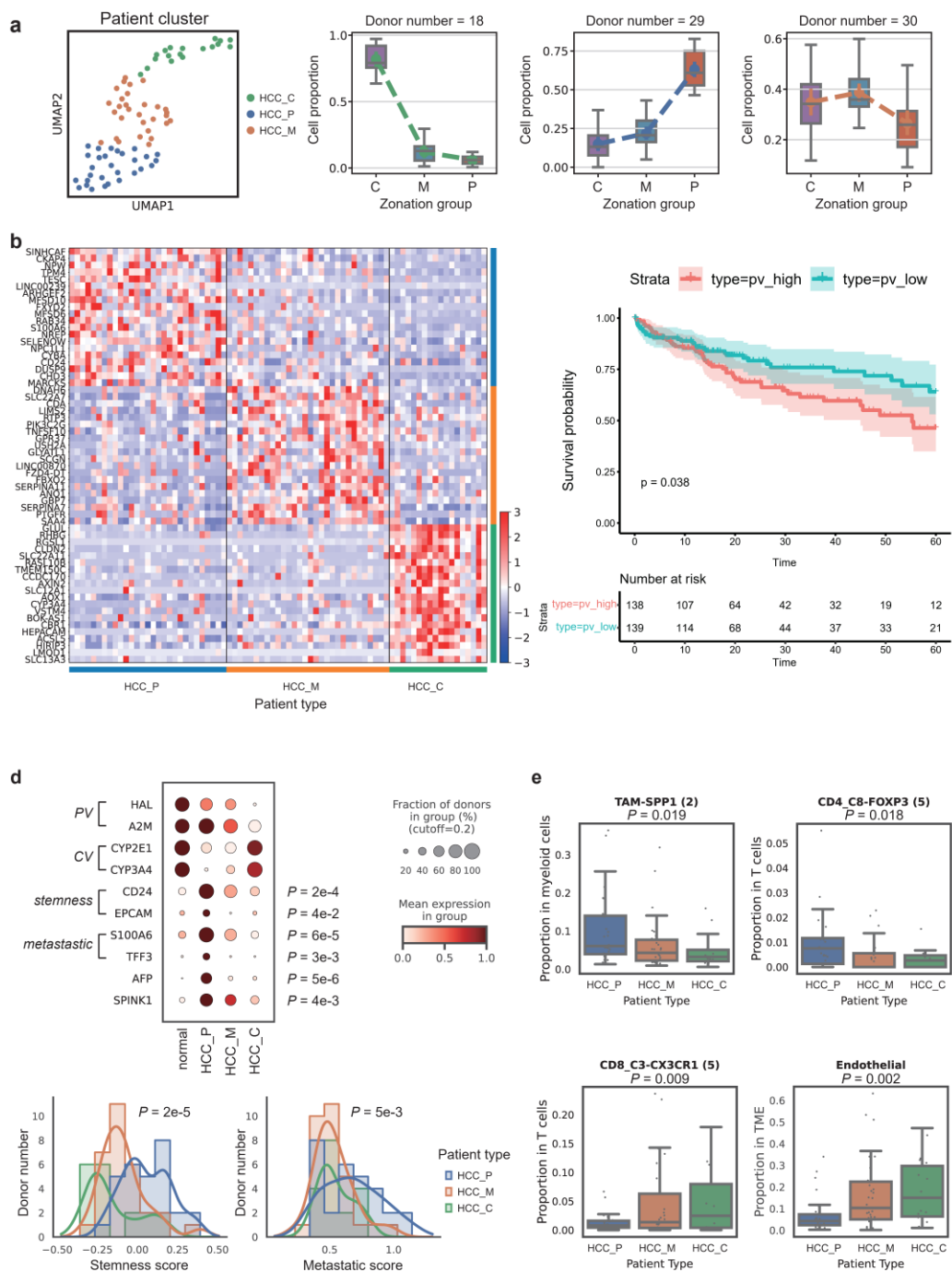
test) and the metastasis scores ( $P = 5e-3$ , ANOVA test) among the three HCC subtypes. The HCC\_P patients showed much higher scores of cancer stemness and metastasis than the other two subtypes (Fig. 6c). Besides, in tumor cells of HCC\_P, we observed the higher expression of *AFP*, which is a marker for fetal liver and a well-known marker for primary liver cancer. Additionally, *SPINK1*, a marker for hepatoblasts<sup>45</sup> and a tumor-promoting factor<sup>46</sup>, exhibited up-regulation in tumor cells of HCC\_P (Fig. 6c). These differences were not observed among the cells in different zones in normal liver, suggesting that these alterations should be associated with pathological changes (Extended Data Fig. 6b, c). Based on these findings, we speculated that the HCC\_P subtype may correspond to a more poorly differentiated HCC phenotype.

Then, we investigated the dys-regulated metabolic processes in these HCC subtypes. Pseudo-bulk of HCC patients and normal donors were scored using a set of curated metabolic gene modules<sup>47</sup>. We found that urea cycle, which is normally conducted in portal and periportal regions<sup>48</sup>, is most significantly down-regulated in HCC\_P patients. The finding suggested that tumor hepatocytes of HCC\_P subtype presented a higher degree of urea cycle disorder (UCD), which correlated with poor prognosis<sup>47</sup>. Conversely, glycolysis was not a normal functional process of portal and periportal hepatocytes<sup>23,49</sup>, but we observed a slight upregulation of the glycolytic modules in tumor cells of HCC\_P patients (Extended Data Fig. 6d).

We further conducted a detailed analysis of the characteristics of the tumor microenvironment (TME) in patients with HCC\_P subtype. Utilizing scCancer2<sup>50</sup>, we transferred TME cell labels from multiple public datasets<sup>15,51,52</sup> to our collected data. Our analysis revealed that HCC\_P patients had a higher abundance of SPP1+ tumor-associated macrophages (TAMs) and FOXP3+ CD4 T cells compared to the other subtypes. In contrast, the proportions of CX3CR1+ CD8 T cells and endothelial cells were lower in HCC\_P patients (Fig. 6e, Extended Data Fig. 6e). We noticed that SPP1+ TAMs were described as a potentially pro-tumorigenic/pro-metastatic subtype

1 in colorectal cancer. Intriguingly, we found that cells annotated as SPP1+ TAMs in  
 2 HCC TME displayed high expression levels of markers for MMP9+ macrophages<sup>9</sup>  
 3 (Extended Data Fig. 6f), which had been previously implicated in promoting HCC  
 4 progression. These observations suggested that patients classified as HCC\_P subtype  
 5 might exhibit a relatively immunosuppressive TME.

6 Taken together, the zonation mapping of tumor cells by LiverCT defines novel  
 7 HCC subtypes. Notably, HCC\_P patients exhibited the worst overall survivals,  
 8 characterized by increased expression of stemness and metastatic factors, along with  
 9 the presence of metabolism dysregulation and an immunosuppressive  
 10 microenvironment.



1  
2 **Fig. 5 | The discovery of HCC subtypes by tumor cell zonation mapping.** a, The UMAP of the  
3 three HCC subtypes defined based on the distributions of LiverCT mapped tumor cells' zonation  
4 states. b, The heatmap the differentially expressed genes among the three subtypes using the  
5 patient-level pseudo-bulk data. c, The survival analysis of the TCGA HCC patients based on the  
6 HCC\_P pseudo-bulk signature (the P value was calculated by log-rank test). d, The expression of  
7 marker genes for central, peri-portal areas and HCC\_P patients (top). The stemness score and

1 metastatic score distribution among different patient types (bottom). e, The proportions of the  
2 TME cell subtypes in the three HCC subtypes.

## 3 **A web-based portal of multidimensional portraits of the** 4 **atlas**

5 To facilitate a convenient browsing, we design a web-based portal for the atlas. It  
6 contains two user-friendly tools, namely *cell mapping* and *cell sorting*, and four  
7 portraits which are the *gene portrait*, *cell portrait*, *zonation portrait* and *disease*  
8 *portrait*, each portrait contains multiple different views (Fig. 6).

9 Detailed information of genes, cells and zonation of normal reference cell map  
10 can be found in the corresponding portraits. The *gene portrait* provides the expression  
11 distribution of the selected gene across cell types. The *cell portrait* portrays the uHAF  
12 tree of uniLIVER as well as the features of the selected cell type quantitatively,  
13 including its number and cell-cell interaction. The *zonation portrait* provides the  
14 expression distribution of the selected gene across zonation and highly expressed  
15 genes (HEGs) of each zone.

16 For mapping new datasets onto the normal reference map, we developed a *cell*  
17 *mapping* tool LiverCT (freely available via uniLIVER website). The method mainly  
18 contains three parts, namely cell type classification, “variant” state identification,  
19 hepatocyte zonation reconstruction. When users input single-cell sequencing data,  
20 LiverCT will provide predicted cell types at Level 1 and Level 2. Additionally, it will  
21 provide two scores: the deviated score and the intermediate score, as well as cells the  
22 recommended thresholds for identifying the cells with deviated states or intermediate  
23 states, respectively.

24 Using LiverCT, we comprehensively annotated the disease data and constructed  
25 the *disease portrait*. The *disease portrait* shows the characteristics of the disease from  
26 two views: (1) *molecular view*; (2) *cellular view*. The *molecular view* presents the



1 **Fig. 6 | The database content and online tools of uniLIVER.** The database consists of four  
2 portraits and two tools: (1) the *gene portrait* showing the gene expression among different cell  
3 types and subtypes within a lineage; (2) the *gell portrait* displaying uHAF tree and the cell-cell  
4 interaction within normal data; (3) the *zonation portrait* showing the gene expression within four  
5 zones in liver lobules; (4) the *cell mapping* page displaying the cell annotation and variant state  
6 identification as well as zonation reconstruction pipeline of LiverCT; (5) the *disease portrait*  
7 showing the characteristics of deviated states and intermediated states; (6) The *cell sorting*  
8 providing a one line tool which allows users to download data in uniLIVER flexibly.  
9

# 1 Discussion

2 In this study, analogy to the genome sequence mapping, we have provided a machine  
3 learning based framework for disease “variant” analysis. As a tool for uniLIVER,  
4 LiverCT get several interesting findings by mapping disease datasets to the normal  
5 reference map. It finds that neutrophils and hepatic stellate cells are strongly deviated  
6 in adjacent tumor, and the intermediate-state tumor cells are associated with  
7 unfavorable outcomes in HCC.

8 The function of hepatocytes along the lobule radial axis is highly heterogeneous,  
9 which in turn results in differences of zonal patterns of drug responses and oncogenic  
10 transformation<sup>48</sup>. Although hepatocytes’ function is impaired by diseases, we posit  
11 that they still exhibit the characteristics of the CV-PV axis at the global transcriptome  
12 level. These characteristics might be influenced by both the microenvironment and  
13 long-term epigenetic phenomena<sup>54,55</sup>. Tumor cell zonation tendency mapping defines  
14 novel HCC subtypes. Among them, the HCC\_P subtype has worst survival with a  
15 SPP1+ macrophage infiltrated suppressive immune microenvironment. Clinically, the  
16 HCC novel subtypes enable different therapy choices. Further investigation is needed  
17 to elucidate the molecular mechanisms by which tumor cells interacts with immune  
18 cells, ultimately resulting in a poorer prognosis in HCC\_P patients.

19 Defining a comprehensive and refined normal reference map is essential but  
20 challenging, as it requires capturing both cellular and population variations<sup>29</sup>. LiverCT  
21 presents a promising opportunity to assess the saturation of the atlas. When  
22 incorporating new healthy datasets, we can evaluate whether any novel deviated states  
23 emerge. If no new cell types are discovered, we can consider the atlas to be ready.  
24 However, if new cell types are detected, we can fine-tune the model until the  
25 identification of previously unobserved cellular states ceases.

26 As data continues to accumulate, there is a recent surge in the development of  
27 large-scale models that have demonstrated state-of-the-art performance across a wide

1 range of downstream tasks<sup>56-59</sup>. These models offer a promising opportunity to create  
 2 a more comprehensive atlas. Also, with the development of spatial transcriptomics  
 3 technology, it is now possible to further portray the spatial microenvironment of a  
 4 cell, which is important to understand the cellular niches of liver.  
 5



# 1 **Methods**

## 2 **Data collection and processing**

3 In the current atlas, we archived 18 human liver datasets, including 6 healthy datasets,  
 4 1 cirrhosis dataset and 11 liver cancer datasets (Supplementary Table 1). For 17  
 5 publicly available datasets provided by dataset generators, we collected the expression  
 6 matrix and processed it using Seurat pipeline<sup>60</sup>. Besides the public datasets, we  
 7 generated ~30K healthy data and use scCancer<sup>61</sup> pipeline to do quality control. The  
 8 gene symbols were unified to the list of 43,878 HUGO Gene Nomenclature  
 9 Committee (HGNC) approved symbols with the toolkit in hECA<sup>53</sup>, with withdrawn  
 10 and alias symbols converted into HGNC approved symbols.

11 In addition, we collected phenotype information at multiple levels including  
 12 donor, sample and cell. At the donor level, we gathered gender, age and fibrotic status  
 13 if available. At the sample level, we categorized the sample status according to its  
 14 location, harmonizing it as normal (N), primary tumor (T), non-tumor (NT), the joint  
 15 area between the tumor and adjacent normal tissues (PJ), hepatic lymph node (HLN),  
 16 metastatic lymph node (MLN), portal vein tumor thrombus (PVTT), Ascites (ASC),  
 17 Blood (BLO) (Supplementary Table 1). At the cell level, we collected the original  
 18 annotations and standardized them to the cell type at level 1 in the uHAF tree  
 19 (Supplementary Table 2).

## 20 **Normal data integration and annotation**

21 To visualize the cells in the normal reference map, the neighbor graph was built based  
 22 on the 30 latent dimensions that were obtained from the scANVI output with the  
 23 default parameter setting of sc.pp.neighbors function. The dimensionality of cells was  
 24 further reduced using Uniform Manifold Approximation and Projection (UMAP) with

1 sc.tl.umap function based on the neighbor graph built above. To determine the Level 1  
2 label of cells, we used two methods. If the original study provided labels of cells, we  
3 would map those labels to the uHAF to obtain the Level 1 label. If not, marker genes  
4 would be used to identify the cell types.

5 To further annotate cells of each resulting Level 1 cluster, a new neighbor graph  
6 was built using 30 latent dimensions of scANVI. Clusters were classified into level 2  
7 labels using marker genes.

## 8 **Normal hepatocyte zonation annotation**

9 We annotated zonation labels for hepatocytes from the uniLIVER normal reference  
10 map. The zonation groups provided by Guilliams et al. (2022) for the human liver  
11 spatial transcriptome were used as the reference<sup>5</sup>.

12 For each Visium sample, we conducted Wilcoxon test to find differential  
13 expressed genes between C-spots (“Central”) and P-spots (“Periportal” + “Portal”).  
14 This step is implemented via rank\_genes\_groups() function in Scanpy<sup>62</sup>. In order to  
15 mitigate the impact of inter-individual variability, only genes showing significant  
16 zonal differences (pvals\_adj < 0.01 for C-markers, and pvals\_adj < 0.05 for P-  
17 markers) in more than 3 samples were considered. To accommodate scRNA-seq data  
18 characteristics, we filtered out genes with a mean log- normalized expression lower  
19 than 0.1 in hepatocytes from our single-cell data.

20 A min-max scaler is applied to each gene in the same sample first to preserving  
21 gradient information. Then, a spot’s score can be calculated as:

$$22 \quad score = \frac{\text{mean}_{gene \in \{P\text{-markers}\}} Expr_{gene}}{\text{mean}_{gene \in \{P\text{-markers}\}} Expr_{gene} + \text{mean}_{gene \in \{C\text{-markers}\}} Expr_{gene}}$$

23 We visualized the original group labels and our defined score on Visium spots  
24 (Extended Data Fig. 2a). Furthermore, we plotted the score distribution for the four  
25 zonation labels (Extended Data Fig. 2b). These results demonstrated that the score can  
26 effectively indicate the location along the CV-PV axis in a healthy liver. This score

1 can be applied equally to spots in spatial transcriptome, as well as cells in single-cell  
2 transcriptome.

3 For hepatocytes from uniLIVER normal reference map, we first conducted a  
4 quality control step. We removed non-viable cells with percentage of mitochondrial  
5 gene counts over 30%. Also, cells with an expressed gene number lower than 1000  
6 were excluded from the subsequent annotation and analysis. These specific thresholds  
7 were determined based on the distribution of QC indicators obtained using Scanpy  
8 function `calculate_qc_metrics()` (Extended Data Fig. 2c).

9 To address distributional biases between spatial and single-cell transcriptomes, as  
10 well as variations in experimental techniques for single-cell sequencing, we  
11 conducted a correction step. Our hypothesis was that the scores for livers from healthy  
12 donors should exhibit a similar distribution. Therefore, we adjusted the mean and  
13 variance of the score distribution within each batch of single-cell transcriptome data  
14 to align with the corresponding distribution observed in the spatial transcriptome data.

15 We employed a normal function to fit the score distribution of each zonation  
16 label. The parameters were determined via maximum likelihood estimation:  $\mu_l = \bar{X}_l$ ,  
17  $\sigma_l^2 = \frac{n-1}{n} S_l^2$ . We transferred the previously fitted distribution to the single-cell data.  
18 Bayesian estimation was utilized to infer the zonation group of each cell, assuming an  
19 equal prior probability for each zonation label:

$$20 \quad label = \underset{l \in L}{\operatorname{argmax}} p(l)p(score|l) = \underset{l \in L}{\operatorname{argmax}} N(score; \mu_l, \sigma_l^2)$$

21 where L represents the set of the four zonation group labels (Central, Mid, Periportal,  
22 Portal).

## 23 **Modeling the effect of demographic covariates on gene** 24 **programs**

25 To model the effect of demographic covariates (gender and age) on gene programs,  
26 we performed the generalized linear mix model (GLMM). We first split cells by level

2 labels, then filtered out genes that were expressed in fewer than 10 cells. Sample-level pseudo-bulks, which were generated by summing gene counts across cells within each level 2 label for each sample, were used to fit the model. Pseudo-bulks were normalized using calcNormFactors function of edgeR with default parameter settings. Then voom<sup>63</sup> was used to fit GLMM for differential expression and perform hypothesis test on fixed effects. Gene expression was modeled as:

$$\log(\text{normcount}) \sim 1 + \text{age} + \text{gender} + (1|\text{donor ID})$$

where the donor is treated as a random effect, and age and gender are modeled as fixed effects. We used the Benjamini-Hochberg procedure to correct the resulting p-values within each covariate. Significant genes (adjusted p-value < 0.05) were selected for gene set enrichment analysis using the enrichGO function in the clusterProfiler package<sup>64</sup>.

## 13 **LiverCT: a machine learning based cell-type mapping**

14 We developed LiverCT (machine learning based Liver Cell Type mapping), to map  
15 new datasets onto the normal reference map. Two-level cell type labeling was  
16 provided by a hierarchical ensemble learning classifier. On the basis of accurate cell  
17 type prediction, LiverCT identified cells in “variant” states, which can be broadly  
18 categorized into two types: deviated states and intermediate states. Specifically, for  
19 hepatocytes, LiverCT further predicted zonal groups along the CV-PV axis at sub-  
20 lobule scale. The workflow of LiverCT is depicted in Extended Data Fig. 1b.

21 **Batch correction.** To mitigate batch effects between the query data and the  
22 normal reference, query datasets were projected to the common latent space of then  
23 normal reference map using scArches<sup>20</sup>, a transfer learning method. The parameter  
24 “encode\_covariates” of the scANVI model was set to True to allow us to fine-tune the  
25 weights of newly introduced edges in the input layer. We conducted 20 epochs during  
26 the fine-tuning process of scArches. Subsequent models operated in this latent space.

27 **Hierarchical ensemble learning cell type classification.** The manually annotated

1 labels served as the reference standard for classification. The process followed a  
 2 hierarchical tree structure as shown in Fig. 2a, to improve the resolution of cell type  
 3 labeling step by step. The query cells were first divided into 8 major cell types. Then  
 4 within each major type, a finer-grained classification was carried out, resulting in 17  
 5 labels at the second level. Both layers of the classification were implemented using an  
 6 ensemble learning model. It consisted of a Multi-layer Perceptron (MLP) classifier<sup>65</sup>,  
 7 an XGBoost classifier<sup>66</sup>, a Logistic Regression classifier<sup>67</sup> using one-vs-rest strategy  
 8 and a Random Forest classifier<sup>68</sup>. A soft voting strategy was implemented to generate  
 9 the predicted probabilities for each cell type. The algorithm was accelerated using  
 10 parallel threading managed by the joblib (<https://github.com/joblib/joblib>) package.

11 ***Deviated states identification.*** We utilized a One-Class Support Vector Machine  
 12 (OCSVM) for unsupervised novelty detection<sup>69</sup>. For each fine-grained label in the  
 13 second level, a OCSVM model was trained. By delineating the contour of the feature  
 14 space occupied by cells in the normal reference, the OCSVM model effectively  
 15 identified cell states that deviated from the normal states. We first used a Radial Basis  
 16 Function (RBF) Kernel to transform the initial observations to a non-linear feature  
 17 space:

$$18 \quad K(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2\right)$$

19 The feature map for the RBF kernel was approximated with the Nyestrom  
 20 method for acceleration<sup>70,71</sup>, using `sklearn.kernel_approximation.Nystroem()`. Then, a  
 21 linear OCSVM was performed in the transformed feature space. The OCSVM model  
 22 was solved using Stochastic Gradient Descent (SGD). This algorithm was chosen due  
 23 to its efficiency in processing large training sets. The optimization problem was  
 24 defined as follows:

$$25 \quad \min_{w,b} \frac{\nu}{2} \|w\|^2 + b\nu + \frac{1}{n} \sum_{i=1}^n \max(0, 1 - (\langle w, x_i \rangle + b))$$

26 where  $w$  and  $b$  represented the linear coefficient and the intercept to be optimized,  
 27 and  $\nu$  was a hyperparameter. The hyperparameters of the model were automatically

1 set based on the distribution of the training data. Specifically, the parameter  $\gamma$  of the  
 2 RBF kernel was set to  $1/(n\_features * var(X))$  as suggested by the sklearn library.  
 3 The parameter  $\nu$  was incrementally increased until 10% of the training data was  
 4 detected as outliers. Cells located outside the frontier-delimited subspace were  
 5 annotated as under “deviated” states. Euclidean distances between query data  
 6 observations to the frontier hypersphere were calculated. These distances were then  
 7 normalized by the 80th percentile value for each model. Subsequently, the normalized  
 8 distances were negated and truncated between -1 and 1, resulting in the deviated  
 9 scores. Higher deviated scores represented larger deviations from the normal  
 10 distribution.

11 ***Intermediate states identification.*** We assumed that intermediate states only  
 12 existed between cell types under the same major type. A special case is hepatocytes  
 13 and cholangiocytes, where an intermediate state between the two has been  
 14 demonstrated to exist in certain disease conditions<sup>36</sup>. Therefore, even though they are  
 15 already distinguished at the major cell type level, we have still identified the  
 16 intermediate state between them. We employed a one-vs-one SVM model<sup>72</sup> to identify  
 17 the classification boundaries between the top two classes to which the cell was most  
 18 likely to belong. The pipeline consisted of a standard scalar and a linear SVM  
 19 optimized using SGD. We calculated Euclidean distances of the samples to the  
 20 separating hyperplane, then used a generalized RBF kernel to transform distances to  
 21 scores between 0 and 1:

$$22 \quad \text{intermediate score} = \exp(-|distance|)$$

23 The higher intermediate score represented that the cell tended to be more  
 24 intermediate between the two types. A threshold was determined as 0.6 manually to  
 25 classify cells with scores above it as being in intermediate states.

26 ***Tumor hepatocyte zonation states mapping.*** In the context of disease data,  
 27 particularly hepatocellular carcinoma (HCC), the expression of many zonal landmark  
 28 genes was found to be absent or exhibited a loss of gradient. As a result, the scoring

1 approach that relied on a subset of genes defined in a healthy liver was unsuitable for  
2 analyzing disease data. To overcome this, we employed a supervised learning  
3 classifier trained on our normal reference, and used it to transfer zonation labels to the  
4 disease data.

5 To address the individual heterogeneity of human hepatocytes, we used donor ID  
6 as batch labels to train a scANVI model. This model generated a 30-dimensional  
7 latent space with batch corrected. For the input features of the scANVI model, we  
8 identified 2000 highly variable genes (HVGs). The selection of HVGs was performed  
9 using the "seurat\_v3" flavor provided by the Scanpy pipeline.

10 We chose the Random Forest algorithm, which employs feature sampling steps to  
11 ensure reliable classifications even when features are missing. This attribute makes it  
12 particularly suitable for analyzing disease state data. We used the low-dimensional  
13 latent vectors as input for training and implemented the algorithm using  
14 `sklearn.ensemble.RandomForestClassifier` with 100 estimators.

15 We used manually annotated zonation groups as the reference standard and  
16 implemented a more general categorization approach by using three classification  
17 labels for the training process. Specifically, we combined the Periportal and Portal  
18 regions, leading to three labels: C for Central, M for Mid, and P for Periportal+Portal.

19 For disease data, we utilized the transfer learning method, scArches, to acquire  
20 latent space representations consistent with the reference. During the scArches  
21 surgery process, 20 epochs of fine-tuning were performed. We then used the trained  
22 Random Forest classifier to predict zonation label for each individual cell.

23 We built LiverCT on Python (3.9.7), using the following packages: numpy  
24 (1.22.4), scipy (1.8.1), pandas (1.4.3), anndata (0.8.0), scanpy (1.9.1), scArches  
25 (0.5.9), joblib (1.1.0), scikit\_learn (1.1.1), xgboost (1.7.6). The code is open-sourced  
26 at Github (<https://github.com/fyh18/LiverCT>).

# 1 Variant states analysis

2 We sampled up to 1000 cells from each donor to maintain the balance of patient cell  
3 number. After quality control, 272,464 cells were left to constitute the core disease  
4 data.

5 In the T cell function analysis section, we use “sc.tl.score\_genes” in scanpy to  
6 add module score to each cell. Pearson correlation was employed to assess the  
7 correlation between T cell function and deviation scores.

8 The intermediate gene signature was derived by filtering genes based on specific  
9 criteria, including a log fold change (FC) greater than 0.5 and an adjusted p-value  
10 (pvals\_adj) lower than 0.01.

# 11 HCC Classification

12 We analyzed cells annotated as hepatocytes by LiverCT from the core disease data.  
13 Only samples from HCC primary tumors were included. Samples with less than 50  
14 hepatocyte-like cells were filtered out. We then calculated the proportion of cells with  
15 the three zonal labels (“C”, “M” and “P”) for each patient, resulting in a  $n \times 3$   
16 matrix where each row represented a donor and each column represented the cell  
17 proportion of a certain zone. We referred to this matrix as “zonal proportion space of  
18 patients”. Subsequently, we performed a 3-cluster spectral clustering within this  
19 space to classify three HCC subtypes, namely HCC\_C, HCC\_M and HCC\_P.

20 We summed up the counts of all HCC tumor hepatocytes for each patient, and  
21 then perform log-normalization to acquire pseudo bulk data.

# 22 Portraits of uniLIVER

23 **Gene portrait.** We provided the expression distribution of the selected gene across cell  
24 types. The ridge plots showed the non-zero expression distributions in different cell



1 types. The number at the right of the ridge plots showed the non-zero percentages of  
2 the expression values.

3 **Cell portrait.** CellChat<sup>73</sup> was employed to infer cell-cell interactions by analyzing the  
4 expression patterns of known ligand-receptor pairs across diverse cell types. We  
5 followed the official workflow with default parameters.

6 **Zonation portrait.** We provided the average expression values of the selected gene  
7 across zonation. Besides, differentially expressed genes of different zones can be  
8 visualized by heatmap.

9 **Disease portrait.** In *Molecular view*, we present the features of disease states and  
10 intermediate states, as well as the characteristics of disease cell types. The former two  
11 are compared within a specific disease, while the latter is compared between different  
12 disease conditions. In deviated state section, deviated score distribution in level 2 is  
13 displayed and we can see the most susceptible cell type. By selecting a cell type,  
14 differentially expressed genes in deviated states compared with normal states are shown.  
15 Similarly, in *intermediate state* section, the ratio of intermediate states between two cell  
16 type is shown and we can see the differentially expressed genes in intermediate states  
17 compared with the other two cell types. The *Disease cell type* section displays  
18 differentially expressed genes (DEGs) between the selected disease and another  
19 condition in the same cell type. Enrichment analysis is conducted based on the DEGs.

20

## 21 **Data availability**

22 The uniLIVER website is publicly accessible via [<https://liver.unifiedcellatlas.org>].

23 The normal reference map and core disease data (processed as data matrix) are  
24 publicly available through the *databrowser* section and can be easily downloaded  
25 from *download* section in the web server. The source codes, trained models and  
26 documents of LiverCT are also provided at the website.

## 27 **Acknowledgements**

1 This publication is part of the Human Cell Atlas –  
 2 <https://www.humancellatlas.org/publications/>. We thank Qiuyu Lian, Qinglin Mei,  
 3 Yiran Shan, Xinqi Li, Qifan Hu and Nan Yan and Yifan Sun for their help on this  
 4 work. This work is funded by the National Key Research and Development Program  
 5 of China (No. 2021YFF1200901) and the National Natural Science Foundation of  
 6 China (Nos. 61721003, 62133006 and 92268104).

# **Author contributions**

8 J.G. conceived the study. Y.W., Y. F. and Y.M. collected datasets involved in the study.  
 9 Y.W. and Y.F. designed the LiverCT algorithm. Y.F. implemented the LiverCT  
 10 algorithm. Y.M. performed the unsupervised annotation experiment. Y.W. and Y.F.  
 11 designed the biological applications. Y.M. imported data into the database. Y.F., Y.W.,  
 12 Y.L. and Y.M. designed the web page of uniLIVER. M.Y., R.Y., A.C. deployed the  
 13 database. G.D., J.D. and Y.C. collected the clinical samples. Z.C., Y.C., W.L., W.G.,  
 14 J.D., X.Z. and Y.W. provided advice on experiments. Y.W., Y. F. and Y.M. wrote the  
 15 manuscript. J.G. supervised the computational analysis. All authors agreed on the  
 16 final version of the manuscript.

# **Conflicts of interests**

18 None declared.

19

20

# References

1. Hashimshony, T., Wagner, F., Sher, N. & Yanai, I. CEL-Seq: Single-Cell RNA-Seq by Multiplexed Linear Amplification. *Cell Reports* **2**, 666-673 (2012).
2. Tang, F. *et al.* mRNA-Seq whole-transcriptome analysis of a single cell. *Nature Methods* **6**, 377-382 (2009).
3. Aizarani, N. *et al.* A human liver cell atlas reveals heterogeneity and epithelial progenitors. *Nature* **572**, 199-204 (2019).
4. Ramachandran, P. *et al.* Resolving the fibrotic niche of human liver cirrhosis at single-cell level. *Nature* **575**, 512-518 (2019).
5. Williams, M. *et al.* Spatial proteogenomics reveals distinct and evolutionarily conserved hepatic macrophage niches. *Cell* **185**, 379-396. e338 (2022).
6. Payen, V. L. *et al.* Single-cell RNA sequencing of human liver reveals hepatic stellate cell heterogeneity. *JHEP Reports* **3**, 100278 (2021).
7. Andrews, T. S. *et al.* Single-cell, single-nucleus, and spatial RNA sequencing of the human liver identifies cholangiocyte and mesenchymal heterogeneity. *Hepatology Communications* **6**, 821-840 (2022).
8. Losic, B. *et al.* Intratumoral heterogeneity and clonal evolution in liver cancer. *Nature Communications* **11**, 291 (2020).
9. Lu, Y. *et al.* A single-cell atlas of the multicellular ecosystem of primary and metastatic hepatocellular carcinoma. *Nature Communications* **13**, 4594 (2022).
10. Ma, L. *et al.* Single-cell atlas of tumor cell evolution in response to therapy in hepatocellular carcinoma and intrahepatic cholangiocarcinoma. *Journal of hepatology* **75**, 1397-1408 (2021).
11. Massalha, H. *et al.* A single cell atlas of the human liver tumor microenvironment. *Zenodo* (2030).
12. Sun, Y. *et al.* Single-cell landscape of the ecosystem in early-relapse hepatocellular carcinoma. *Cell* **184**, 404-421. e416 (2021).
13. Qi, Z. *et al.* Integrated multiomic analysis reveals comprehensive tumour heterogeneity and novel immunophenotypic classification in hepatocellular carcinomas. *Gut* **68**, 2019 (2019).
14. Zhang, M. *et al.* Single-cell transcriptomic architecture and intercellular crosstalk of human intrahepatic cholangiocarcinoma. *Journal of Hepatology* **73**, 1118-1130 (2020).
15. Zheng, C. *et al.* Landscape of Infiltrating T Cells in Liver Cancer Revealed by Single-Cell Sequencing. *Cell* **169**, 1342-1356.e1316 (2017).
16. Xue, R. *et al.* Liver tumour immune microenvironment subtypes and neutrophil heterogeneity. *Nature*, 1-7 (2022).
17. Ma, L. *et al.* Multiregional single-cell dissection of tumor and immune cells reveals stable lock-and-key features in liver cancer. *Nature Communications* **13**, 7533 (2022).
18. Liu, Y. *et al.* Identification of a tumour immune barrier in the HCC microenvironment that determines the efficacy of immunotherapy. *Journal of Hepatology* **78**, 770-782 (2023).
19. Chen, S. *et al.* hECA: The cell-centric assembly of a cell atlas. *iScience* **25**, 104318 (2022).
20. Lotfollahi, M. *et al.* Mapping single-cell data to reference atlases by transfer learning.

- 1 *Nature biotechnology* **40**, 121-130 (2022).
- 2 21. Gayoso, A. *et al.* A Python library for probabilistic analysis of single-cell omics data.
- 3 *Nature Biotechnology* **40**, 163-166 (2022).
- 4 22. Luecken, M. D. *et al.* Benchmarking atlas-level data integration in single-cell genomics.
- 5 *Nature Methods* **19**, 41-50 (2022).
- 6 23. Yang, Q., Zhang, S., Ma, J., Liu, S. & Chen, S. In Search of Zonation Markers to Identify
- 7 Liver Functional Disorders. *Oxidative Medicine and Cellular Longevity* **2020**, 9374896
- 8 (2020).
- 9 24. Paris, J. & Henderson, N. C. Liver zonation, revisited. *Hepatology* **76** (2022).
- 10 25. Halpern, K. B. *et al.* Single-cell spatial reconstruction reveals global division of labour in
- 11 the mammalian liver. *Nature* **542**, 352-356 (2017).
- 12 26. Sasse, D., Katz, N. & Jungermann, K. FUNCTIONAL HETEROGENEITY OF RAT-
- 13 LIVER PARENCHYMA AND OF ISOLATED HEPATOCYTES. *FEBS LETTERS* **57**, 83-
- 14 88 (1975).
- 15 27. Squair, J. W. *et al.* Confronting false discoveries in single-cell differential expression.
- 16 *Nature Communications* **12**, 5692 (2021).
- 17 28. Crowell, H. L. *et al.* muscat detects subpopulation-specific state transitions from multi-
- 18 sample multi-condition single-cell transcriptomics data. *Nature Communications* **11**, 6077
- 19 (2020).
- 20 29. Sikkema, L. *et al.* An integrated cell atlas of the lung in health and disease. *Nature Medicine*
- 21 **29**, 1563-1577 (2023).
- 22 30. Zhang, Y. & Zhang, Z. The history and advances in cancer immunotherapy: understanding
- 23 the characteristics of tumor-infiltrating immune cells and their therapeutic implications.
- 24 *Cellular & Molecular Immunology* **17**, 807-821 (2020).
- 25 31. Chu, Y. *et al.* Pan-cancer T cell atlas links a cellular stress response state to immunotherapy
- 26 resistance. *Nature Medicine* **29**, 1550-1562 (2023).
- 27 32. Aran, D. *et al.* Comprehensive analysis of normal adjacent to tumor transcriptomes. *Nature*
- 28 *Communications* **8**, 1077 (2017).
- 29 33. Kim, J. *et al.* Transcriptomes of the tumor-adjacent normal tissues are more informative
- 30 than tumors in predicting recurrence in colorectal cancer patients. *Journal of Translational*
- 31 *Medicine* **21**, 209 (2023).
- 32 34. Li, H. *et al.* Dysfunctional CD8 T Cells Form a Proliferative, Dynamically Regulated
- 33 Compartment within Human Melanoma. *Cell* **176**, 775-789.e718 (2019).
- 34 35. Yuan, X. *et al.* Single-cell profiling of peripheral neuroblastic tumors identifies an
- 35 aggressive transitional state that bridges an adrenergic-mesenchymal trajectory. *Cell*
- 36 *Reports* **41**, 111455 (2022).
- 37 36. Sia, D., Villanueva, A., Friedman, S. L. & Llovet, J. M. Liver Cancer Cell of Origin,
- 38 Molecular Class, and Effects on Patient Prognosis. *Gastroenterology* **152**, 745-761 (2017).
- 39 37. Fan, B. *et al.* Cholangiocarcinomas can originate from hepatocytes in mice. *The Journal of*
- 40 *Clinical Investigation* **122**, 2911-2915 (2012).
- 41 38. Niu, Y. *et al.* Loss-of-Function Genetic Screening Identifies Aldolase A as an Essential
- 42 Driver for Liver Cancer Cell Growth Under Hypoxia. *Hepatology* **74** (2021).

39. Gowhari Shabgah, A. *et al.* Shedding more light on the role of Midkine in hepatocellular carcinoma: New perspectives on diagnosis and therapy. *Iubmb Life* **73**, 659-669 (2021).
40. Lian, Q. *et al.* HCCDB: a database of hepatocellular carcinoma expression atlas. *Genomics, proteomics & bioinformatics* **16**, 269-275 (2018).
41. Chang, K. *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nature Genetics* **45**, 1113-1120 (2013).
42. Barkal, A. A. *et al.* CD24 signalling through macrophage Siglec-10 is a target for cancer immunotherapy. *Nature* **572**, 392-396 (2019).
43. Ziming, J. *et al.* HCCDB v2.0: Decompose the Expression Variations by Single-cell RNA-seq and Spatial Transcriptomics in HCC. *bioRxiv*, 2023.2006.2015.545045 (2023).
44. Tsui, Y.-M., Chan, L.-K. & Ng, I. O.-L. Cancer stemness in hepatocellular carcinoma: mechanisms and translational potential. *British Journal of Cancer* **122**, 1428-1440 (2020).
45. Wesley, B. T. *et al.* Single-cell atlas of human liver development reveals pathways directing hepatic cell fates. *Nature Cell Biology* **24**, 1487-1498 (2022).
46. Chen, F. *et al.* Targeting SPINK1 in the damaged tumour microenvironment alleviates therapeutic resistance. *Nature Communications* **9**, 4315 (2018).
47. Wu, T. *et al.* Discovery of a Carbamoyl Phosphate Synthetase 1-Deficient HCC Subtype With Therapeutic Potential Through Integrative Genomic and Experimental Analysis. *Hepatology* **74** (2021).
48. Ben-Moshe, S. & Itzkovitz, S. Spatial heterogeneity in the mammalian liver. *Nature Reviews Gastroenterology & Hepatology* **16**, 395-410 (2019).
49. Manco, R. & Itzkovitz, S. Liver zonation. *JOURNAL OF HEPATOLOGY* **74**, 466-468 (2021).
50. Zeyu, C. *et al.* scCancer2: data-driven in-depth annotations of the tumor microenvironment at single-level resolution. *bioRxiv*, 2023.2008.2022.554137 (2023).
51. Guo, X. *et al.* Global characterization of T cells in non-small-cell lung cancer by single-cell sequencing. *Nature Medicine* **24**, 978-985 (2018).
52. Zhang, L. *et al.* Single-Cell Analyses Inform Mechanisms of Myeloid-Targeted Therapies in Colon Cancer. *Cell* **181**, 442-459.e429 (2020).
53. Chen, S. *et al.* hECA: The cell-centric assembly of a cell atlas. *Iscience* **25**, 104318 (2022).
54. Brosch, M. *et al.* Epigenomic map of human liver reveals principles of zoned morphogenic and metabolic control. *Nature Communications* **9**, 4150 (2018).
55. Ben-Moshe, S. *et al.* Spatial sorting enables comprehensive characterization of liver zonation. *Nature Metabolism* **1**, 899-911 (2019).
56. Theodoris, C. V. *et al.* Transfer learning enables predictions in network biology. *Nature* **618**, 616-624 (2023).
57. Haotian, C., Chloe, W., Hassaan, M. & Bo, W. scGPT: Towards Building a Foundation Model for Single-Cell Multi-omics Using Generative AI. *bioRxiv*, 2023.2004.2030.538439 (2023).
58. Yang, F. *et al.* scBERT as a large-scale pretrained deep language model for cell type annotation of single-cell RNA-seq data. *Nature Machine Intelligence* **4**, 852-866 (2022).
59. Minsheng, H. *et al.* Large Scale Foundation Model on Single-cell Transcriptomics. *bioRxiv*,

1        2023.2005.2029.542705 (2023).

2        60. Stuart, T. *et al.* Comprehensive integration of single-cell data. *Cell* **177**, 1888-1902. e1821

3        (2019).

4        61. Guo, W. *et al.* scCancer: a package for automated processing of single-cell RNA-seq data

5        in cancer. *Briefings in Bioinformatics* **22** (2020).

6        62. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data

7        analysis. *Genome Biology* **19**, 15 (2018).

8        63. Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: precision weights unlock linear model

9        analysis tools for RNA-seq read counts. *Genome Biology* **15**, R29 (2014).

10       64. Wu, T. *et al.* clusterProfiler 4.0: A universal enrichment tool for interpreting omics data.

11       *The Innovation* **2**, 100141 (2021).

12       65. Popescu, M.-C., Balas, V., Perescu-Popescu, L. & Mastorakis, N. Multilayer perceptron

13       and neural networks. *WSEAS Transactions on Circuits and Systems* **8** (2009).

14       66. Chen, T. & Guestrin, C. in *Proceedings of the 22nd ACM SIGKDD International*

15       *Conference on Knowledge Discovery and Data Mining* 785–794 (Association for

16       Computing Machinery, San Francisco, California, USA, 2016).

17       67. Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized Linear

18       Models via Coordinate Descent. *JOURNAL OF STATISTICAL SOFTWARE* **33**, 1-22 (2010).

19       68. Breiman, L. Random Forests. *Machine Learning* **45**, 5-32 (2001).

20       69. Schölkopf, B., Platt, J. C., Shawe-Taylor, J., Smola, A. J. & Williamson, R. C. Estimating

21       the Support of a High-Dimensional Distribution. *Neural Computation* **13**, 1443-1471

22       (2001).

23       70. Williams, C. K. I. & Seeger, M. in *Proceedings of the 13th International Conference on*

24       *Neural Information Processing Systems* 661–667 (MIT Press, Denver, CO, 2000).

25       71. Yang, T., Li, Y.-F., Mahdavi, M., Jin, R. & Zhou, Z.-H. in *NIPS*.

26       72. Cortes, C. & Vapnik, V. Support-vector networks. *Machine Learning* **20**, 273-297 (1995).

27       73. Jin, S. *et al.* Inference and analysis of cell-cell communication using CellChat. *Nature*

28       *Communications* **12**, 1088 (2021).

29

# 1 **Additional information**

2 Extended Data Fig 1-6

3 Supplementary Table 1-7

1 Table1 | The collected datasets in uniLIVER

Study	No. of donors	Gender (M:F)	Age (yr)	Sequencing method	Final #cells	Source
Aizarani et al. 2019 <sup>3</sup>	9	n/a	n/a	mCel-Seq2	9,466	GSE124395
Ramachandran et al. 2019/healthy <sup>4</sup>	5	4M:1F	n/a	10X	34,601	GSE136103
Guilliams et al. 2022 <sup>5</sup>	34	19M:15F	28-77	10X	167,510	GSE192742
Payen et al. 2021 <sup>6</sup>	2	1M:1F	3-39	10X	26,685	GSE158723
Andrews et al. 2021 <sup>7</sup>	5	2M:3F	18-60	10X	59,977	GSE185477
Gu et al. 2022	4	2M:2F	47-66	10X	32,886	In house
Losic et al. 2020 <sup>8</sup>	2	1M:1F	66-67	10X	49,674	GSE112271
Lu et al. 2022 <sup>9</sup>	10	9M:1F	48-65	10X	71,915	GSE149614
Ma et al. 2021 <sup>10</sup>	37	23M:13F	35-81	10X	48,318	GSE116113
Massalha et al. 2020 <sup>11</sup>	6	2M:4F	40-74	MARS-seq	4,691	GSE146409
Sun et al. 2021 <sup>12</sup>	18	17M:1F	42-76	MIRALCS	16,498	CNP0000650
Zhang et al. 2019 <sup>13</sup>	10	9M:1F	32-84	10X/SMART-seq2	73,261	GSE140228
Zhang et al. 2020 <sup>14</sup>	6	4M;2F	n/a	10X	37,814	GSE138709/ GSE142784
Zheng et al. 2017 <sup>15</sup>	6	4M;2F	26-64	SMART-seq2	5,063	GSE98638
Ramachandran et al. 2019/cirrhosis <sup>4</sup>	5	3M:2F	n/a	10X	26,279	GSE136103
Xue et al. 2022 <sup>16</sup>	124	94M:30F	31-88	10X	1,337,829	PRJCA007744
Ma et al. 2022 <sup>17</sup>	7	n/a	n/a	10X	112,506	GSE189903
Liu et al. 2023 <sup>18</sup>	6	6F	48-64	10X	83,793	skrx2fz79n

2