

RESEARCH

Open Access



# Application of a single-cell-RNA-based biological-inspired graph neural network in diagnosis of primary liver tumors

Dao-Han Zhang<sup>1†</sup>, Chen Liang<sup>1†</sup>, Shu-Yang Hu<sup>1†</sup>, Xiao-Yong Huang<sup>1,2,3†</sup>, Lei Yu<sup>1,2,3†</sup>, Xian-Long Meng<sup>1,2,3</sup>, Xiao-Jun Guo<sup>1,2,3</sup>, Hai-Ying Zeng<sup>4</sup>, Zhen Chen<sup>5</sup>, Lv Zhang<sup>5</sup>, Yan-Zi Pei<sup>1</sup>, Mu Ye<sup>1</sup>, Jia-Bin Cai<sup>1</sup>, Pei-Xin Huang<sup>2</sup>, Ying-Hong Shi<sup>1,2,3</sup>, Ai-Wu Ke<sup>2,3</sup>, Yi Chen<sup>2</sup>, Yuan Ji<sup>4</sup>, Yujiang Geno Shi<sup>1,2</sup>, Jian Zhou<sup>1,2,3</sup>, Jia Fan<sup>1,2,3,6</sup>, Guo-Huan Yang<sup>1\*</sup>, Qi-Man Sun<sup>1\*</sup>, Guo-Ming Shi<sup>1,2,3,5,6\*</sup>  and Jia-Cheng Lu<sup>1,2,3\*</sup>

## Abstract

Single-cell technology depicts integrated tumor profiles including both tumor cells and tumor microenvironments, which theoretically enables more robust diagnosis than traditional diagnostic standards based on only pathology. However, the inherent challenges of single-cell RNA sequencing (scRNA-seq) data, such as high dimensionality, low signal-to-noise ratio (SNR), sparse and non-Euclidean nature, pose significant obstacles for traditional diagnostic approaches. The diagnostic value of single-cell technology has been largely unexplored despite the potential advantages. Here, we present a graph neural network-based framework tailored for molecular diagnosis of primary liver tumors using scRNA-seq data. Our approach capitalizes on the biological plausibility inherent in the intercellular communication networks within tumor samples. By integrating pathway activation features within cell clusters and modeling unidirectional inter-cellular communication, we achieve robust discrimination between malignant tumors (including hepatocellular carcinoma, HCC, and intrahepatic cholangiocarcinoma, iCCA) and benign tumors (focal nodular hyperplasia, FNH) by scRNA data of all tissue cells and immunocytes only. The efficacy to distinguish iCCA from HCC was further validated on public datasets. Through extending the application of high-throughput scRNA-seq data into diagnosis approaches focusing on integrated tumor microenvironment profiles rather than a few tumor markers, this framework also sheds light on minimal-invasive diagnostic methods based on migrating/circulating immunocytes.

<sup>†</sup>Dao-Han Zhang, Chen Liang, Shu-Yang Hu, Xiao-Yong Huang and Lei Yu contributed equally to this work.

\*Correspondence:

Guo-Huan Yang  
yang.guohuan@zs-hospital.sh.cn  
Qi-Man Sun  
sun.qiman@zs-hospital.sh.cn  
Guo-Ming Shi  
shi.guoming@zs-hospital.sh.cn  
Jia-Cheng Lu  
jclu@fudan.edu.cn

Full list of author information is available at the end of the article



© The Author(s) 2024. **Open Access** This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

**Keywords** Single-cell transcriptome, Graph neural network, Diagnostic model, Primary liver tumors, Tumor microenvironment

## Introduction

Primary liver tumors can be classified into two major types: primary liver cancer and benign lesion. Primary liver cancer is the sixth most commonly diagnosed cancer and the third leading cause of cancer death worldwide in 2020, with approximately 906,000 new cases and 830,000 deaths in 2020 [1, 2]. Hepatocellular carcinoma (HCC) and intrahepatic cholangiocarcinoma (iCCA) are two major pathological subtypes of primary liver cancer. Liver benign lesions, like angiomyolipoma (AML), focal nodular hyperplasia (FNH) and hepatocellular adenomas, are generally asymptomatic and do not require specific treatment but rather image-based monitoring [3, 4]. While for HCC and iCCA, timely treatment is critical to improve the survival. However, accurate and rapid diagnosis is challenging in clinical practice [5].

Serum biomarkers such as alpha fetoprotein (AFP) in combination with classical radiology were commonly used for diagnosis of HCC [6–8]. Elevated AFP (cut-off of >20 ng/mL) was detected in only 39% of HCC [9], making it controversial as a surveillance biomarker for HCC [10]. Additionally, identifying the property of the lesions can sometimes be challenging [11], although magnetic resonance imaging (MRI) and contrasted computed tomography (CT) are routinely used for diagnosis of liver tumors [2, 4, 8, 12, 13]. Furthermore, radiology tests largely rely on the interpretation of well-trained radiologists [5, 14].

Pathology is a gold standard for tumor diagnosis [4, 13, 15, 16]. Biopsies of liver lesions can clarify their properties and molecular typing, which can provide valuable guidance for the treatment and prognosis prediction. However, adequate and high-quality tissue specimens, strict procedure of staining slides, and well-trained pathologists are needed to make correct diagnosis [17, 18]. Moreover, the routine histological features on haematoxylin-eosin (H&E) staining are often insufficient to make definitive diagnosis, multiple immunohistochemical staining often need to be performed. Therefore, exploring novel methods to diagnose liver tumors is important.

Recently, single-cell RNA sequencing (scRNA-seq) have depicted the integrated profiling of individual cells in tumor microenvironment. Machine learning (ML) and artificial neural networks (ANN) enables single-cell technology to better characterize distinct cell subsets, quantify inter-cellular communication, dissect cell fate branch points [19]. Such technologies help to better understand properties of tumor cells and their interactions with the microenvironment, making it a potential tool for tumor

diagnosis [20]. Currently, scRNA-seq technologies are mainly used in basic research at cell-level, rather than clinical diagnoses at patient-level [21]. Despite researchers have used omics data and simple ANN in diagnosis of carcinoma of unknown primary (CUP) in clinical setting [22–24], ANN models based on scRNA-seq data remain largely unexplored.

Hence, we reported a pilot exploration of tumor diagnosis based on the scRNA-seq of human liver tumors using a biologically-inspired graph neural network (GNN). We first successfully distinguished primary liver cancers including HCC and iCCA from benign lesions such as FNH. We then constructed a pathologist-independent automated workflow characterizing no expertise or manual work (i.e. cell annotation) for raw scRNA-seq data processing. Our framework is also applicable to scRNA-seq data of immunocytes in the absence of tumor cells, validating its ability to diagnose based on the systemic features of tumor microenvironments. The capability to distinguish two subtypes of primary liver cancer (HCC and iCCA) by our framework is also validated in both internal and public datasets.

## Results

### Diagnosis algorithms based on traditional scRNA analysis protocol

We initially adopted the traditional scRNA-seq analysis protocol to differentiate between benign and malignant tumor samples [20]. We aggregated all patient data, performed dimensional reduction and clustering, and manually annotated the epithelial&hepatocytes, which include both benign and malignant tumor cells. We then identified the top 10 differentially expressed genes between epithelial&hepatocytes of benign and malignant tumor samples, and used their averaged expressions in epithelial&hepatocytes annotated from each sample to train a Multilayer Perceptron (MLP). However, this method performed poorly, achieving an Area Under the Curve (AUC) value of only 0.52, an accuracy (ACC) of  $48 \pm 3\%$ , and an F1 score (F1) of  $0.50 \pm 0.02$ . These results indicate that the traditional scRNA-seq analysis approach, which focus on figuring out tumor gene biomarkers expressed in tumor cell clusters through manually annotation, was not sufficient to reliably distinguishing between benign and malignant tumor samples. Therefore, new data-driven methods distilling more stable and global features from original expression data is needed, to create more reliable and objective diagnostic pipeline.

### The gossip flow (GF) framework

Inspired by intercellular signal transduction process within tumor microenvironment such as antigen presenting and leukocyte recruiting [25], we developed the Gossip Flow (GF) model, a GNN-based model for the diagnosis of liver tumors, based on intercellular interactions, pathway enrichment, and cell clustering. The entire diagnostic framework can be automatically performed based on scRNA-seq data. Our model is named Gossip Flow because it was adapted from GNN architectures used to analyse social interactions [26]. GF does not attempt to directly distinguish tumor cells from all the cells collected from the sample but rather pictures and amplifies the difference through inter-cellular interactions.

GF starts by dimensional reduction and clustering. Based on the result of unsupervised clustering, original expression matrix is further processed into pathway enrichment score (PES) matrix and inter-cluster communication features (CNT) matrix through Gene Set Variation Analysis (GSVA) pathway enrichment analysis and CellChat ligand-receptor interaction probability estimation, respectively. A cell activation and interaction graph structure is then generated from these two matrices. PES matrix serves as the node features ( $h_{i(i \neq 0)}^0$ , enrichment scores of  $K$  representative pathways of cell cluster  $i$ ) while the CNT matrix represents edge weights ( $e_{ij}$ , the average directional ligand-receptor interaction probability from cell cluster  $i$  to  $j$ ), ensuring the message propagation direction in the following directed graph convolutional network (DGCN) layers in alliance with the direction of the ligand-receptor message-passing pathways in vivo. The graph structure serves as input for a  $L$  layer DGCN with a master node (node 0,  $h_0^0$  padded with zeros) added to receive projections from all other nodes ( $e_{i0} = 1$ ,  $e_{0i(i \neq 0)} = 0$ ) and integrating features for readout [27]. The feature vector of the master node of the final layer ( $h_0^L$ ) then projects to the output node via a fully connected (FC) layer. All the DGCN layers share the same set of trainable parameters ( $W_1$  represents intra-node message propagation, while  $W_2$  represents inter-node message propagation, in every layer) to avoid overfitting. Consisting of extremely few trainable parameters compared with normal ANN models, GF can be reasonably applied in prediction tasks with relatively small samples.

A graphical workflow of the GNN diagnosis framework based on liver tissue scRNA-seq is shown in Fig. 1.

### GF robustly classified primary liver cancer (iCCA&HCC) from benign tumors (FNH)

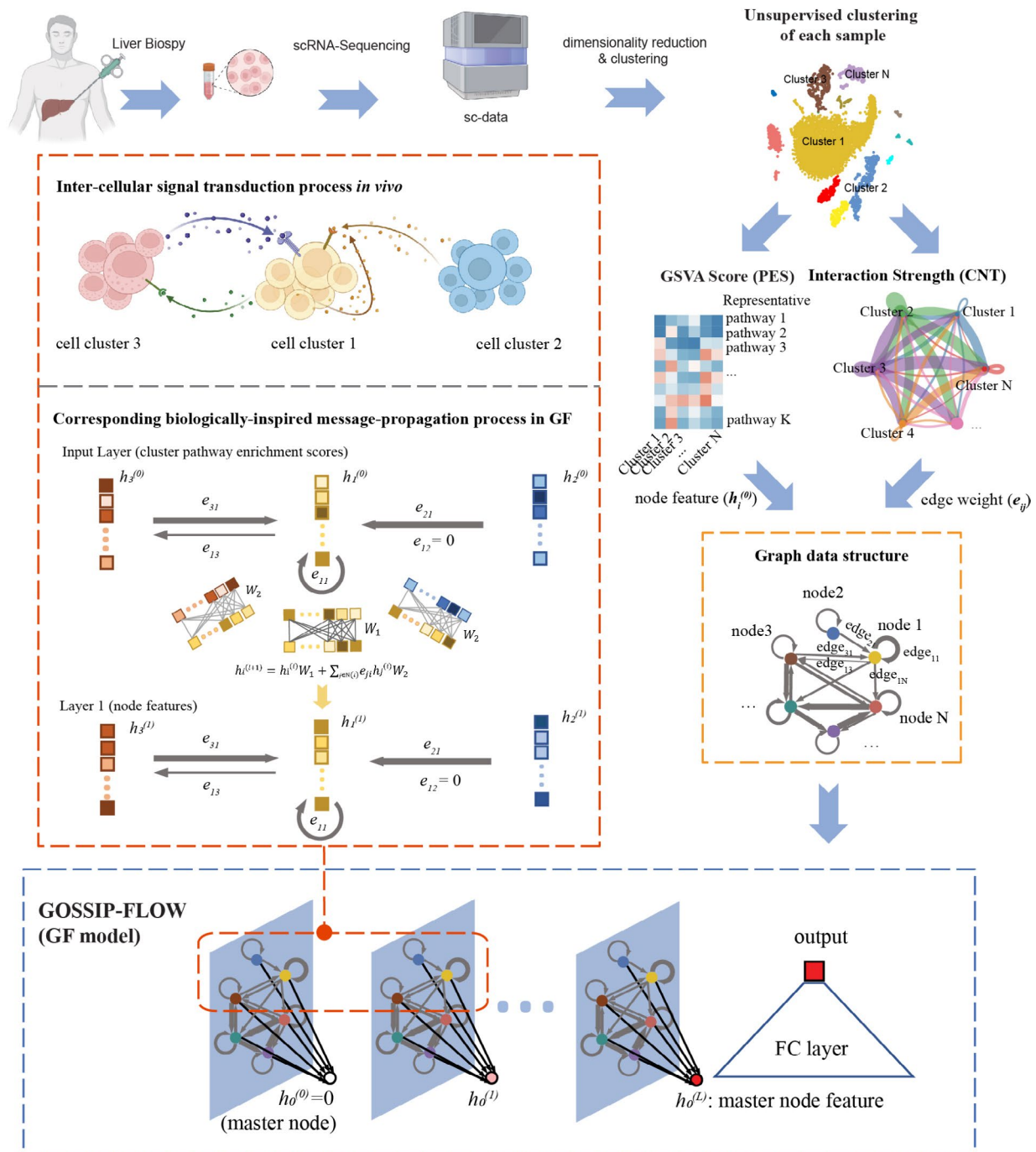
Referring to a suitable number of markers depicting of the microenvironment [28–35], we selected the top 10 statistically significant pathways in each of the 6 major cell types we manually annotated ( $k=10$ ). Considering

that most cell-cell interactions like antigen presenting are unlikely to exceed three layers of trans-cellular signal transduction [36], we set the message propagation layer to three ( $L=3$ ). Based on these two roughly estimated hyperparameters ( $k=10$ ,  $L=3$ ), our GF model achieved an AUC value of 0.75, an ACC of  $70 \pm 2\%$ , and an F1 score of  $0.70 \pm 0.02$ , which is significantly more satisfactory compared to our previous diagnostic framework that followed the traditional scRNA analysis pipeline. (Fig. 2A–D)

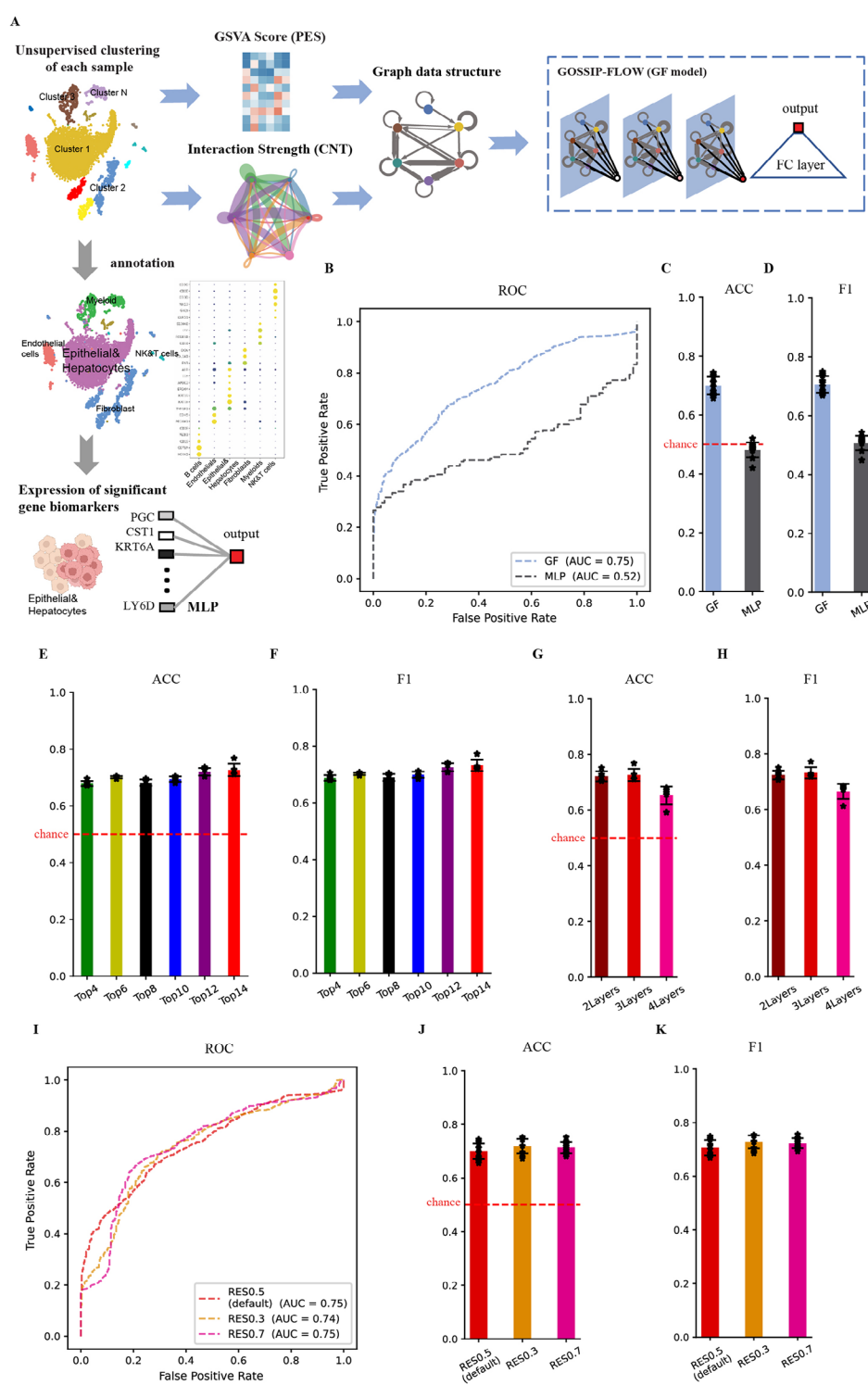
To figure out the general robustness of our framework, we are curious whether our model is sensitive to hyperparameters tuning. Major hyperparameters of GF model are number of top significant pathways selected from each cell type ( $k$ ) during choosing representative pathway selecting and number of DGCN message-propagation layers ( $L$ ). We tested GF models using  $k$  around 10 ( $k \in \{4, 6, 8, 10, 12, 14\}$ ), and  $L$  around 3 ( $L \in \{2, 3, 4\}$ ). Among all the GF models we tested, the best AUC value of 0.79 was achieved when  $k$  was set to 14, and  $L$  was set to 3. An ACC of  $73 \pm 2\%$  (mean  $\pm$  standard deviation (SD)) was obtained for distinguishing malignant tumors from benign tumors, accompanied by an F1 of  $0.74 \pm 0.03$ . Besides, there are no significant difference in accuracies and F1 scores if top 12 or 14 significant pathways were selected in each 6 major cell types, indicating such numbers of pathways are adequate for depicting microenvironment. (Fig. 2E&F and Table S1) GO numbers and names of pathways after excluding repeated and not available ones in GSVA are shown in Table 1, while selecting details are shown in Methods section and Supplementary Material 2. Similarly, 2 or 3 layers of message propagation resulted in similar performance, while 4 layers of message propagation do not achieve as good performance, possibly due to few biological processes involving 4 layers of inter-cellular communication in tumor microenvironments in vivo. (Fig. 2G&H and Table S2)

Some of the parameters in data preprocess, though normally set as default, may still affect the robustness of the framework, especially the CellChat may be affected by the number of input clusters. Hence, we also conducted a sensitivity test on the models' capability of dealing with testing samples of varied number of cell clusters yielded by altering Principal Components Analysis (PCA) resolution (RES). Detailed numbers of clusters of each sample under different resolutions see Supplementary Material 3. Despite the number of clusters changed, our model maintained relatively high performance. (Fig. 2G–I)

In all, GF framework is relatively robust across hyperparameters within a reasonable range, both at the level of connectivity calculation (CellChat), input pathway numbers, and the layer of message propagation.



**Fig. 1** The architecture of the Gossip Flow (GF) framework. The figure illustrates the workflow and architecture of the Gossip Flow (GF) framework for processing scRNA-seq data from tissue samples. **Sample Collection and Sequencing:** Tissue samples are collected from each patient and subjected to scRNA-seq. **Dimensionality Reduction and Clustering:** The gene expression matrix generated from the scRNA-seq data undergoes dimensionality reduction and clustering analysis, resulting in multiple cell clusters for each sample. **Pathway Enrichment Scores (PES) Matrix:** For each cell cluster, GSVA pathway enrichment scores are calculated for  $K$  selected functional pathways, forming the PES matrix, which serves as the node features of the input graph structure. **Inter-cluster communication features (CNT) Matrix:** The unidirectional intercellular communication probability between  $N$  cell clusters is calculated via CellChat based on overexpressed ligands and receptors, forming the CNT matrix, which serves as the edge weights of the input graph structure. **Graph Construction:** Shown in the yellow dotted frame. The PES matrix and CNT matrix are used to construct the input graph structure. A master node (node 0) is added to receive projections from all other nodes, integrating features for readout. **Directed Graph Convolutional Network (DGCN):** Shown in the blue dotted frame. The graph structure is input into a DGCN with  $L$  layers. All DGCN layers share the same set of trainable parameters ( $W_1$  &  $W_2$ ) to avoid overfitting. The master node collects global features. The red dotted frame depicts part of the message-propagation details of node 1 in the first layer of DGCN as an example, which is biologically-inspired by intercellular signal transduction process *in vivo*. **Output Projection:** The feature vector of the master node in the final layer ( $h_0^{(L)}$ ) is projected to the output node via a fully connected (FC) layer. \* Created with BioRender.com



**Fig. 2** GF robustly classified primary liver cancer (ICCA&HCC) from benign tumours (FNH). **A**. Schematic figures of GF framework (blue arrow) and control diagnostic MLP model following traditional single-cell RNA analysis pipeline (grey arrow). **B-D**. Receiver operating characteristic (ROCs), accuracy (ACCs) and F1 scores (F1s) of GF framework and control MLP model following traditional single-cell RNA analysis pipeline (MLP). Each LOOCV test was repeated 10 times. **E, F**. ACCs and F1s of 3-layer directional GF models with different top k differentiative pathways of 6 cell groups. Each LOOCV test was repeated 5 times. **G, H**. ACCs and F1s of testing results of GF models with different DGCN layers and message propagating directions. Each LOOCV test was repeated 5 times. **I-K**. ROCs, ACCs and F1s of the GF models tested on data generate from different PCA resolution (RES = 0.3, 0.5, 0.7). Each LOOCV test was repeated 10 times. All the ACCs and F1s are presented as means  $\pm$  standard deviations (SDs). All the error bars depict SDs, while red horizontal dotted lines represent the level of random chance (50%)



**Table 1** Representative pathways: selected top 14 differentially activated pathways of GO datasets for each 6 distinct cell types within benign and malignant tumor samples

ID	Description	Differen- tiate cell types
GO:0001667	ameboidal-type cell migration	Endo, Fib
GO:0042113	B cell activation	B
GO:0072562	blood microparticle	Epi
GO:0030055	cell-substrate junction	B, Endo, Epi, Fib, Mye
GO:0062023	collagen-containing extracellular matrix	Endo, Epi
GO:0002181	cytoplasmic translation	B, Epi, Mye
GO:0022625	cytosolic large ribosomal subunit	Mye
GO:0022626	cytosolic ribosome	B, Epi, Mye
GO:0005788	endoplasmic reticulum lumen	Fib
GO:0043542	endothelial cell migration	Endo
GO:0003158	endothelium development	Endo
GO:0050673	epithelial cell proliferation	Endo
GO:0045229	external encapsulating structure organization	Fib
GO:0005201	extracellular matrix structural constituent	Fib
GO:0101002	ficolin-1-rich granule	Mye
GO:0006091	generation of precursor metabolites and energy	Epi
GO:0007159	leukocyte cell-cell adhesion	Mye, NKT
GO:0002443	leukocyte mediated immunity	B
GO:0002449	lymphocyte mediated immunity	B
GO:1,903,131	mononuclear cell differentiation	NKT
GO:0006936	muscle contraction	Fib
GO:0003012	muscle system process	Fib
GO:0030099	myeloid cell differentiation	NKT
GO:0050867	positive regulation of cell activation	B, NKT
GO:1,903,039	positive regulation of leukocyte cell-cell adhesion	B, NKT
GO:0050870	positive regulation of T cell activation	NKT
GO:0022407	regulation of cell-cell adhesion	NKT
GO:1,903,706	regulation of hemopoiesis	NKT
GO:0052547	regulation of peptidase activity	Fib
GO:1,901,342	regulation of vasculature development	Endo
GO:0044391	ribosomal subunit	B, Epi, Mye
GO:0005840	ribosome	Epi, Mye
GO:0002040	sprouting angiogenesis	Endo
GO:0003735	structural constituent of ribosome	B, Epi, Mye
GO:0042110	T cell activation	B, Endo, Mye, NKT
GO:0030217	T cell differentiation	NKT
GO:0070820	tertiary granule	Mye
GO:0090130	tissue migration	Endo
GO:0031983	vesicle lumen	Epi
GO:0042060	wound healing	Fib

Abbreviations B, B cells; Epi, epithelial cells and hepatocytes; Endo, endothelial cells; Fib, fibroblast; Mye, myeloid cells; NKT, NK cells and T cells

**GF integrates both intercellular communication features and pathway activation features while predicting**

To verify that our GF model successfully integrated both inter-cluster communication features (CNT matrix) and pathway activation features (PES matrix) calculated from raw scRNA-seq data during diagnosis, we validate the performance of our model on partially polluted testing data. To assess the significant of inter-cellular communication (CNT matrix) to our diagnostic framework, we replace the edge weight of input graph data with Gaussian noise generated randomly of the same shape, while the node feature still valid PES matrix (PES). (schematic figure see Fig. 3A). An AUC of only 0.46, an ACC of  $49\pm1\%$ , and an F1 of  $0.51\pm0.01$  were achieved on PES data, indicating that inter-cluster communication features are essential for the efficacy of the GF model (Fig. 3B-D).

Similarly, we tested graph data preserving only edge weight preserved, while node features replaced by Gaussian noise generated randomly rather than valid PES matrix (CNT). An AUC of 0.56, an ACC of  $51\pm2\%$ , and an F1 of  $0.53\pm0.01$  were achieved, indicating that the pathway activation features are also essential for the efficacy of the GF model (Fig. 3B-D).

As designed, both the pathway enrichment features and inter-cluster communication features are essential for the high performance of GF models.

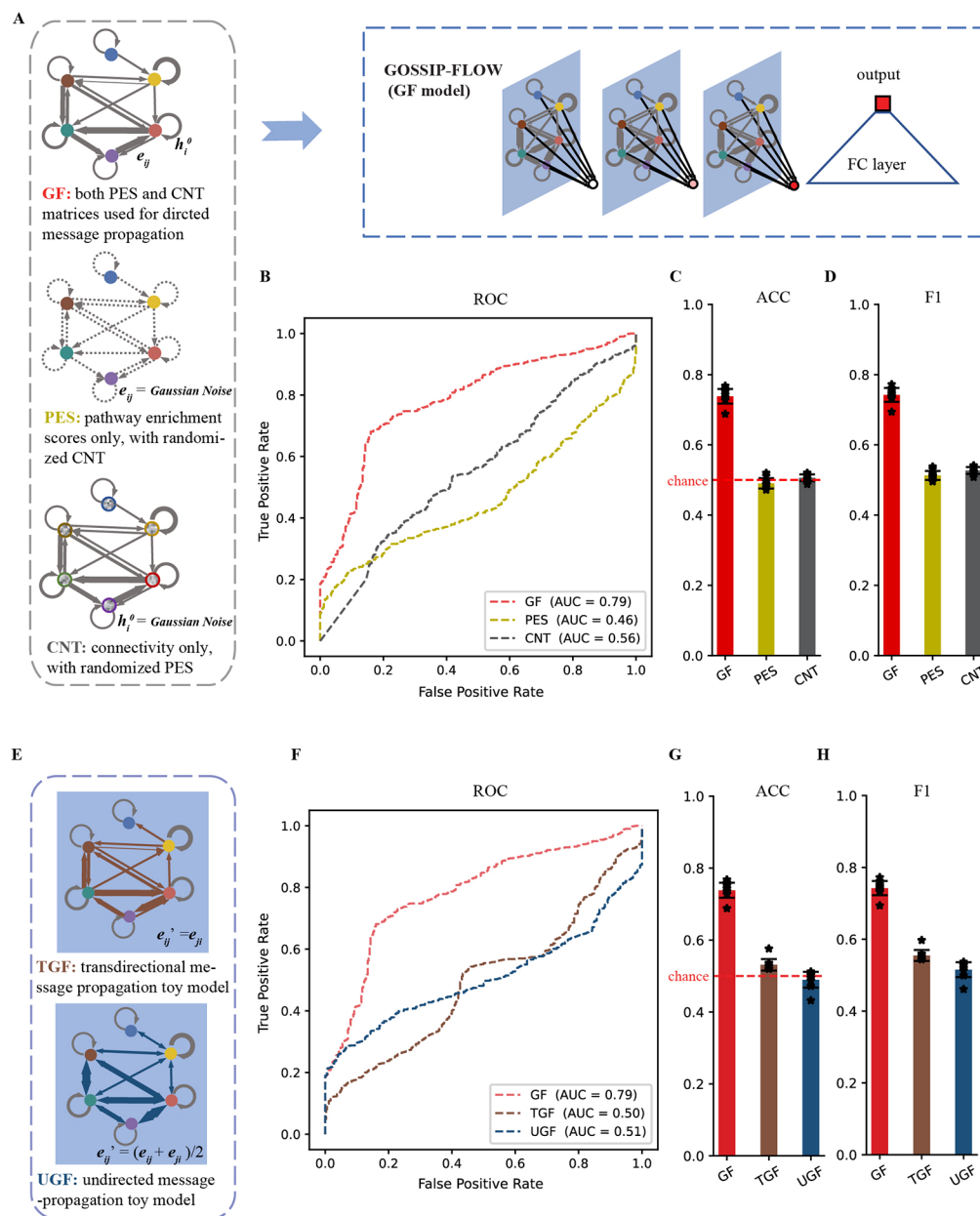
**GF requires a biologically plausible message propagation direction**

To further investigate how GF collects features during message propagation, we constructed toy models with hyperparameters identical to those of GF. In contrast, the message was propagated in an undirected manner among all the nodes except for the master node, which kept receiving inputs from other nodes for readout (undirected GF, UGF). (schematic figure see Fig. 3E). An AUC of only 0.51, an ACC of  $49\pm2\%$ , and an F1 of  $0.52\pm0.02$  were achieved by UGF. We also reversed the message propagation direction of GF (trans-directed GF or TGF). An AUC of only 0.5, an ACC of  $53\pm2\%$ , and an F1 of  $0.55\pm0.02$  were achieved (Fig. 3F-H).

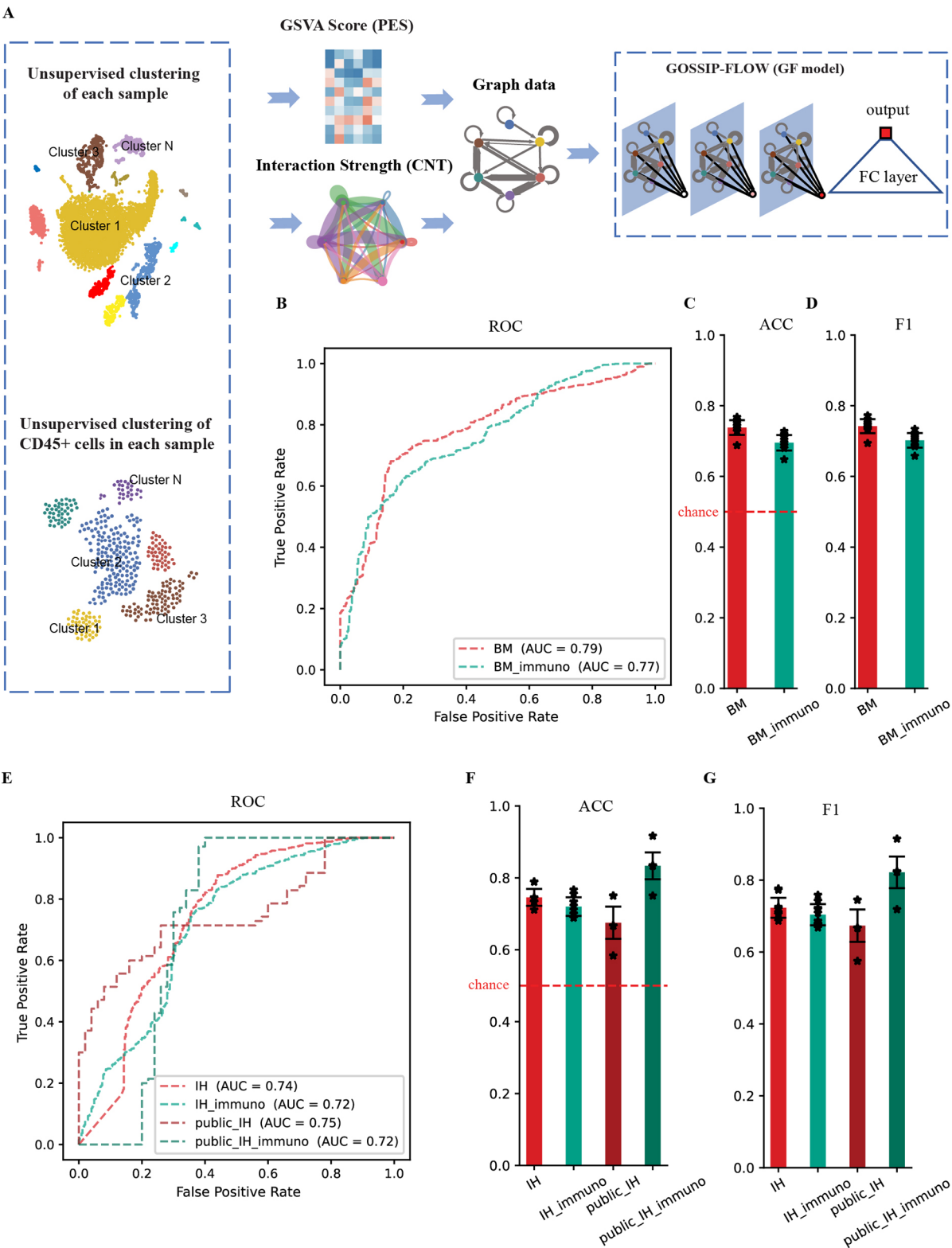
The toy models failed to perform diagnosis when the message propagation process was set either undirectedly or trans-directedly, which further suggests that GF relies on biologically plausible message propagation processes while collecting and analysing sparsely distributed features within the transcriptome.

**GF can be applied on transcriptome data of immunocytes only**

Emulating a clinical scenario where a biopsy failed to obtain an adequate number of tumor cells, such as a biopsy missing the small target getting only paratumor



**Fig. 3** The GF model is a biologically inspired network that successfully integrates pathway enrichment features and unidirectional intercellular communication probabilities while performing diagnoses. **A**. Schematic figures showing GF models tested on normal data, CNT polluted data (PES) and PES polluted data (CNT). **B–D**. ROCs, ACCs and F1s of GF models tested on normal data and data partially polluted (PES and CNT). Each LOOCV test was repeated 10 times. Without edge weights calculated from intercellular communication probabilities, PES reached an AUC of 0.46, an ACC of  $0.49 \pm 0.01$ , and an F1 of  $0.51 \pm 0.01$ . Similarly, without pathway enrichment scores, CNT reached an AUC of 0.56, an ACC of  $0.51 \pm 0.02$ , and an F1 of  $0.53 \pm 0.01$ . Integrating both pathway enrichment features and inter-cluster communication features is of vital significance for the efficacy of the GF model. **E**. Schematic figures showing toy GF models with transdirectional message-propagation direction (TGF) and undirected message-propagation direction (UGF). Other hyperparameters were set identical to those of the best GF model. **F–H**. ROCs, ACCs and F1s of GF models, TGF toy model and UGF toy model. UGF achieved an AUC of merely 0.51, an ACC of  $0.49 \pm 0.02$ , and an F1 of  $0.52 \pm 0.02$ , while TGF achieved an AUC of merely 0.50, an ACC of  $0.53 \pm 0.02$ , and an F1 of  $0.55 \pm 0.02$ , indicating that GF relies on a biologically plausible message propagation direction. All the models were tested on ten rounds of LOOCV tests. All the ACCs and F1s are presented as means  $\pm$  standard deviations (SD). All the error bars depict SDs, while red horizontal dotted lines represent the level of random chance (50%)



**Fig. 4** (See legend on next page.)



(See figure on previous page.)

**Fig. 4** The GF model can be applied to scRNA data of immunocytes and distinguish subtypes of primary liver tumours. **A.** Schematic figures showing GF models applied on scRNA data of all cells and immunocytes only. **B-D.** ROCs, ACCs and F1s of GF model applied on transcriptome data consisting of whole cells and immunocytes (BM\_immuno). An AUC of 0.77, an ACC of  $0.70 \pm 0.02$  and an F1 of  $0.70 \pm 0.02$  were achieved when testing on immunocytes, verifying the capacity of the GF framework to capture the systemic features of tumour microenvironments and its robustness in tackling tissue samples with high heterogeneity. **E-G.** ROCs, ACCs and F1s of the best GF model in terms of distinguishing subtypes (iCCA and HCC) of primary liver tumours when applied on transcriptome data consisting of whole cells (IH)/immunocytes (IH\_immuno). An AUC of 0.74, an ACC of  $0.75 \pm 0.02$  and an F1 of  $0.72 \pm 0.03$  were achieved by IH model of whole cells. An AUC of 0.72, an ACC of  $0.72 \pm 0.03$  and an F1 of  $0.70 \pm 0.03$  were achieved by IH model of immunocytes. An AUC of 0.75, an ACC of  $0.67 \pm 0.04$  and an F1 of  $0.67 \pm 0.05$  were achieved when testing on public data of whole cells. An AUC of 0.72, an ACC of  $0.83 \pm 0.04$  and an F1 of  $0.82 \pm 0.04$  were achieved on public data of immunocytes. For all the models tested on internal datasets via cross-validation, tests were performed through ten rounds of LOOCV tests. Each test performed on public data was repeated 10 times. All the ACCs and F1s are presented as means  $\pm$  standard deviations (SDs). All the error bars depict SDs, while red horizontal dotted lines represent the level of random chance (50%)

tissues infiltrated by immunocytes, we applied the GF framework on scRNA-seq data consisting solely of immunocytes selected from each sample.

Remarkably, despite lacking transcriptome data from tumor cells, GF framework on scRNA-seq data of immunocytes (BM\_immuno) yielded promising results: An AUC of 0.77, an ACC of  $70 \pm 2\%$ , and an F1 of  $0.70 \pm 0.02$  were achieved. (Fig. 4B-D).

Our findings substantiate our hypothesis that GF does not rely solely on tumor cell markers for predictions, unlike many conventional molecular pathology diagnostics. Instead, it captures differences between the microenvironments of benign and malignant tumors in an integrative manner. This is evidenced by the framework's sustained diagnostic performance, even when provided with transcriptome data exclusively from immunocytes, absent from any input from tumor cells.

#### GF can be applied to the classification of iCCA from HCC

We further applied the GF framework to distinguish iCCA from HCC (IH). We performed the similar sensitivity test as previous. Performances of models with different number of pathways, number of layers and opposite message propagation directions are shown in Table S3&4. GO numbers and names of pathways after excluding repeated and not available ones in GSVA are shown in Table 2, while selecting details are shown in Methods section and Supplementary Material 4. The best AUC of 0.74 was achieved when the number of DGCN layers in the model was set to 2. An ACC of  $75 \pm 2\%$  and an F1 of  $0.72 \pm 0.03$  were obtained by this model. Similarly, we tested the GF framework on only immunocytes (IH\_immuno) selected from each sample. An AUC of 0.72, an ACC of  $72 \pm 3\%$ , and an F1 of  $0.70 \pm 0.03$  were achieved when testing on immunocytes (Fig. 4E-G).

To further validate the generalization of our model, we tested the GF model trained on our internal dataset directly on public datasets (public\_IH). An AUC value of 0.75 and an ACC of  $67 \pm 4\%$  were obtained for distinguishing iCCA from HCC, accompanied by an F1 of  $0.67 \pm 0.05$ . We also tested the GF framework on only immunocytes selected from each sample of public data (public\_IH\_immuno). An AUC of 0.72, an ACC of

$83 \pm 4\%$ , and an F1 of  $0.82 \pm 0.04$  were achieved when testing on immunocytes. (Fig. 4E-G).

#### Discussion

Tailored for distilling more robust representations of tumor microenvironments through scRNA-seq, our GF framework successfully demonstrates the feasibility of using ANN models to automatically diagnose benign and malignant tumors. Furthermore, our framework can also effectively distinguish HCC from iCCA. Even importantly, we apply our framework on scRNA data of only immunocytes to identify the property of liver tumors. Our framework represents a novel and automated diagnostic method for liver tumor, other than serum biomarkers, radiologic and pathologic methods.

Our study proposes a data-processing pipeline to distill diagnostic information from scRNA data. Despite the high-dimensional and high-resolution advantages of scRNA-seq, there are challenges in feeding the original scRNA expression matrix into appropriate diagnostic models. On one hand, scRNA expression matrix typically contain a vast number of genes and cells but a relatively small number of samples, which can lead to overfitting. Additionally, at the single cell and gene level, noise is introduced during the biological data generation process, either from sample preparation or experimental operations [37]. Consequently, many gene biomarkers exhibit poor stability and are rarely replicated in other studies, making it challenging to transfer findings from one dataset to another due to overfitting and the low signal-to-noise ratio in many gene expressions [38]. Therefore, using these low reproducible gene biomarkers is not a proper solution for diagnosis directly on scRNA data. Hence, we employ dimensional-reduction, clustering, and pathway enrichment, to compress the sparse and non-stable (low reproducible) original expression matrix into a non-sparse and more stable format, which is more suitable for modelling. Furthermore, by calculating enrichment scores and inter-cluster communication features, we distil biologically-plausible information based on prior knowledge. Non-Euclidean nature serves as another challenge for modelling on scRNA data. After pre-processing, scRNA-seq data at the patient level still

**Table 2** Representative pathways: selected top 12 differentially activated pathways of GO datasets for each 6 distinct cell types within icCA and HCC samples

ID	Description	Differentiate cell types
GO:0072562	blood microparticle	Epi
GO:0045296	cadherin binding	Epi
GO:0005911	cell-cell junction	Epi
GO:0031589	cell-substrate adhesion	Fib
GO:0030055	cell-substrate junction	B, Endo, Epi, Fib
GO:0051087	chaperone binding	NKT
GO:0062023	collagen-containing extracellular matrix	Endo, Fib
GO:0002181	cytoplasmic translation	B, Fib
GO:0022625	cytosolic large ribosomal subunit	B, Fib
GO:0022626	cytosolic ribosome	B, Epi, Fib, Mye
GO:0022627	cytosolic small ribosomal subunit	B
GO:0030139	endocytic vesicle	Mye
GO:0030666	endocytic vesicle membrane	Mye
GO:0005788	endoplasmic reticulum lumen	Fib
GO:0043542	endothelial cell migration	Endo
GO:0045229	external encapsulating structure organization	Fib
GO:0005201	extracellular matrix structural constituent	Fib
GO:0031072	heat shock protein binding	NKT
GO:0007599	hemostasis	Endo
GO:0005178	integrin binding	Endo
GO:0007159	leukocyte cell-cell adhesion	Mye
GO:0002443	leukocyte mediated immunity	NKT
GO:0050900	leukocyte migration	Mye
GO:0051346	negative regulation of hydrolase activity	Epi
GO:0002683	negative regulation of immune system process	NKT
GO:0050867	positive regulation of cell activation	B
GO:0001819	positive regulation of cytokine production	Mye, NKT
GO:0032103	positive regulation of response to external stimulus	Mye
GO:0050878	regulation of body fluid levels	Endo
GO:0022407	regulation of cell-cell adhesion	Mye, NKT
GO:0052547	regulation of peptidase activity	Epi
GO:0050863	regulation of T cell activation	NKT
GO:0044391	ribosomal subunit	B
GO:0005840	ribosome	B
GO:0003735	structural constituent of ribosome	B
GO:0042110	T cell activation	Mey
GO:0051082	unfolded protein binding	NKT
GO:0031983	vesicle lumen	Endo, Epi, Mye, NKT
GO:0042060	wound healing	Endo

**Abbreviations** B, B cells; Epi, epithelial cells and hepatocytes; Endo, endothelial cells; Fib, fibrocytes; Mye, myeloid cells; NKT, NK cells and T cells

vary in size. Such non-Euclidean data is inappropriate for traditional ANNs such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), which require fixed-size inputs and Euclidean spatio-temporal patterns [39]. GNNs, however, can handle graphs of various sizes and structures, making them ideal for processing graph structure converted from scRNA data [40]. Therefore, our implementation of GNN draw more effective diagnosis from scRNA data, compared with traditional scRNA analytical pipeline.

Our findings also highlight the diagnostic value of inter-cellular communication in liver tumor microenvironment, which is hard to explore using other diagnostic methods. There exist wide communications resulting direct, dual effects of various tissue and immune cells on tumor cells within the tumor microenvironment. Such as CD4+cytotoxic T cells can directly induce apoptosis in tumor cells through the release of IFN- $\gamma$  [41]. Conversely, it also shows how tumor cells can impair the cytotoxicity of CD8+cytotoxic T cells by secreting factors like lactic acid and TGF- $\beta$  [36]. Complex communication networks govern tumor progression and therapeutic responses in the liver tumor microenvironment, which also reflect the essential property of the tumor [25]. Notably, the direction of message propagation is another decisive factor accounting for our model's efficacy. In traditional graph neural network tasks using datasets like Cora-ML, PubMed, CiteSeer, and DBLP, prediction accuracies typically show minor differences regardless of message propagation direction. DGCN and classic undirected graph convolutional networks (UGCNs) demonstrated stable accuracies, with less than a 5% decrease when message direction was ignored [42]. Our biologically-inspired GF models exhibited a striking contrast. They lost effectiveness when intentionally reversing message propagation from ligand-to-receptor to receptor-to-ligand (TGF), or when substituting DGCN with an undirected GCN (UGF). This striking phenomenon underscores the critical importance of tracing upstream intercellular communication within tumor microenvironment. Our observations suggest that propagating scRNA-seq data within graph structures of the ligand-receptor direction, which mimic intercellular signal cascade pathways in vivo [43, 44], led to integration and amplification of diagnostic information encapsulated within discrete features of tumor microenvironments. Therefore, such biologically-inspired design not only enhances the interpretability and effectiveness of diagnostic predictions but also underscores the significance of considering biological context in computational diagnostic frameworks based on scRNA-seq data.

Our framework represents a paradigm shift in the field of molecular pathology diagnostics. Instead of relying on a few tumor indices or signatures as most of the previous

molecular diagnostic models do, our framework pictures the microenvironment as an integrated system. This approach allows our framework to exhibit robust performance even on scRNA data from immunocytes only. Previous studies have reported phenotypes of immunocytes varied in different tumor types, such as distinct “intermediate” exhausted CD8<sup>+</sup>T cells and “nonclassic” plasmacytes in extrahepatic cholangiocarcinoma [45], and higher proportions of regulatory T cells with resting dendritic cells in HCC tissues [46]. Phenotypes and inflammatory mediators within the tumor microenvironment also altered when adenomatous stage progressing into the cancerous stage in colorectal tumors [47]. Our findings further suggest the state of immunocytes carries key information indicating the property of tumor. Considering that immune cells possess the ability to migrate to tissues surrounding tumor, future work needs to extend the idea of a “comprehensive and systemic understanding” into similar diagnostic models based on immunocytes from alternative tissue samples. These include biopsies that collect necrotic tumor tissue infiltrated by immunocytes, and paratumoral tissue samples in cases where the targeted tissues are too small or located in hard-to-reach anatomical positions (i.e. adjacent to large blood vessels or the pericardium). Additionally, we may implement our diagnostic framework in tumor liquid biopsies from immunocytes collected in peripheral blood samples [48]. These immunocytes could also carry information from the tumor microenvironment, aiding in diagnosis.

Powered by more training samples in the future, distilling more systemic and integrative microenvironment features is promising when more biological-plausible features, such as detailed inter-cellular ligand-receptor-specific messages, fed into GF framework. Furthermore, integrating our ANN-based framework into multimodal networks is relatively straightforward. By combining features distilled from scRNA-seq data with information from radiology, serum markers, and pathology data, a more systemic and reliable diagnosis can be achieved. However, further validation, such as sensitivity testing on large-scale population datasets, is needed to ensure the practicality of these diagnostic methods.

While the scRNA-seq-based diagnostic framework presented in this study shows promising results, it is important to acknowledge its potential limitations compared to standard pathological approaches. Currently, scRNA-seq technology is more expensive and requires more infrastructure than traditional histopathological examination and immunohistochemistry. Additionally, the turnaround time from sample collection to final diagnosis may be longer for the scRNA-seq-based approach. The relatively small training dataset in this pilot study could also limit the model's ability to fully capture the complex tumor microenvironment. These practical

constraints may restrict the widespread clinical adoption of the scRNA-seq-based diagnostic method, especially compared to the standard pathological evaluations by experienced clinicians, which remain the primary method for diagnosing primary liver tumors in most clinical settings. Future research is needed to optimize the scRNA-seq-based approach (such as developing more straightforward statistical approaches) and clearly demonstrate its advantages over existing methods before it can be widely implemented in routine clinical practice.

## Methods

### *Patient sample collection*

The Ethics Committee of Zhongshan Hospital, Fudan University granted the study's ethical approval (B2022-480R), and written informed consent was obtained from each patient. We selected 25 patients with primary liver tumors who underwent surgical resection or percutaneous liver biopsy at Zhongshan Hospital between March 2021 and October 2022. Tumor tissues were collected and processed within 90 min after surgery. 7 cases of HCC, 11 cases of iCCA, and 7 cases of FNH were identified by well-experienced pathologists in Zhongshan Hospital. The first two were malignant tumors, while the latter was benign. This was a non-interventional study, so blinding and randomization principles were not implemented. Detailed clinical information of all patients is available in Table S5.

### *Single-cell process of samples*

For the quality check and counting of single cell suspension, the cell survival rate is generally above 85%. The cells that have passed the test are washed and resuspended to prepare a suitable cell concentration of 700~1200 cells/ $\mu$ l for 10x Genomics Chromium™. The system is operated on the machine. GEMs (Gel Bead in Emulsion) were constructed for single cell separation according to the number of cells to be harvested. After GEMs were normally formed, GEMs were collected for reverse transcription in a polymerase chain reaction (PCR) machine for labelling. The GEMs were oil-treated, and the amplified cDNA was purified by magnetic beads, and then subjected to cDNA amplification and quality inspection. The 3' Gene Expression Library was constructed with the quality-qualified cDNA. After fragmentation, adaptor ligation, sample index PCR, etc., the library is finally quantitatively examined. The final library pool was sequenced on the Illumina Novaseq 6000 instrument using 150-base-pair paired-end reads.

### *scRNA-seq data pre-processing*

R 4.2.1 was used for the whole down-streaming scRNA-seq analysis. We assessed cell quality using the following criteria: (1) the number of total count per cell (library

size) was below 50,000; (2) the number of detected genes was above 300 and below 7,500; (3) the percentage of mitochondrial genes was below 10; (4) the percentage of hemoglobin genes was below 0.1. After quality control, a total of 242,598 cells were retained for downstream analysis. Detailed preprocessing results is available in Fig. S1 & Table S6.

We merged the expression matrices of each sample and performed data integration using the harmony package (version 0.1.1). Plotting of the integrated data revealed good mixing of immune and stromal cells between different samples, therefore, significant batch effects were excluded. To expand our input data and improve the model's training, we resampled each sample, randomly selecting 5,000 cells from each sample and using them to construct a new Seurat object. This process was repeated ten times, resulting in 250 new samples.

#### **Public data collection**

To further validate the generality of our model, some public scRNA-seq datasets of primary liver malignant tumor from GEO database were used for test set. Raw data of GSE162616, GSE166635, GSE189175, GSE125449 and GSE138709, were downloaded and integrated for further process.

#### **Cell unsupervised clustering**

To extract useful information from single-cell data for model construction, each sample's single-cell data was performed dimensionality reduction and clustering analysis using the Seurat package (version 4.3.0). After scaling the sample data, we used the FindVariableFeatures function to select 2,000 highly variable genes, and then used the RunPCA function to calculate the top 50 principal components (PCs) of the data. We then selected the first 20 PCs and performed unsupervised clustering on the cells using the FindNeighbors and FindClusters functions (default resolution=0.5, and we also set resolution to 0.3 and 0.7 for sensitivity test). UMAP was used for dimensionality reduction of the sample data.

#### **Cell annotation**

Cell annotation was conducted as part of our analysis to categorize cells into 6 major cell types. While this step was not explicitly integrated into the GF framework, it is a widely accepted method for characterizing the liver tumor microenvironment which we used to identify differentially expressed genes associated with each annotated major cell type and elucidated representative genes and pathways.

The 6 major cell types were annotated in the liver tumor microenvironment using known cell markers reported in previous articles [28–35]. NK&T cells were identified by the presence of CD3D, CD3E, CD3G,

KLRD1, GNLY, and NKG7. B&Plasma cells were identified by the presence of MS4A1, CD79A, CD19, MZB1, and CD38. Myeloid cells were identified by the presence of CD14, FCGR3A, LYZ, and S100A8. Endothelial cells were identified by the presence of PECAM1, CDH5, and TM4SF1. Fibroblasts were identified by the presence of FN1, COL1A1, and DCN. Epithelial&hepatocyte were identified by the presence of KRT18, KRT19, EPCAM, APOC3, TTR, and ALB. NK&T cells, B&Plasma cells, and Myeloid cells were defined as immune cells, while other cells were non-immune cells. (Fig. S2)

#### **Control diagnostic framework based on traditional scRNA analysis pipeline**

After unsupervised clustering and cell annotation, we utilized the FindMarkers function to obtain top 10 differentially expressed genes in tumor cell (epithelial&hepatocyte type mentioned above) of benign and malignant tumor samples. Detailed names of gene were shown in Supplementary Material 5, which were selected based on the absolute value of avg\_log2FC. Subsequently, we extract the average expression of these 10 genes in epithelial&hepatocyte cell types of every sample and feed them into a 2-layer MLP. MLP was trained and validated in exactly the same methods shown in Training, internal and external validating details in [Methods](#) section, ensuring strict control.

#### **Features and pathways selecting**

In order to investigate the biological functional differences between different cell types in benign and malignant tumors, we identify differentially expressed genes associated with each annotated major cell type and elucidated representative pathways. We utilized the FindMarkers function to obtain differentially expressed genes in 6 distinct cell populations mentioned above of benign and malignant tumor samples. Subsequently, we performed pathway analysis on these differentially expressed genes and selected top  $k$  ( $k \in \{4, 6, 8, 10, 12, 14\}$  for our sensitivity analysis) statistical-significant (with smallest  $p$  value) pathways for each major cell type in GO datasets as representative pathways of tumor-microenvironment [22]. After discarding duplicated pathways and excluded those pathways not available in GSVA pathway enrichment package, we acquire the selected pathways we used to generate pathway enrichment scores (PES) matrix of each sample. For later analysis, we fixed  $k=14$  according to the accuracy in sensitivity test, finally obtaining 40 representative pathways ( $K=40$ ). We also conducted the same analysis between HCC and ICCA samples. We fixed  $k=12$  and eventually obtained 39 representative pathways ( $K=39$ ). Detailed information about these pathways (including selecting result in each step) can be found in Supplementary Material 2 & 4.



### Calculation of pathway enrichment score (PES) matrix

To obtain the feature matrix for each sample in our GNN model, we performed pathway enrichment analysis using the GSVA package (version 1.46.0). We downloaded the C5 human gene set (GO gene set) from the msigdb package (version 7.5.1) and used the GSVA function to calculate GSVA pathway enrichment scores for functional pathways of each cell cluster in each sample. We then pick the enrichment score of our previously selected pathways to generate PES matrix for each sample. Each column of PES matrix represents the PES of one cell cluster, which will be used as the node feature of layer 0 ( $h_i^0$ ) of input graph data.

### Calculation of cluster communication weight

To investigate the intercellular communication among different cell populations in single-cell samples, we performed cell communication analysis using the CellChat package (version 1.6.1). We imported the Secreted Signaling dataset from CellPhoneDB.Human to analyze the interactions of receptor-ligand pairs between cell populations. We identified ligands or receptors overexpressed in a cell group during data preprocessing, and projected the gene expression data onto a protein-protein interaction (PPI) network. We identified interactions between overexpressed ligands and receptors, and calculated the intercellular communication probability between cell populations using the computeCommunProb function. The communication probability helped infer the biological significance of cell-cell communication strength. Finally, we calculated the average strength of all receptor-ligand interactions between different cell clusters (strength of ligands on cluster  $i$  and receptors on cluster  $j$  denotes as  $e_{ij}$ ) and constructed a cluster communication weight for each sample.

### Model architecture

The input cluster communication weight and cluster feature matrix were analysed by PyTorch (version 1.13.1) and PyG (version 2.2.0).

The GF model was built up with two parts. First is the  $L$  layers of DGCN with initial node features set as cluster feature matrix for each cell clusters and edge weights set as cluster communication weight.

Before feeding into the network, cluster communication weight was normalized by dividing the mean edge weight of the total cluster communication weight matrix, while the node features were normed within each feature vector of a node via calculating z-scores.

The following calculation show the detailed message-propagation process from other nodes ( $N$ ) to nodes  $i$  in layer  $l$  of the DGCN:

$$h_i^{(l+1)} = h_i^{(l)} W_1 + \sum_{j \in N(i)} e_{ji} h_j^{(l)} W_2$$

We also use a latent master node (node 0,  $h_0^0$  padded with **zeros**), which is connected to every input node in the graph with the edge weight set as mean edge weight of the total inter-cluster communication features. The feature vector of master node ( $h_0^l$ ) serves as a global scratch space that each node writes to but not read from ( $e_{i0} = 1$ ,  $e_{0i(i \neq 0)} = 0$ ). This allows information to collect from long distances during the propagation phase. And final output was projected from the feature vector of master node only ( $h_0^L$ ), through a single layer of fully connect layer.

In each layer of DGCN, size of the hidden layer matches exactly with number of pathways selected (input dimension), and all the layers share one set of trainable parameters ( $W_1$  and  $W_2$ ). For GF models distinguishing benign and malignant tumors, 40 pathways were selected, resulting in 3241 trainable parameters in total. For GF models distinguishing iCCA and HCC, 39 pathways were selected resulting in 3082 trainable parameters in total. These light-weighted structures enabled the iteration of GF models on limited scRNA-seq data, preventing immediate overfitting.

### Data augmentation

For better performance, we randomly discard one node in each sample to generate a new graph and added to our training set, as the simplest augmentation strategy for graph data.

Since different categories had different ratio of positive and negative samples, such data augmentation was repeated for different times respectively to eliminate the unfavorable effect caused by unbalanced training data. (i.e. if testing set consists of benign samples from one of the patients, each benign sample in training set repeated data augmentation for 17 times while each malignant sample repeated for 5 times). After data augmentation, GF models distinguishing benign and malignant tumors was trained on approximately 1080 to 1190 samples, while GF models distinguishing iCCA and HCC was trained on approximately 660 to 770 samples. Data augmentation fulfilled the potential of scRNA-seq data, which satisfied the criterion for successful iterations of these ANN models.

Data augmentation were only performed on samples of cases used for training. Any possible leakage of information in both resampling stage and data augmentation stage was strictly prohibited in both internal and external validation process, by limiting all the data acquired from one case is limited to a fold in each cross-validation.



### Training, internal and external validating details

For internal validation, GF adopted leave-one-out-cross-validation (LOOCV) test in our own datasets. 5 samples (or data of different dimensional reduction resolution, polluted PES and CNT data derived from these 5 samples) from each case will be separated and considered as the test set, while the full/immunocytes scRNA-seq data of  $5 \times (n-1)$  samples resampled from the rest  $n-1$  cases are considered as the training set. At the level of the cases, train and test sets were completely separated to avoid any possible information leakage.

During training, BCEWithLogitsLoss for binary classification was calculated for back propagation. We use Adam optimizer with an initial learning rate of 0.001 and a batch size of 4. Total epoch number was 100 and the learning rate was adjusted in each epoch following cosine annealing.

After 100 epochs, testing was performed. The output probability was calculated through sigmoid function. The predicted category was defined as probability exceeding 0.5. All the predicted label and probability were documented for future analysis. Considering class imbalances in training samples, both ACC and F1 were used as metrics of assessment in each LOOCV test. Each LOOCV test was performed 10 times and area under curve (AUC) was calculated from total results, in order to further assess the performance and reliability of GF models.

As for external validation, models are trained on our own datasets containing all the cases, while the independent public datasets were used as test sets. Each test also performed 10 times. Other training details were set identical to internal validation.

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12967-024-05670-1>.

Supplementary Material 1

Supplementary Material 2

Supplementary Material 3

Supplementary Material 4

Supplementary Material 5

Supplementary Material 6

Supplementary Material 7

### Acknowledgements

This study was supported by the Program of Shanghai Academic Research Leader (22XD1402700), the Key Disease Joint Research Program of Xuhui District (XHLHGG202103), the National Key Research and Development Program of China (2019YFC1316000), the National Natural Science Foundation of China (82273234, 82403964), Beijing Mutual Care Public Welfare Foundation (GDXZ-08-05), Sanming Project of Medicine in Shenzhen (SZSM202003009), Fellowship from the China Postdoctoral Science Foundation (2024M750532), National Science and Technology Major Project of China (2023ZD0511500), the Outstanding Resident Clinical Postdoctoral Program of Zhongshan

Hospital Affiliated to Fudan University (to Lu Jia-Cheng); Youth Fund of Zhongshan Hospital Affiliated to Fudan University (to Lu Jia-Cheng). We thank the Medical Science Data Center of Fudan University for the data analysis support.

### Author contributions

Dao-Han Zhang and Chen Liang finished the data preprocessing, the construction of the diagnostic framework and wrote original draft. Jia-Cheng Lu, Xiao-Yong Huang, Lei Yu, Xian-Long Meng, and Xiao-Jun Guo collected tumor samples and analysed the data of single cell sequencing. Hai-Ying Zeng, Zhen Chen, Lv Zhang, Yan-Zi Pei, Shu-Yang Hu, Mu Ye, Qi-Man Sun, Guo-Huan Yang, Jia-Bin Cai, Pei-Xin Huang, Lv Zhang, Ying-Hong Shi, and Ai-Wu Ke collected clinical data of patients enrolled. Yuan Ji was in charge the pathological diagnosis. Jia-Cheng Lu, Guo-Huan Yang, Qi-Man Sun, Yujiang Geno Shi, Guo-Ming Shi, Jian Zhou, and Jia Fan led this plan and discussed the manuscript.

### Data availability

Data relevant to this study are accessible from the authors under restricted access according to our clinical trial protocol, which enables us to share de-identified information with researchers from other institutions but prohibits us from making it publicly available. Access can be granted upon reasonable request. Any data provided must be kept confidential and cannot be shared with others unless approval is obtained. Source data and code to recreate the figures in the manuscript will be publicly released with code upon publication of the manuscript.

### Code availability

Code and source data to replicate the main findings of this study can be found on Supplementary Material 6 & 7.

### Declarations

#### Conflict of interest

All authors claim that there is no conflict of interest.

#### Author details

<sup>1</sup>Department of Liver Surgery and Transplantation, Zhongshan Hospital, Fudan University, Shanghai 200032, China

<sup>2</sup>Liver Cancer Institute, Fudan University, Shanghai 200032, China

<sup>3</sup>Key Laboratory of Carcinogenesis and Cancer Invasion, Ministry of Education of the People's Republic of China, Shanghai 200032, China

<sup>4</sup>Department of Pathology, Zhongshan Hospital, Fudan University, Shanghai 200032, China

<sup>5</sup>Clinical Research Unit, Institute of Clinical Science, Zhongshan Hospital of Fudan University, Shanghai 200032, China

<sup>6</sup>Department of Liver Surgery, Shanghai Geriatric Medical Center, Fudan University, Shanghai 200032, China

Received: 14 April 2024 / Accepted: 12 September 2024

Published online: 01 October 2024

### References

1. Sung H, Ferlay J, Siegel RL, et al. Global Cancer statistics 2020: GLOBOCAN estimates of incidence and Mortality Worldwide for 36 cancers in 185 countries. *CA Cancer J Clin*. 2021;71(3):209–49.
2. Valle JW, Kelley RK, Nervi B, Oh D-Y, Zhu AX. Biliary tract cancer. *Lancet*. 2021;397(10272):428–44.
3. Singal AG, Lok AS, Feng Z, Kanwal F, Parikh ND. Conceptual model for the hepatocellular carcinoma screening continuum: current status and research agenda. *Clin Gastroenterol Hepatol*. 2022;20(1):9–18.
4. National Comprehensive Cancer Network Guidelines for Hepatobiliary Cancers. (Version 1.2022). <https://www.nccn.org/guidelines/guidelines-process/transparency-process-and-recommendations/GetFileFromFileManagerGuid?FileManagerGuidId=bc919db1-70ff-4f66-9ddd-a0a5c6506878>. Published 2022. Accessed May 4, 2023, 2023.
5. Hu C, Xia T, Cui Y, et al. Trustworthy multi-phase liver tumor segmentation via evidence-based uncertainty. *Eng Appl Artif Intell*. 2024;133:108289.

6. Tzartzeva K, Obi J, Rich NE, et al. Surveillance imaging and alpha fetoprotein for early detection of hepatocellular carcinoma in patients with cirrhosis: a meta-analysis. *Gastroenterology*. 2018;154(6):1706–18. e1701.
7. Li X, Zhou Y, Jiang J. A retrospective study of long-term clinical outcomes in patients with Chronic Hepatitis C Treated with Interferon and Ribavirin. *Discov Med*. 2023;35(178):868–76.
8. Lv J, Xu Y, Xu L, Nie L. Quantitative functional evaluation of liver fibrosis in mice with dynamic contrast-enhanced photoacoustic imaging. *Radiology*. 2021;300(1):89–97.
9. Collier J, Sherman M. Screening for hepatocellular carcinoma. *Hepatology*. 1998;27(1):273–8.
10. Wang W, Wei C. Advances in the early diagnosis of hepatocellular carcinoma. *Genes Dis*. 2020;7(3):308–19.
11. Grazioli L, Ambrosini R, Frittoli B, Grazioli M, Morone M. Primary benign liver lesions. *Eur J Radiol*. 2017;95:378–98.
12. Su T-H, Wu C-H, Liu T-H, Ho C-M, Liu C-J. Clinical practice guidelines and real-life practice in hepatocellular carcinoma: a Taiwan perspective. *Clin Mol Hepatol*. 2023;29(2):230.
13. Cohen D, Kesler M, Muchnik Kurash M, Even-Sapir E, Levine C. A lesson in humility: the added values of PET-MRI over PET-CT in detecting malignant hepatic lesions. *Eur J Nucl Med Mol Imaging*. 2023;1–11.
14. LeGout JD, Bolan CW, Bowman AW, et al. Focal nodular Hyperplasia and focal nodular hyperplasia-like lesions. *Radiographics*. 2022;42(4):1043–61.
15. Klenk C, Gawande R, Uslu L, et al. Ionising radiation-free whole-body MRI versus (18)F-fluorodeoxyglucose PET/CT scans for children and young adults with cancer: a prospective, non-randomised, single-centre study. *Lancet Oncol*. 2014;15(3):275–85.
16. Acs B, Rantalainen M, Hartman J. Artificial intelligence as the next step towards precision pathology. *J Intern Med*. 2020;288(1):62–81.
17. Lew M, Hissong EM, Westerhoff MA, Lamps LW. Optimizing small liver biopsy specimens: a combined cytopathology and surgical pathology perspective. *J Am Soc Cytopathol*. 2020;9(5):405–21.
18. Liu JT, Glaser AK, Bera K, et al. Harnessing non-destructive 3D pathology. *Nat Biomedical Eng*. 2021;5(3):203–18.
19. Chan LK, Tsui YM, Ho DW, Ng IO. Cellular heterogeneity and plasticity in liver cancer. *Semin Cancer Biol*. 2022;82:134–49.
20. Fan J, Slowikowski K, Zhang F. Single-cell transcriptomics in cancer: computational challenges and opportunities. *Exp Mol Med*. 2020;52(9):1452–65.
21. Bao S, Li K, Yan C, Zhang Z, Qu J, Zhou M. Deep learning-based advances and applications for single-cell RNA-sequencing data analysis. *Brief Bioinform*. 2022;23(1):bbab473.
22. He B, Zhang Y, Zhou Z, et al. A neural network framework for predicting the tissue-of-origin of 15 common cancer types based on RNA-Seq data. *Front Bioeng Biotechnol*. 2020;8:737.
23. He B, Dai C, Lang J, et al. A machine learning framework to trace tumor tissue-of-origin of 13 types of cancer based on DNA somatic mutation. *Biochim et Biophys Acta (BBA)-Molecular Basis Disease*. 2020;1866(11):165916.
24. He B, Sun H, Bao M, et al. A cross-cohort computational framework to trace tumor tissue-of-origin based on RNA sequencing. *Sci Rep*. 2023;13(1):15356.
25. Mellman I, Chen DS, Powles T, Turley SJ. The cancer-immunity cycle: indication, genotype, and immunotype. *Immunity*. 2023;56(10):2188–205.
26. Monti F, Frasca F, Eynard D, Mannoni D, Bronstein MM. Fake news detection on social media using geometric deep learning. *arXiv preprint arXiv:190206673*. 2019.
27. Kipf TN, Welling M. Semi-Supervised Classification with Graph Convolutional Networks. Paper presented at: International Conference on Learning Representations.
28. Zhang M, Yang H, Wan L, et al. Single-cell transcriptomic architecture and intercellular crosstalk of human intrahepatic cholangiocarcinoma. *J Hepatol*. 2020;73(5):1118–30.
29. Xue R, Zhang Q, Cao Q, et al. Liver tumour immune microenvironment subtypes and neutrophil heterogeneity. *Nature*. 2022;612(7938):141–7.
30. Lin C-I, Merley A, Sciuto TE, et al. TM4SF1: a new vascular therapeutic target in cancer. *Angiogenesis*. 2014;17:897–907.
31. Ma L, Hernandez MO, Zhao Y, et al. Tumor cell biodiversity drives microenvironmental reprogramming in liver cancer. *Cancer Cell*. 2019;36(4):418–30. e416.
32. Shubinsky G, Schlesinger M. The CD38 lymphocyte differentiation marker: new insight into its ectoenzymatic activity and its role as a signal transducer. *Immunity*. 1997;7(3):315–24.
33. Sinha D, Kumar A, Kumar H, Bandyopadhyay S, Sengupta D. dropClust: efficient clustering of ultra-large scRNA-seq data. *Nucleic Acids Res*. 2018;46(6):e36–36.
34. Sun Y, Wu L, Zhong Y, et al. Single-cell landscape of the ecosystem in early-relapse hepatocellular carcinoma. *Cell*. 2021;184(2):404–21. e416.
35. Zhang Q, He Y, Luo N, et al. Landscape and dynamics of single immune cells in hepatocellular carcinoma. *Cell*. 2019;179(4):829–45. e820.
36. Ye Z, Cheng P, Huang Q, Hu J, Huang L, Hu G. Immunocytes interact directly with cancer cells in the tumor microenvironment: one coin with two sides and future perspectives. *Front Immunol*. 2024;15:1388176.
37. Huang H, Wu N, Liang Y, Peng X, Shu J. SLNL: a novel method for gene selection and phenotype classification. *Int J Intell Syst*. 2022;37(9):6283–304.
38. Chen S, Zeng J, Huang L, et al. RNA adenosine modifications related to prognosis and immune infiltration in osteosarcoma. *J Translational Med*. 2022;20(1):228.
39. Bronstein MM, Bruna J, LeCun Y, Szlam A, Vandergheynst P. Geometric deep learning: going beyond euclidean data. *IEEE Signal Process Mag*. 2017;34(4):18–42.
40. Zafeiriou S, Bronstein M, Cohen T, et al. Guest Editorial: Non-euclidean Machine Learning. *IEEE Trans Pattern Anal Mach Intell*. 2022;44(2):723–6.
41. Liu S, Yu Y, Zhang M, Wang W, Cao X. The involvement of TNF- $\alpha$ -related apoptosis-inducing ligand in the enhanced cytotoxicity of IFN- $\beta$ -stimulated human dendritic cells to tumor cells. *J Immunol*. 2001;166(9):5407–15.
42. Tong Z, Liang Y, Sun C, Rosenblum DS, Lim A. Directed graph convolutional network. *arXiv Preprint arXiv:200413970*. 2020.
43. Gong S, Liang X, Zhang M, et al. Tumor Microenvironment-activated hydrogel platform with programmed release property evokes a Cascade-Amplified Immune response against Tumor Growth, Metastasis and Recurrence. *Small*. 2022;18(50):2107061.
44. Wu H, Fu X, Zhai Y, Gao S, Yang X, Zhai G. Development of effective tumor vaccine strategies based on immune response cascade reactions. *Adv Healthc Mater*. 2021;10(13):2100299.
45. Xu L, Lu Y, Deng Z, et al. Single-cell landscape of immunocytes in patients with extrahepatic cholangiocarcinoma. *J Translational Med*. 2022;20(1):210.
46. Chen Q-F, Li W, Wu P-H, Shen L-J, Huang Z-L. Significance of tumor-infiltrating immunocytes for predicting prognosis of hepatitis B virus-related hepatocellular carcinoma. *World J Gastroenterol*. 2019;25(35):5266.
47. Shi Y, Li Z, Zheng W, et al. Changes of immunocytic phenotypes and functions from human colorectal adenomatous stage to cancerous stage: update. *Immunobiology*. 2015;220(10):1186–96.
48. Xu X, Huang X, Sun J et al. 3D-stacked multistage inertial microfluidic chip for high-throughput enrichment of circulating tumor cells. *Cyborg Bionic Syst*. 2022.

## Publisher's note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.