

## ARTICLE OPEN



# Cell graph neural networks enable the precise prediction of patient survival in gastric cancer

Yanan Wang<sup>1,15</sup>, Yu Guang Wang<sup>2,3,4,15</sup>, Changyuan Hu<sup>1</sup>, Ming Li<sup>5</sup>, Yanan Fan<sup>4</sup>, Nina Otter<sup>6</sup>, Ikuan Sam<sup>7</sup>, Hongquan Gou<sup>7</sup>, Yiqun Hu<sup>7</sup>, Terry Kwok<sup>1,8</sup>, John Zalcberg<sup>9,10</sup>, Alex Boussioutas<sup>11</sup>, Roger J. Daly<sup>12</sup>, Guido Montúfar<sup>3,6</sup>, Pietro Liò<sup>12</sup>, Dakang Xu<sup>7</sup>, Geoffrey I. Webb<sup>13,14</sup> and Jiangning Song<sup>1,14</sup>

Gastric cancer is one of the deadliest cancers worldwide. An accurate prognosis is essential for effective clinical assessment and treatment. Spatial patterns in the tumor microenvironment (TME) are conceptually indicative of the staging and progression of gastric cancer patients. Using spatial patterns of the TME by integrating and transforming the multiplexed immunohistochemistry (mIHC) images as Cell-Graphs, we propose a graph neural network-based approach, termed *Cell-Graph Signature* or  $CG_{Signature}$ , powered by artificial intelligence, for the digital staging of TME and precise prediction of patient survival in gastric cancer. In this study, patient survival prediction is formulated as either a binary (*short-term* and *long-term*) or ternary (*short-term*, *medium-term*, and *long-term*) classification task. Extensive benchmarking experiments demonstrate that the  $CG_{Signature}$  achieves outstanding model performance, with Area Under the Receiver Operating Characteristic curve of  $0.960 \pm 0.01$ , and  $0.771 \pm 0.024$  to  $0.904 \pm 0.012$  for the binary- and ternary-classification, respectively. Moreover, Kaplan–Meier survival analysis indicates that the “digital grade” cancer staging produced by  $CG_{Signature}$  provides a remarkable capability in discriminating both binary and ternary classes with statistical significance ( $P$  value  $< 0.0001$ ), significantly outperforming the AJCC 8th edition Tumor Node Metastasis staging system. Using Cell-Graphs extracted from mIHC images,  $CG_{Signature}$  improves the assessment of the link between the TME spatial patterns and patient prognosis. Our study suggests the feasibility and benefits of such an artificial intelligence-powered digital staging system in diagnostic pathology and precision oncology.

npj Precision Oncology (2022)6:45; <https://doi.org/10.1038/s41698-022-00285-5>

## INTRODUCTION

Gastric cancer (GC) accounted for 768,793 deaths in 2020, representing the fourth deadliest cancer globally<sup>1</sup>. The 5-year survival rate of GC is around 20%<sup>2</sup>. A more accurate prognosis can greatly assist clinical decision-making, especially regarding which patients would benefit from aggressive treatment. The Tumor-Node-Metastasis (TNM) staging system<sup>3</sup> is the most prevalent cancer staging system primarily used in hospitals and medical centers worldwide, which reflects information on the primary tumor, affected lymph nodes, and metastasis. Many current treatment recommendations and guidelines are based on the TNM stages. However, significant differences in clinical outcomes have been observed in GC patients with the same TNM stage and similar treatment regimens<sup>4–6</sup>. These findings indicate the TNM staging system has limitations and accordingly, cannot be used to accurately predict the prognosis of cancer patients. As such, new strategies that can provide more tailored staging information and improve prognosis predictions are highly desirable.

Recent years have seen numerous data-driven, machine learning-based studies of cancer prognosis. For instance, Yu et al.<sup>7</sup> introduced prognosis prediction of lung adenocarcinoma and squamous cell carcinoma of stage I, and their model can distinguish the shorter-term survivors from longer-term survivors ( $p = 0.003$  and  $p = 0.023$ ). Mobadersany et al.<sup>8</sup> presented a survival convolutional neural network (SCNN), and their developed histology image-based SCNN reached comparable performance on astrocytomas of grades III and IV with histology grading or molecular subtyping. In another study, Jiang et al.<sup>9</sup> proposed the GC-SVM classifier as a powerful survival predictor using the data of immunomarkers and could predict the adjuvant chemotherapy benefit of gastric cancer patients with stages II and III. Wulczyn et al.<sup>10</sup> conducted a survival prediction study involving multiple cancers based on deep learning, and as a result, their model was capable of making significant survival predictions for five out of ten cancers and could effectively stratify cancer patients of stages II and III. Jiang et al.<sup>11</sup> developed a convolutional neural network-based classifier from H&E images to predict the prognosis of stage III colon cancer patients. Dimitriou et al.<sup>12</sup> introduced a K-nearest

<sup>1</sup>Biomedicine Discovery Institute and Department of Biochemistry and Molecular Biology, Monash University, Melbourne, VIC 3800, Australia. <sup>2</sup>Institute of Natural Sciences, School of Mathematical Sciences, Key Laboratory of Scientific and Engineering Computing of Ministry of Education (MOE-LSC), and Center for Mathematics of Artificial Intelligence Institute, Shanghai Jiao Tong University, Shanghai 200240, China. <sup>3</sup>Max Planck Institute for Mathematics in Sciences, Leipzig 04103, Germany. <sup>4</sup>School of Mathematics and Statistics, The University of New South Wales, Sydney, NSW 2052, Australia. <sup>5</sup>Key Laboratory of Intelligent Education Technology and Application of Zhejiang Province, Zhejiang Normal University, Jinhua 321004, China. <sup>6</sup>Department of Mathematics, Department of Statistics, University of California, Los Angeles, CA 90095, USA. <sup>7</sup>Faculty of Medical Laboratory Science, Ruijin Hospital, School of Medicine, Shanghai Jiao Tong University, Shanghai 200025, China. <sup>8</sup>Biomedicine Discovery Institute and Department of Microbiology, Monash University, Melbourne, VIC 3800, Australia. <sup>9</sup>School of Public Health and Preventive Medicine, Monash University, Melbourne, VIC 3004, Australia. <sup>10</sup>Department of Medical Oncology, The Alfred Hospital, Melbourne, VIC 3004, Australia. <sup>11</sup>The Alfred Hospital and Central Clinical School, Monash University, Melbourne, VIC 3004, Australia. <sup>12</sup>Cambridge Centre for AI in Medicine, Department of Computer Science and Technology, University of Cambridge, Cambridge CB3 0FD, United Kingdom. <sup>13</sup>Department of Data Science and Artificial Intelligence, Monash University, Melbourne, VIC 3800, Australia. <sup>14</sup>Monash Data Futures Institute, Monash University, Melbourne, VIC 3800, Australia. <sup>15</sup>These authors contributed equally: Yanan Wang, Yu Guang Wang. ✉email: pl219@cam.ac.uk; dakang\_xu@163.com; Geoff.Webb@monash.edu; Jiangning.Song@monash.edu

neighbor-based method to predict the mortality of stage II colorectal cancer patients using immunofluorescence images. Although these prognosis prediction studies achieved promising performance using H&E staining histology or immunohistochemistry staining images, they were often restricted to specific subtypes or stages of the corresponding cancers. Moreover, these studies did not consider any spatial information from the tumor microenvironment (TME).

Cell distribution in the TME is not random but is rather associated with the underlying functional state<sup>13–15</sup>. Therefore, the exploration of the TME of cancer samples would offer critical insights into the key spatial patterns associated with the growth, cancer progression, and thus patient prognosis<sup>16</sup>. The recent advent of the multiplexed immunohistochemistry (mIHC) staining technique enables systematic investigation of the TME<sup>17,18</sup> and supports extraction of enriched spatial information from the TME, including the cell location, cell types, cell and nucleus morphological information, and related optical information<sup>16,19</sup>. Researchers have applied the mIHC technique to analyze the TME of pancreatic cancer and found that spatial distribution of cytotoxic T cells in proximity to cancer cells correlates with increased overall patient survival<sup>16</sup>. Barua et al.<sup>19</sup> applied a statistical scoring based method, G-cross function, to measure the patterns of two different cell types, such as T-reg and CD8, and found that high infiltration of T-reg in the core tumor area is an independent predictor of worse overall survival (OS) in patients of non-small cell lung cancer. However, these studies only considered the spatial features of limited cell types and only used handcrafted features. Therefore, comprehensive and quantitative methods that assess the relationships between spatial features descriptive of cell distribution and prognosis are currently lacking.

Inspired by the concept of the Cell-Graph<sup>13,20</sup> and the success of graph neural networks (GNN)<sup>21–23</sup>, especially their applications to the analysis of biology data<sup>24,25</sup>, we hypothesize that intricate spatial distribution information of the TME is informative for the prediction of the OS of GC patients and a GNN model can effectively capitalize on useful patterns generated by Cell-Graphs. To validate this hypothesis, we have developed a novel GNN-based approach for predicting the prognosis of GC patients using Cell-Graph data, which we call the *Cell-Graph Signature* or  $CG_{Signature}$ . The overall workflow is illustrated in Figure 1 and Supplementary Fig. 1. In this study, we formulate prognosis prediction as a classification problem by predicting the patient's survival time interval rather than a continuous time frame or a risk score and develop a workflow to perform the following threefold tasks. Firstly, it extracts comprehensive spatial and morphological information from mIHC images. Secondly, it further uses the extracted spatial information to stratify patients into either binary (*short-term* and *long-term*) or ternary (*short-term*, *medium-term*, and *long-term*) classes. Finally, it conducts the Kaplan–Meier survival analysis to verify the clinical significance of the  $CG_{Signature}$ .

$CG_{Signature}$  represents a powerful survival predictor under comprehensive and extensive benchmarking tests of gastric cancer across all subtypes and stages. Specifically,  $CG_{Signature}$  can effectively stratify short-term, medium-term, and long-term GC survivors at the early diagnosis stage, and achieved area under the receiver-operating characteristic curve (AUROC) values of  $0.960 \pm 0.01$  in terms of binary classification, and  $0.771 \pm 0.024$  to  $0.904 \pm 0.012$  in terms of ternary classification, respectively. In the follow-up survival analysis,  $CG_{Signature}$  outperformed the AJCC 8th edition TNM staging system on the testing cohort in terms of Harrell's Concordance Index<sup>26</sup>, Hazard Ratio (HR), and *p* value.

## RESULTS

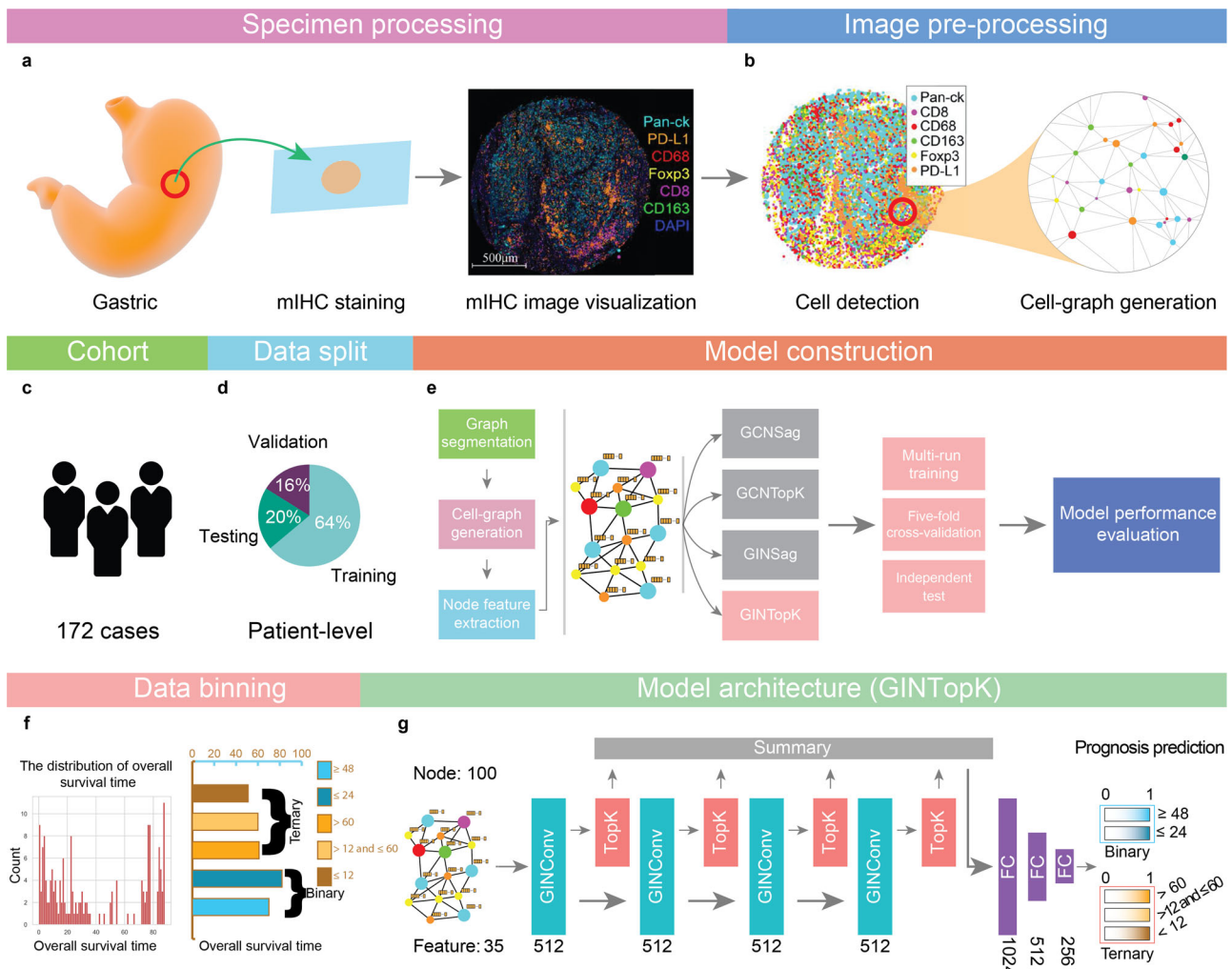
### Clinical characteristics and data-binning of the patient cohort

We collected the gastric cancer patient data between 2008 and 2010 from Shanghai Ruijin Hospital, affiliated with the School of Medicine, Shanghai Jiao Tong University. After removing patients labeled as “lost to follow-up”, we ended up with 172 patients in the cohort for further analysis. All the patients were diagnosed with gastric adenocarcinoma; no neoadjuvant therapies were applied prior to curative gastrectomy. All samples were extracted from surgical samples and were formalin-fixed and paraffin-embedded. The clinical characteristics of this cohort are illustrated in Supplementary Table 1. This cohort contains 124 males, 47 females, and one case without gender information. With respect to the survival status, 113 cases were recorded as “deceased” while 59 patients as “alive”. The OS time of the cohort ranges from 0 to 88 months with a median follow-up time of 27 months. Here, the survival time of 0 months indicates the death occurred before the first discharge of the patient from the hospital. The statistical summary of their TNM (the AJCC 8th edition) stages is provided in Supplementary Table 1. In particular, the patient numbers of the TNM stages of I, II, III, and IV are 14, 52, 95, and 3, respectively. Two data-binning strategies were applied to segment the patient OS into binary- or ternary-class datasets. More specifically, patients with OS time shorter than 24 months and longer than 48 months were categorized as short-term, and long-term in the binary-class dataset. The patients whose OS time is between 24 months and 48 months were removed from the training dataset but used in subsequent survival analyses. More details can be found in the section on “Survival analysis and performance comparison with the TNM staging system”. In the ternary-class dataset, patients were classified into short-term, medium-term, and long-term classes, using the thresholds of 12 and 60 months. Here, the thresholds for binary and ternary class data-binning were chosen by considering the relative class balance and the commonly used follow-up periods of 12, 36, and 60 months. To extend the classification margin, we used 24 and 48 months as the thresholds rather than 36 months to divide patients into short-term and long-term classes in binary data binning. We did not optimize the data binning threshold, which can be conducted when more data become available. Model training and subsequent analysis were performed using these two datasets.

### Workflow overview

Figure 1 illustrates an overall workflow and the model architecture of the proposed  $CG_{Signature}$  approach. As shown in Fig. 1a, the mIHC technique was used to stain the GC tissue samples. Specifically, the nuclear counterstain, DAPI, was used for cell nuclei staining, and six antibodies of Pan-CK, CD8, CD68, CD163, Foxp3, and PD-L1 were used as annotation indicators for six different types of cells. After digitalization, cell locations, types, and related optical and morphological features were extracted using the digital pathology software. After this procedure, we obtained the CSV files in which each row corresponds to each cell with the node features shown in Table 1. Based on these CSV files as the input, we developed a workflow (details can be seen in Algorithm 1) to process the raw data and build the GNN-based model to predict the patient OS interval using the features extracted from mIHC images.

The key steps of the workflow are as follows: (1) Image pre-processing: Sub-sampling and Cell-Graph generation were performed at this step. Specifically, each mIHC image was firstly segmented into multiple non-overlapping regions with no more than 100 cells. For each region, we built a graph where each cell was represented as a node and the reciprocal of the Euclidean distance of each cell-cell pair was used to establish edges between them with a distance of fewer than 20  $\mu\text{m}$ . Detailed information can be found in the section “Cell-Graph construction” of Methods.



**Fig. 1** An overall workflow of graph neural network-based prognosis prediction using Cell-Graphs. **a** Specimen processing: the tumor tissues were extracted from gastric cancer, and stained with seven different biomarkers including DAPI, Pan-CK, CD8, CD68, CD163, Foxp3, and PD-L1. **b** Image pre-processing: sub-sampling and cell-graph construction were conducted for image pre-processing. **c** An illustration for the cohort, 172 gastric cancer patients were collected. **d** Data split. The training, validation, and testing datasets were split with the percentages of 64%, 16%, and 20%, respectively. **e** Model construction: four different GNN model architectures, including GCNSag, GCNTopK, GINSag, and GINTopK, were constructed and compared. Multi-run model training, fivefold cross-validation, and independent tests were conducted to evaluate the performance of the constructed GNN models. **f** Data binning: overall survival time ranged from 0 to 88 months, and two data-binning strategies were applied to generate binary- and ternary-class datasets. **g** Model architecture: The four models shared the same architecture but employed different types of convolutional unit and pooling layer, which consists of four consecutive convolutional layer and pooling layer blocks, followed by a summary layer and three fully connected layers, prior to the generation of the final classification outcome. The architecture of the best-performing GINTopK model is illustrated herein, which outperformed the other three model architectures and also achieved the best performance on the test dataset. The corresponding number of hidden layers or feature dimensions is indicated at the bottom of each box. Here, FC stands for “fully connected layer”.

Then, we extracted a total of 35 features (as shown in Table 1) for each cell as the node attributes, including five optical features for each biomarker and five morphological features for each cell. Such generated cell-based graph is referred to as Cell-Graph<sup>13,20</sup>. There are ~90 Cell-Graphs constructed for each mIHC image (for each patient). Cell-Graphs originated from the same mIHC image share the same label with the corresponding patient. (2) Data split: After Cell-Graph construction, the whole dataset was partitioned into the training, validation, and test sets with the ratio of 0.64:0.16:0.20 at the patient level. In addition, we also generated the files for performing fivefold cross-validation by generating five non-overlapping training-validation subsets and evaluating the model performance on these fivefold subsets. (3) Hyperparameter optimization: We utilized the Hyperopt toolkit<sup>27</sup> from the Ray software package<sup>28</sup> to tune the hyperparameters of GNN models. The optimized hyperparameters were then used for the follow-up

model training and performance evaluation. (4) Model performance evaluation and data visualization: To comprehensively assess the capability and reliability of our GNN model, we evaluate model performance using multi-run model training, fivefold cross-validation, and independent tests. The test results were visualized by generating the receiver-operating characteristic (ROC) curves, confusion matrix, and boxplots of Accuracy, F1-Score, and Matthews Correlation Coefficient (MCC). Performance metrics are defined in Section “Metrics of model performance evaluation” in the Supplementary material.

#### Performance benchmarking of different GNN models for prognosis prediction

We constructed four different types of GNN models and examined their performance in predicting the OS of gastric cancer patients,



**Table 1.** The list of node attributes and their variable types.

Feature name	Feature type
DAPI positive	Boolean
DAPI positive nucleus	Boolean
DAPI positive cytoplasm	Boolean
DAPI nucleus Intensity	Float
DAPI cytoplasm intensity	Float
PD-L1 (Opal 520) positive	Boolean
PD-L1 (Opal 520) positive nucleus	Boolean
PD-L1 (Opal 520) positive cytoplasm	Boolean
PD-L1 (Opal 520) nucleus intensity	Float
PD-L1 (Opal 520) cytoplasm intensity	Float
CD68 (Opal 540) positive	Boolean
CD68 (Opal 540) positive nucleus	Boolean
CD68 (Opal 540) positive cytoplasm	Boolean
CD68 (Opal 540) nucleus intensity	Float
CD68 (Opal 540) cytoplasm intensity	Float
Foxp3 (Opal 570) positive	Boolean
Foxp3 (Opal 570) positive nucleus	Boolean
Foxp3 (Opal 570) positive cytoplasm	Boolean
Foxp3 (Opal 570) nucleus intensity	Float
Foxp3 (Opal 570) cytoplasm intensity	Float
CD8 (Opal 620) positive	Boolean
CD8 (Opal 620) positive nucleus	Boolean
CD8 (Opal 620) positive cytoplasm	Boolean
CD8 (Opal 620) nucleus intensity	Float
CD8 (Opal 620) cytoplasm intensity	Float
Pan-CK (Opal 690) positive	Boolean
Pan-CK (Opal 690) positive nucleus	Boolean
Pan-CK (Opal 690) positive cytoplasm	Boolean
Pan-CK (Opal 690) nucleus intensity	Float
Pan-CK (Opal 690) cytoplasm intensity	Float
Cell area ( $\mu\text{m}^2$ )	Float
Cytoplasm area ( $\mu\text{m}^2$ )	Float
Nucleus area ( $\mu\text{m}^2$ )	Float
Nucleus perimeter ( $\mu\text{m}$ )	Float
Nucleus roundness	Float

Each type of feature is comprised of three Boolean variables and two float variables. These Boolean variables were identified by the pathology software based on the float values of Nucleus Intensity and Cytoplasm Intensity of each biomarker. Moreover, five different morphology features were extracted as the node attributes.

including GINTopK, GINSAG, GCNTopK, and GCNSAG. Here, GIN<sup>23</sup>, and GCN<sup>29</sup> are two graph convolution computational units (differences can be seen in Supplementary Fig. 2), whereas TopKPooling<sup>22,30,31</sup> and SAGPooling<sup>31,32</sup> are two graph pooling computational units. The graph convolutional and pooling layers are the core components of the GNN architecture. Fivefold cross-validation was conducted to assess the model of each GNN model on both binary- and ternary-classification tasks. The results are averaged on ten repetitions of fivefold cross-validation for GINTopK on binary classification (as shown in Supplementary Fig. 3) to circumvent the randomness of the model during training. In this procedure, Accuracy, F1-Score, MCC, and AUROC were calculated to evaluate the performance. Figure 2a illustrates the performance results of binary classification on fivefold cross-validation. As we observe, the median values of both Accuracy

and F1-score for the four GNNs ranged from 0.83 to 0.92, while the median values of MCC ranged from 0.66 to 0.84, respectively. Figure 2b shows the performance results of ternary classification on fivefold cross-validation. We can see that the ternary-class classification models achieved the median values of Accuracy ranging from 0.76 to 0.82, F1-score from 0.64 to 0.72, and MCC from 0.46 to 0.5, respectively. According to the results shown in Fig. 2a, b, GINTopK slightly outperformed the other three GNN models on both binary- and ternary classifications. Therefore, GINTopK was selected as the best-performing GNN model and employed for subsequent performance benchmarking and survival analysis.

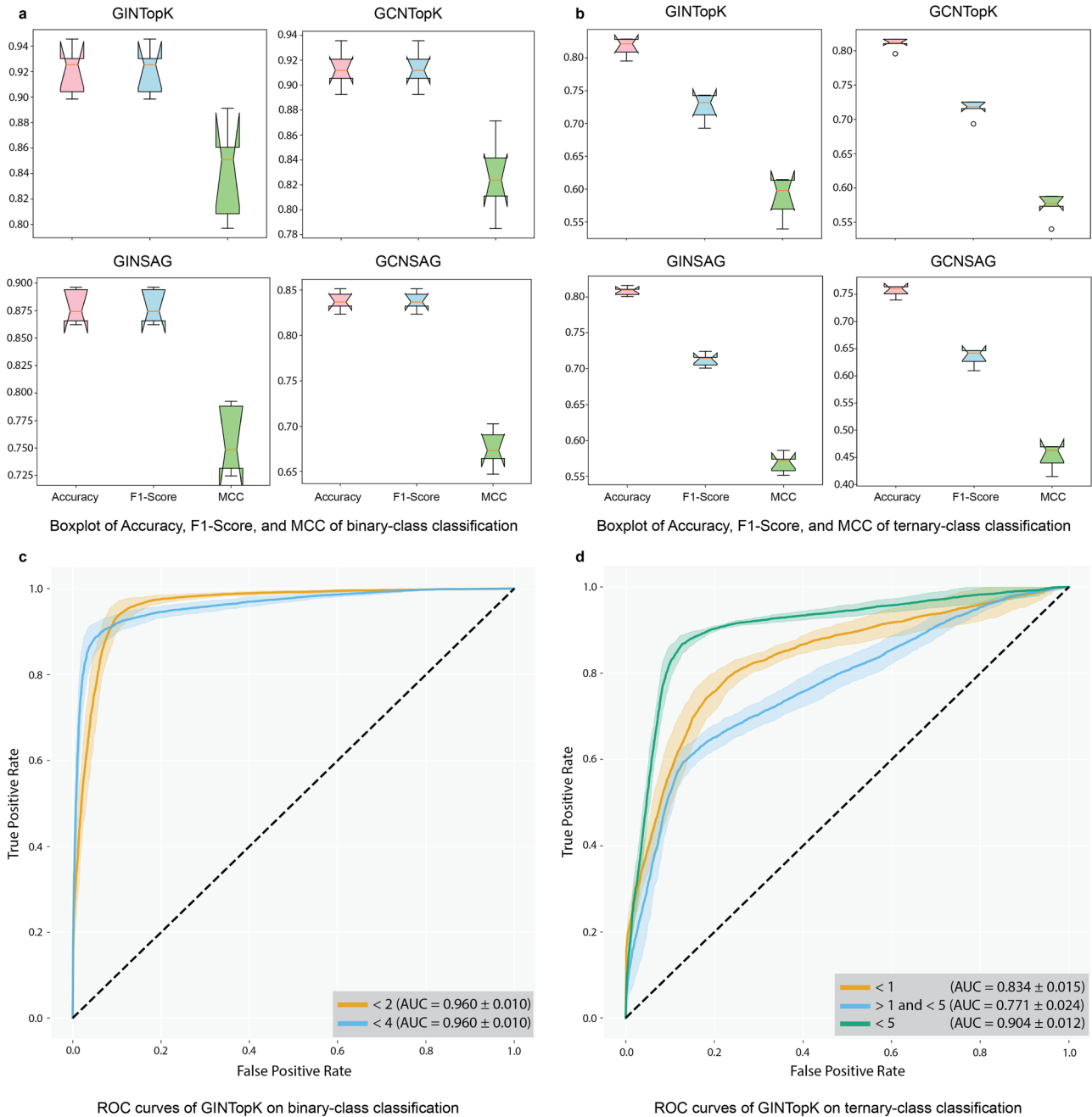
ROC curves of GINTopK on the binary- and ternary-classification tasks are illustrated in Fig. 2c, d, respectively. The binary-class GINTopK model achieved the AUROC value of  $0.96 \pm 0.01$  on fivefold cross-validation. In contrast, the ternary-class GINTopK classifier reached the AUROC values of  $0.834 \pm 0.015$ ,  $0.771 \pm 0.024$ , and  $0.904 \pm 0.012$  for *short-term* (<12 months), *medium-term* (>12 and <60 months), and *long-term* (>60 months) on fivefold cross-validation, respectively (Fig. 2d). Moreover, the performance results of the binary class GINTopK model on ten repetitions of fivefold cross-validation are displayed in Supplementary Fig. 3. We can see that the median values of both Accuracy and F1-score were within the range of 0.90-0.93 (MCC values ranged from 0.80 to 0.86), thereby suggesting the stability of our proposed GINTopK model.

In Fig. 3, the performance results of the GINTopK model on the independent test are visualized using ROC curves and a confusion matrix. It can be seen that the model achieved similar performance to that on fivefold cross-validation in terms of AUROC values on both binary- and ternary-classification tasks. In terms of the confusion matrix, 96% and 89% of the short-term and long-term patients could be accurately predicted using the binary-classification model. The true positive percentages of the ternary-class model were 81%, 59%, and 85%, corresponding to the short-term, medium-term, and long-term classes (Fig. 3).

Taken together, the outstanding performance of the GINTopK model on both cross-validation and independent test indicate that our proposed GNN approach is capable of effectively capturing the underlying prognostic patterns from the well-constructed Cell-Graphs. The captured prognostic patterns by GNN model are characteristic of the spatial information of cell locations and types of the TME, which incorporates more potentially informative features than the TNM staging system.

### Ablation studies and prognostic value of different types of cell features

To examine the effect of node features of different cell types on model performance, we further performed ablation studies to assess the contribution of features to the binary- and ternary-classification performance by removing each type of feature in an iterative manner. Thirty-five node features of seven types were used in this study, including DAPI, Pan-CK, CD8, CD68, Foxp3, PD-L1, and morphological features. We first evaluated the performance of the GNN model trained using all these features, and then, evaluated the performance of the models trained using the remaining features after removing each type of feature from the all-feature set in turn. For each iteration, we trained the models five times with random initialization of the weights using the same dataset and calculated the mean and standard deviation of Accuracy. The results are shown in Table 2, where the feature contribution was measured by the accuracy change compared with that of the all-feature model. Note that when a type of feature is removed, and accuracy increase means that including the feature type reduced accuracy, and an accuracy decrease means that the feature type played an important role in attaining the all-feature accuracy.



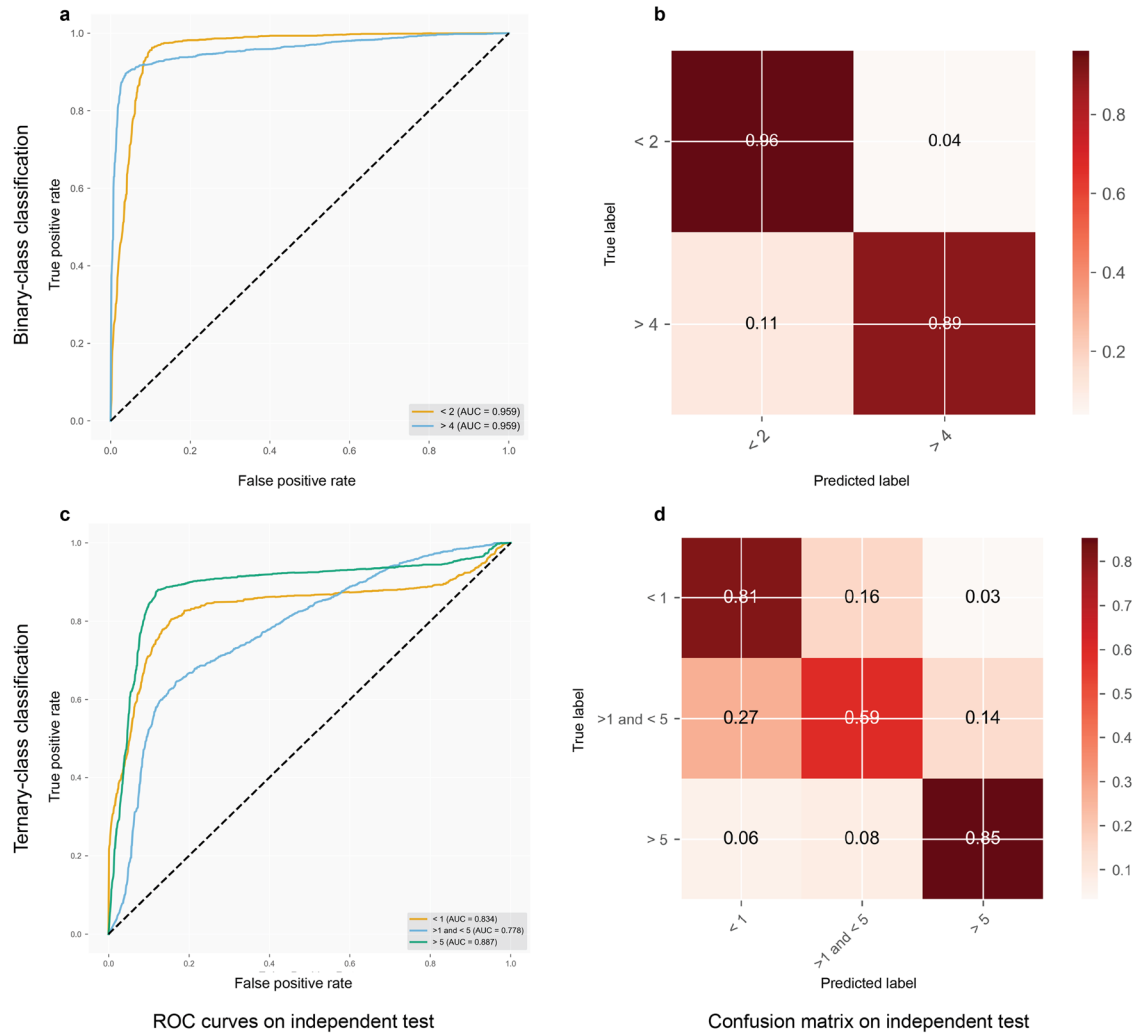
**Fig. 2 Model performance of four GNNs on fivefold cross-validation.** **a, b** show the Boxplots of performance metrics of Accuracy, F1-score, and MCC on fivefold cross-validation. **c, d** illustrate the ROCs of GINTopK binary- and ternary models on fivefold cross-validation. In the boxplot, the center line marks the mid-point of the data; the top and bottom lines show the maximum and minimum non-outlier data; the upper and lower bounds of the box indicate the third quartile and first quartile of the data; the height of the notch indicates the 95% confidence interval of the median point; small circles represent outliers.

According to Table 2, the variant models trained using these feature subsets and all-feature sets achieved comparable Accuracy values in both binary and ternary classifications. In the binary classification, the DAPI features and morphological features made more important contributions to the model performance compared with other types of features (e.g., the Accuracy dropped by 0.035 and 0.025, respectively), which reflects the nucleus differences in optical and morphology of the TME. Thus, the inclusion of these two types of features helped to better distinguish the long-term from short-term patients. In the case of ternary classification, we can see that the GNN models trained without the DAPI and morphology features achieved the lowest

Accuracy, which is consistent with the observation in the binary classification. In summary, our proposed GNN model can effectively learn distinguishable spatial features from the TME and further enhance the performance of prognosis stratification by combining the DAPI and morphology features.

**Survival analysis and performance comparison with the TNM staging system**

To further investigate the prognostic values and clinical importance of the predictions produced by  $CG_{Signature}$ , we conducted the Kaplan–Meier survival analysis using the patient-level results



**Fig. 3 Performance assessment of the GINTopK model in terms of ROC curves and confusion matrix on the independent test.** The left column shows the ROC curves of **a** binary- and **c** ternary classification, while the right column displays the confusion matrix of the model predictions on the **b** binary- and **d** ternary classification tasks.

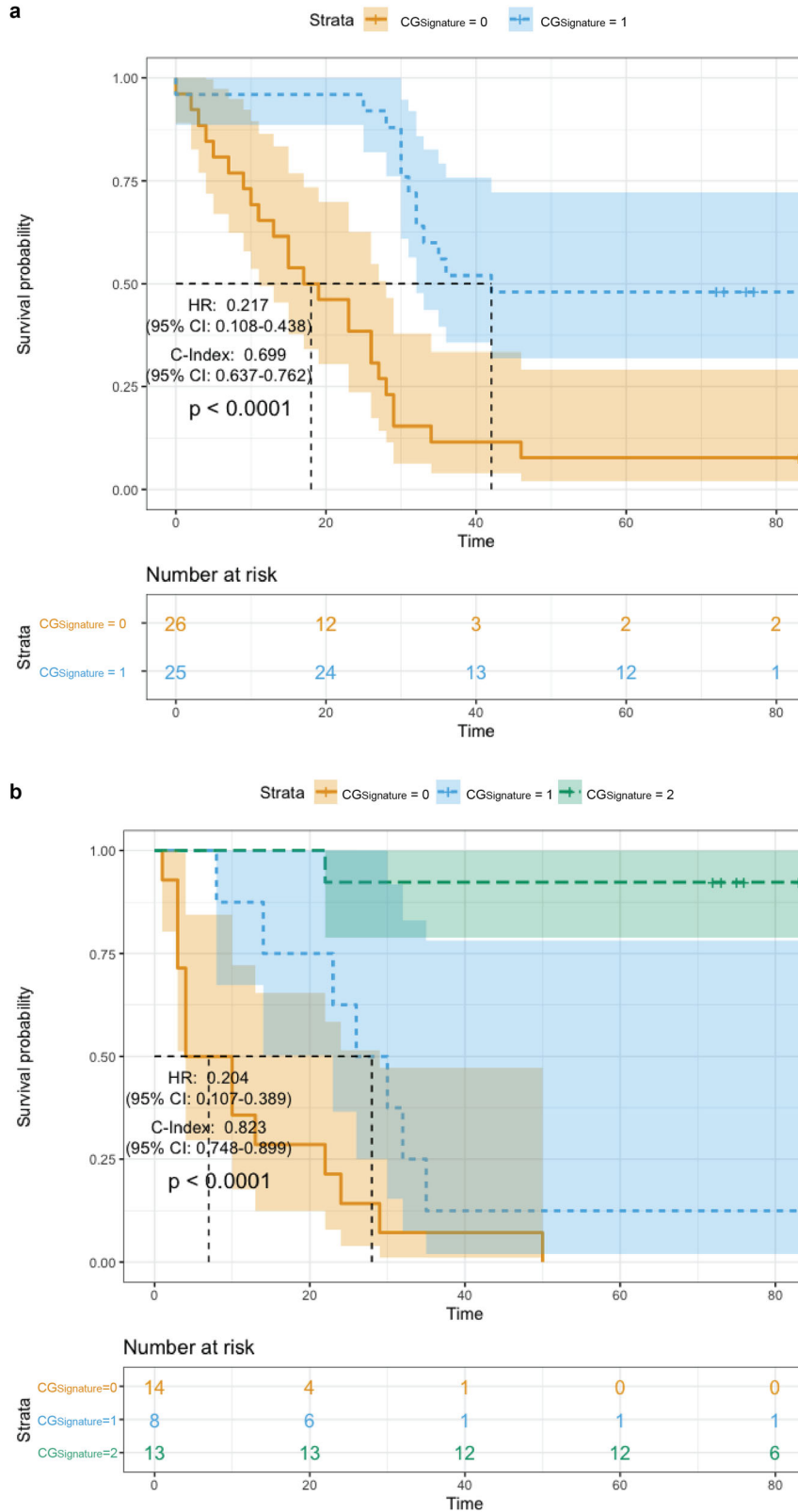
Feature sets	ACC of binary	ACC of ternary
All-features	0.917 ± 0.012	0.719 ± 0.020
No-DAPI	0.882 ± 0.022	0.710 ± 0.018
No-PD-L1	0.921 ± 0.014	0.718 ± 0.006
No-CD68	0.911 ± 0.011	0.717 ± 0.010
No-FOXP3	0.918 ± 0.017	0.719 ± 0.016
No-CD8	0.910 ± 0.006	0.732 ± 0.023
No-Pan-CK	0.927 ± 0.023	0.714 ± 0.016
No-morphology	0.892 ± 0.009	0.706 ± 0.021

The relative importance and contribution of the features were measured by the accuracy change compared with that of the all-feature model.

of both binary- and ternary classifications. For each patient, we first collected the predicted results of all the subsampled Cell-Graphs. Next, we calculated the class percentages of these predictions, and took the class with the maximum percentage as the final patient-level prediction of the corresponding patient.

Using these patient-level predicted results ('digital grade') of binary classification (with predicted class labels of  $CG_{Signature} = 0$  and  $CG_{Signature} = 1$ ) and ternary classification (with predicted class labels of  $CG_{Signature} = 0$ ,  $CG_{Signature} = 1$ , and  $CG_{Signature} = 2$ ), we conducted the survival analysis and plotted their Kaplan-Meier curves, shown in Fig. 4. More specifically, when using the binary-class predictions, the median survival time of patient test cohorts predicted as  $CG_{Signature} = 0$  and  $CG_{Signature} = 1$  were about 18 months and 42 months, respectively. The HR was 0.217 (95% CI: 0.108–0.438), the C-Index was 0.699 (95% CI: 0.637–0.762), and the  $p$  value was  $< 0.0001$ , indicating that  $CG_{Signature}$  has statistically significant prognostic power in separating the two groups of patient cohorts. When using the ternary-class predictions, the median survival time of patient cohorts predicted as  $CG_{Signature} = 0$  and  $CG_{Signature} = 1$  were  $\sim 7$  months and 28 months, respectively. The endpoint survival rate of  $CG_{Signature} = 2$  was  $\sim 92.3\%$  (Fig. 4b). The HR and C-Index were 0.204 (95% CI: 0.107–0.389) and 0.823 (95% CI: 0.748–0.899), respectively, with the  $p$  value  $< 0.0001$ .

We further conducted the univariate and multivariate Cox regression analyses based on the predictions of  $CG_{Signature}$  and the AJCC 8th edition TNM stages. The results are shown in Table 3. In the TNM staging system, there were eight groups of  $I_{A}$ ,  $I_{B}$ ,  $II_{A}$ ,  $II_{B}$ ,  $III_{A}$ ,  $III_{B}$ ,  $III_{C}$ ,  $IV_{A}$ , and  $IV_{B}$ . As no patients of stage  $IV$  were included in



**Fig. 4** Kaplan–Meier survival analysis of patient overall survival based on the “digital grade” (patient-level predictions) produced by **CG<sub>signature</sub>**. **a** Kaplan–Meier survival analysis results based on the binary-classification. **b** Kaplan–Meier survival analysis results based on the ternary classification.

**Table 3.** Univariate and multivariable Cox regression analysis of overall survival (Cox proportional hazards regression model) based on the predictions of binary- and ternary classification by  $CG_{Signature}$ .

Variable	Univariate analysis			Multivariate analysis			
	C-Index <sup>1</sup> (95% CI)	HR <sup>2</sup> (95% CI)	<i>p</i> value	C-Index <sup>1</sup> (95% CI)	HR <sup>2</sup> (95% CI)	<i>p</i> value	
Binary-class test cohort	pT	-	-	-	-	0.363	
	pN	-	-	-	-	0.378	
	TNM-2 <sup>3</sup>	0.659 (0.577–0.740)	5.276 (2.147–12.966)	<1e-4***			0.998
	TNM-3 <sup>4</sup>	-	-	-			0.998
	TNM-6 <sup>5</sup>	0.714 (0.623–0.805)	1.873 (1.388–2.529)	8.1e-4***			0.135
	$CG_{Signature}$ <sup>5</sup>	0.699 (0.637–0.762)	0.217 (0.108–0.438)	<1e-4***	0.804 (0.728–0.881)	0.037 (0.003–0.422)	0.007**
	$CG_{Signature}$ + TNM-2	0.740 (0.661–0.819)	2.412 (1.650–3.525)	<1e-4***			0.146
Ternary-class test cohort	pT	-	-	-	-	0.097	
	pN	-	-	-	-	0.053	
	TNM-2 <sup>3</sup>	-	-	-			0.998
	TNM-3 <sup>4</sup>	0.632 (0.510–0.753)	3.169 (1.335–7.522)	0.019*			0.998
	TNM-6 <sup>6</sup>	0.681 (0.535–0.827)	1.708 (1.212–2.407)	0.028*			0.066
	$CG_{Signature}$ <sup>7</sup>	0.823 (0.748–0.899)	0.204 (0.107–0.389)	<1e-4***	0.883 (0.820–0.947)	0.190 (0.087–0.414)	2.94e-05***

<sup>1</sup>C-Index: concordance index; <sup>2</sup>HR: Hazard ratio; <sup>3</sup>I+II vs. III; <sup>4</sup>I vs. II vs. III; <sup>5</sup>I, IIA, IIB, IIIA, IIIB, IIIC; <sup>6</sup>Low vs. high; <sup>7</sup>Low vs. medium vs. High.

The classification results were compared with Harrell's Concordance Index (C-Index), Hazard Ratio (HR), and *p* value. For the convenience of survival analysis comparison, the variables of TNM stages were regrouped into TNM-2 (I+II vs. III), TNM-3 (I vs. II vs. III), and TNM-6 (I, IIA, IIB, IIIA, IIIB, IIIC), while " $CG_{Signature}$ +TNM-2" denotes a four-class variable by combining the classes of TNM-2 and binary-class  $CG_{Signature}$ .

the binary-class test cohort and only one patient of stage IV was included in the ternary-class test cohort, we excluded the patients of stage IV and those without OS information. Finally, 51 patients (including 20 uncategorized patients) and 35 patients were retained for binary- and ternary-class survival analysis, respectively. The detailed statistical information of the testing cohorts can be found in Supplementary Table 2.

To make a fair comparison, three specific criteria were adopted to aggregate the TNM stages into TNM-2 (I, II vs. III), TNM-3 (I vs. II vs. III), and TNM-6 (I vs. IIA vs. IIB vs. IIIA vs. IIIB vs. IIIC). The survival analysis results are provided in Table 3 and Figs. S4–S9. According to the univariate analysis results shown in Table 3 and Supplementary Fig. 4, the C-Index of the binary-class  $CG_{Signature}$  was 0.699 (*p* value < 0.0001), outperforming TNM-2 with an increase of 0.04. We further combined the TNM-2 with binary-class  $CG_{Signature}$  for survival analysis (Supplementary Fig. 6), which achieved the highest C-Index of 0.748 (*p* value < 0.0001), which was higher than TNM-6 by 0.034 (Supplementary Fig. 5). In ternary-class univariate Cox regression analysis, we compared the results of TNM-3, TNM-6, and the ternary-class  $CG_{Signature}$ . More specifically, C-Index of the ternary  $CG_{Signature}$  was 0.823 (*p* value < 0.0001, Fig. 4b), which was superior to the TNM-3 (Supplementary Fig. 8) and TNM-6 (Supplementary Fig. 9) with an increase of 0.191 and 0.142, respectively. Results of multivariate Cox regression analyses on both binary- and ternary-class test cohorts indicate that our proposed  $CG_{Signature}$  (i.e., the only one with *p* value < 0.05) can serve as an independent risk factor compared with other factors derived from TNM stages (all had *p* value > 0.05). Univariate and multivariate analyses demonstrate the  $CG_{Signature}$  is capable of discriminating and stratifying gastric cancer patients into groups of different prognoses better than the TNM staging system. Moreover, we note that the prognostic power can be even further enhanced by integrating the  $CG_{Signature}$  predictions and the TNM stages for survival analysis, such as the  $CG_{Signature}$  + TNM – 2 in Table 3 and Supplementary Fig. 6.

To summarize, by combining the spatial information from the mlHC images,  $CG_{Signature}$  has demonstrated outstanding

performance in survival analysis and achieved a better or at least comparable performance when compared with the TNM staging system. The results suggest that effective prognostic features can indeed be captured by  $CG_{Signature}$ , which suggests a powerful method complementary to the current TNM staging system.

### Framelet decomposition for cell-graph

To examine the capacity of Cell-Graph to capture useful spatial features from mlHC images, we conducted a framelet decomposition on the whole mlHC images. The framelet transforms (including framelet decomposition and reconstruction) have proved an important tool for distilling multi-resolution information in low-pass and high-passes from the graph data<sup>33–37</sup>.

We extracted low-pass and high-pass information of six types of features, corresponding to six different biomarkers DAPI, PAN-CK, CD8, CD68, FOXP3, and PD-L1. Tables S3–S11 show the low-pass and high-pass coefficients of the framelet decomposition on mlHC images of short-term, medium-term, and long-term survivors on the entire mlHC images. For the selected samples, no significant differences were observed from the low-pass channel. However, major differences can be observed from the high-pass channel on the selected samples. More specifically, remarkable signal differences can be seen from the high-pass channel-1 and channel-2 in terms of the features of Cell Area and Nucleus Perimeter (summarized in Supplementary Tables 6–11). These differences highlight the important prognostic value of cell morphological information of the TME, which is consistent with the prognostic value of different types of cell features.

### DISCUSSION

In this study, we developed a GNN-based approach, Cell-Graph Signature ( $CG_{Signature}$ ), which is capable of predicting the prognosis of gastric cancer patients from Cell-Graphs extracted from mlHC images. Extensive benchmarking tests on multi-run model training, fivefold cross-validation, and independent tests



demonstrate that  $CG_{Signature}$  can accurately predict the prognosis on both binary- and ternary-class classification tasks. We designed and compared the performance of four different GNN architectures, including GINSag, GCNTopK, GCNSag, and GINTopK. As a result, GINTopK achieved the best performance when compared with the other three GNN architectures (GINSag, GCNTopK, and GCNSag) on the same data sets. Feature ablation studies showed that the nucleus optical feature (DAPI) and cell morphological features are essential node features and contributed most to the prognosis prediction, which indicates the potential pivotal roles of nuclear and cell morphology in gastric cancer progression. In survival analysis,  $CG_{Signature}$  clearly outperformed the AJCC 8th TNM staging system in terms of C-Index (0.823, 95% CI: 0.748–0.899) using the ternary-classification model. In particular, we notice that  $CG_{Signature}$  achieved better or comparable performance with the TNM staging system when using the binary-classification model. These results of survival analysis indicate that  $CG_{Signature}$  provides more prognostic power than the existing TNM staging system and can help pinpoint patients who may benefit from more tailored and personalized therapy. Moreover, wavelet decomposition results suggest that Cell-Graphs can indeed capture certain important spatial features informative for classifying patient survival. Although many previous studies of prognosis prediction also achieved promising results, the majority of such studies were only limited to a specific subtype or stage of cancer. Nevertheless, in this study, we show that the proposed  $CG_{Signature}$  method is applicable to gastric cancer patients of all subtypes across all TNM stages. Moreover,  $CG_{Signature}$  achieved better performance when stratifying test patient cohorts into different groups of prognosis, which has proven a powerful prognostic predictor for gastric cancer.

One caveat of the current study is that we could only obtain mIHC image data for a limited number of patients, and accordingly, the performance of the  $CG_{Signature}$  was only benchmarked on a gastric cancer patient dataset with a limited size. In addition, we were not able to collect a sufficient number of stage IV patients, and thus the model performance needs further verification and improvement when more data become available in the future. Thus, in future studies, it would be important to evaluate the performance of GNNs based on Cell-Graph data from mIHC images in much larger and/or multi-center patient cohorts, as well as additional tumor types (in addition to gastric cancer), when more data become available. Exploration of the prognostic value of the  $CG_{Signature}$  method on datasets of other cancer types would surely be needed to verify its utility and capability. Additionally, future extension of the capability of  $CG_{Signature}$  by using whole-slide images and other biomarkers in mIHC/mIF staining, for example, holds great potential for a more comprehensive analysis of the TME<sup>17</sup>; this will in turn serve to better inform the training of more accurate GNN models. The continuing development of cutting-edge, robust, and broadly applicable Cell-Graph-based biomarker discovery algorithms is valuable and desirable to better inform and transform the medical care of cancer patients.

## METHODS

### Dataset

The gastric cancer samples were collected and stained with mIHC technique and prepared as two batches of tissue microarray<sup>38</sup>, in which all the samples were arranged in the matrix configuration. Then the two tissue microarrays were scanned by a digital microscope (brand: Vectra Polaris) under the magnification of  $\times 40$  with each pixel representing 0.5  $\mu\text{m}$ . Totally, 181 mIHC images of cancer tissues were curated as the initial datasets. After excluding patients whose follow-up data were not available, 172 mIHC images were retained and used for model training and benchmarking. The OS time of the patients ranges from 0 to 88 months, as shown in Fig. 1f. Detailed clinical characteristics and a statistical summary

of the cohort are provided in Supplementary Table 1. Fifty-nine patients were still alive at the time of the last follow-up. All the images were stained using multiplexed immunohistochemistry of seven colors and reagents to identify the specific cell types. In this study, cells were stained with antibodies of Pan-CK, Foxp3, CD8, PD-L1, CD68, CD163, and DAPI. Detailed information on these antibodies can be found in Supplementary Table 12. The dataset was randomly partitioned into the training, validation, and test subsets with the ratios of 0.64, 0.16, and 0.20 at the patient level. In addition, datasets for fivefold cross-validation were also prepared. The gastric cancer patient cohort was obtained from Shanghai Jiao Tong University, Ruijin Hospital. This study was approved by the Ethics Committee of Ruijin Hospital, Shanghai Jiao Tong University School of Medicine (ID: 2021-194). Written informed consent was obtained from all patients.

### Label generation

In this study, the survival prediction was formulated as a classification problem in the form of either binary- or ternary classification. To explore the prognostic value of the Cell-Graphs extracted from the gastric cancer TME, the survival time of the cohort was categorized into two and three classes, and used as labels for training binary- and ternary-classification models based on GNNs. In binary classification, 82 patients with a survival time of fewer than 24 months were annotated as short-term while 70 patients with a survival time of longer than 48 months were annotated as long-term. 20 patients with survival times between 24 and 48 months were removed from the training data set, and denoted as uncategorized patients. For the ternary classification, 12 months and 50 months were, respectively, used as the thresholds to divide patients into short-, medium-, and long-term, with the corresponding patient numbers of 51, 60, and 61, respectively.

### Cell segmentation

After digitization, the mIHC images were pre-processed using the pathology software HALO (Indica Labs) for cell segmentation and feature extraction. The extracted information was subsequently saved as a CSV file in which each row represents the features of a cell (as shown in Table 1), including the cell locations, optical features of stained cells, and morphology features. Thirty-five such features were selected as the node features for each cell. Detailed information can be found in the “Node attributes” section.

### Sub-sampling

Each mIHC staining image contains around 7000–13,000 cells. In particular, we conducted the sub-sampling when generating the Cell-Graphs. By treating each cell as a node in the Cell-Graph, we limited the graph size to no more than 100 nodes. A non-overlap sliding window was then applied to extract the local regions that contained  $\sim 100$  cells from the mIHC images. As a result, we obtained 16951 Cell-Graphs, which would be used for GNN model training and testing. The extracted Cell-Graphs from one mIHC image was annotated with the same label as that of the corresponding mIHC image. The performance of the GNN models was firstly assessed at the Cell-Graph level; After that, the prediction outputs of all Cell-Graphs were aggregated to generate the votes for the final prediction outcome at the patient level.

### Cell-Graph construction

According to the previous study on the TME<sup>19</sup>, we assumed that the maximum effective distance was 20  $\mu\text{m}$  between immune and tumor cells [], which is equivalent to 40 pixels in the magnification of this study. We calculated the Euclidean distance between any pair of cells, and used this distance to define the edge weight between them according to the equations (1) and (2) shown below.

For the  $i$ th and  $j$ th cells with Cartesian coordinates  $(x_i, y_i)$  and  $(x_j, y_j)$  (which use pixel as the unit) in the same mIHC image, their Euclidean distance can be calculated as follows:

$$d(i, j) = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}. \quad (1)$$

The weight between the  $i$ th and  $j$ th cells is assigned as follows:

$$w_{ij} := \begin{cases} 40/d(i, j), & d(i, j) \leq 40 \text{ pixel}, \\ 0, & \text{otherwise}. \end{cases} \quad (2)$$

where 0 denotes that there is no interaction between the cell  $i$  and  $j$ . After sub-sampling, a number of Cell-Graphs (up to 100 nodes) were extracted and annotated, with the weight (2) of the edge between a given pair of cells.

### Node attributes

GNN is a powerful deep learning approach that can efficiently extract features from graph-structured data. In the present study, we focused on distilling five morphology features and 30 optical features generated by six staining biomarkers as the attributes of the node for each cell, including DAPI, PAN-CK, CD8, CD68, FOXP3, and PD-L1. The five morphology features include cell area, cytoplasm area, nucleus area, nucleus perimeter, and nucleus roundness. The optical features of each biomarker are comprised of the positive, positive nucleus, positive cytoplasm, nucleus intensity, and cytoplasm intensity. As a result, a total of 35 features were extracted for each cell. The detailed list of the features and their data types is listed in Table 1. All the features were linearly normalized to the range of [0, 1] prior to training the GNN models.

### Architecture of the designed GNNs

Graph-structured data are usually represented in the form of  $(x_i, A_i)$ , where  $x_i$  denotes the feature of the node for the  $i$ th graph sample while  $A_i$  represents its adjacency matrix. A GNN has a similar network architecture to that of the traditional convolutional neural network. To address the classification task in this study, we designed the GNN model architecture of  $CG_{Signature}$  which includes four computational units, each with two-layer graph convolution plus one-layer graph pooling followed by three-layer fully connected layers (MLP), before generating the prediction output (Fig. 1).

The graph convolutional layer is responsible for extracting an array of features from the last output array, which mimics the role of CNN convolution. It changes the dimension  $d$  of the feature array but does not change the number of nodes  $N$ . The output of graph convolutional layers is passed on to the graph pooling which compresses the node number by a fractional proportion while in this process usually the key structural information and node features are preserved. The MLP readout will then output the label class. Graph convolution communicates the structural information of the data to the deep network model via the message passing between the neighborhood nodes, which contributes as the key to successfully capturing the geometric feature of the data. In this work, we adopted the GINConv<sup>23</sup> as the graph convolution and TopKPool<sup>22</sup> as the graph pooling method, respectively. The convolutional layer for GIN can be aggregated by

$$\mathbf{X}^{\text{output}} = \text{MLP}((\mathbf{A} + (1 + \epsilon) \cdot \mathbf{I}) \cdot \mathbf{X}^{\text{in}}), \quad (3)$$

where  $\mathbf{X}^{\text{in}} \in \mathbb{R}^{N \times d}$  is the  $d$ -feature matrix on the nodes of the graph with  $N$  nodes for the input layer, and  $\mathbf{A} \in \mathbb{R}^{N \times N}$  is the adjacency matrix of the graph.  $W$  is the filter weight parameter matrix with the size of  $m \times n$  to be learned by the GNNs, where  $n$  is the number of hidden neurons. GINConv is a special neural message passing operator for GNN aggregation.

Our GNN model was trained by connecting multiple layers of graph convolution activated by a ReLU (Rectifier Linear Unit)<sup>39</sup>. The graph pooling, which is used between two consecutive layers, serves to reduce the dimensionality of the feature map so that the network has appropriate amounts of parameters to circumvent over-fitting<sup>40</sup>. Here we used TopKPooling<sup>22</sup> for graph pooling.

There exist different types of GNN models in the machine learning literature<sup>41</sup>. Specifically, we tested the performance of the GINConv+TopKPool model with the other three popular GNN models, i.e., GINConv+SAGPool, GCNConv+TopKPool, and GCNConv+SAGPool. The results showed that the chosen model (GINConv+TopKPool) achieved the highest AUROC value and stable training performance. Refer to Figure 2a, b for a detailed illustration of the results.

### Hyperparameter optimization

We fine-tuned the hyperparameters for the GNN models with the assistance of HyperOPT<sup>27</sup> and Ray<sup>28</sup>, where the network architecture and batch size were fixed. The hyperparameters were searched within the range as shown in Table 4. More specifically, the best-performing model used the following hyperparameters: learning rate  $5 \times 10^{-4}$ , weight decay rate  $10^{-4}$ , number of hidden neurons 512, pooling ratio 0.5, number of hidden layers 4, batch size 256, and maximal number of epochs 200 with the early stopping strategy.

**Table 4.** Search space for hyperparameters of GNN models.

Hyperparameter	Searching space
Learning rate	$10^{-4}$ , $5 \times 10^{-4}$ , $10^{-3}$
Weight decay ( $L_2$ )	$10^{-4}$ , $5 \times 10^{-4}$ , $10^{-3}$
Hidden units	256, 512
Pooling ratio	0.5, 0.65, 0.75

### Prediction aggregation to assess the patient-level performance

The model performance was evaluated at the Cell-Graph level. After the model was optimized, the patient-level performance of the model was calculated by aggregating the prediction results produced by the optimized model. In particular, we fed Cell-Graphs of the test dataset to the optimal model to predict the label for each of them. Since hundreds of Cell-Graphs were subsampled from the mIHC images of the patient, hundreds of the predictions were also made for a given patient. To generate the patient-level prediction for a patient, we calculated the proportion of Cell-Graphs belonging to a specific class, and then classified the patient as the group that received the largest proportion of the Cell-Graphs.

### Framelet analysis to facilitate interpretation of the model prediction

From the mathematical perspective, the *framelet system*<sup>33–35,37</sup> refers to a set of functions that provide a multi-scale representation of graph-structured data, which has a similar property to the traditional wavelets in the Euclidean space. Using the framelet transforms, we can decompose the graph features into low-pass and high-pass frequencies as the extracted features to train network models, via the framelet-based graph convolution.

Suppose  $\{(\lambda_\ell, u_\ell)\}_{\ell=1}^N$  are the pairs of the eigenvalue and eigenvector for the graph Laplacian  $\mathcal{L}$  of a graph  $\mathcal{G}$  with  $N$  nodes. The (undecimated) framelets at the scale level  $j = 1, \dots, J$  for graph  $\mathcal{G}$  with the above scaling functions can be defined, for  $n = 1, \dots, r$ , as follows:

$$\begin{aligned} \varphi_{j,p}(v) &= \sum_{\ell=1}^N \hat{\alpha} \left( \frac{\lambda_\ell}{2^j} \right) \overline{u_\ell(p)} u_\ell(v) \\ \psi_{j,p}^n(v) &= \sum_{\ell=1}^N \hat{\beta}^{(n)} \left( \frac{\lambda_\ell}{2^j} \right) \overline{u_\ell(p)} u_\ell(v), \end{aligned} \quad (4)$$

where  $\varphi_{j,p}$  and  $\psi_{j,p}^n$  are the low-pass and high-pass framelets translated at the graph node  $p$ . In the framelet analysis above, we have shown the low-pass and high-pass framelet coefficients  $v_{j,p}$  and  $w_{j,p}^n$  for a signal  $f$  on graph  $\mathcal{G}$ . They are the projections  $\langle \varphi_{j,p}, f \rangle$  and  $\langle \psi_{j,p}^n, f \rangle$  of the graph signal onto framelets at the scale  $j$  and node  $p$ . The construction of the framelet system and the framelet transforms rely on the filter bank (a collection of filters) to calculate framelet coefficients. Here we used the filter bank of the Haar-type filters for the experiments<sup>33,37</sup>. The dilation factor is  $2^j$  with the dilation (base) 2 for a natural number  $j$ , where  $j$  indicates the scale level and  $2^j$  is the scale of the framelet. A bigger value of  $j$  indicates that the corresponding framelet coefficient carries more detailed information about the graph signal.

The above framelet system is a *tight frame*, which provides an exact representation of any  $L_2$  function on the graph. This guarantees that the framelet coefficients have a unique representation of a graph signal. Accordingly, the framelet coefficients can fully reflect the feature of the signal. Moreover, the coefficients decompose the signal at multi scales and can be used to observe whether a particular scale or the high-pass or low-pass frequencies contain a more important feature of the data.

### Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## DATA AVAILABILITY

The data used for the main analyses presented here is available for non-commercial use and can be accessible by request.

## CODE AVAILABILITY

All the related scripts and code are publicly available and can be downloaded at [https://github.com/docurdt/Cell-Graph\\_Signature.git](https://github.com/docurdt/Cell-Graph_Signature.git).

## MATERIALS AVAILABILITY

Hardware: all the experiments were performed using PyTorch Geometric<sup>42</sup> on a server with Intel(R) Core(TM) i9-9820X 230 CPU 3.30GHz, NVIDIA GeForce RTX 2080Ti, and NVIDIA TITAN V GV100.

Received: 5 November 2021; Accepted: 17 May 2022;

Published online: 23 June 2022

## REFERENCES

- Sung, H. et al. Global cancer statistics 2020: Globocan estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA Cancer J. Clin.* **71**, 209–249 (2021).
- Etemadi, A. et al. The global, regional, and national burden of stomach cancer in 195 countries, 1990–2017: a systematic analysis for the global burden of disease study 2017. *Lancet Gastroenterol. Hepatol.* **5**, 42–54 (2020).
- Lababede, O. & Meziane, M. A. The eighth edition of TNM staging of lung cancer: reference chart and diagrams. *Oncologist* **23**, 844 (2018).
- Bang, Y.-J. et al. Adjuvant capecitabine and oxaliplatin for gastric cancer after d2 gastrectomy (classic): a phase 3 open-label, randomised controlled trial. *Lancet* **379**, 315–321 (2012).
- Noh, S. H. et al. Adjuvant capecitabine plus oxaliplatin for gastric cancer after d2 gastrectomy (classic): 5-year follow-up of an open-label, randomised phase 3 trial. *Lancet Oncol.* **15**, 1389–1396 (2014).
- Sasako, M. et al. Gastric cancer working group report. *Jpn J. Clin. Oncol.* **40**, i28–i37 (2010).
- Yu, K. H. et al. Predicting non-small cell lung cancer prognosis by fully automated microscopic pathology image features. *Nat. Commun.* **7**, 1–10 (2016).
- Mobadersany, P. et al. Predicting cancer outcomes from histology and genomics using convolutional networks. *Proc. Natl. Acad. Sci. USA* **115**, E2970–E2979 (2018).
- Jiang, Y. et al. Immunomarker support vector machine classifier for prediction of gastric cancer survival and adjuvant chemotherapeutic benefit. *Clin. Cancer Res.* **24**, 5574–5584 (2018).
- Wulczyn, E. et al. Deep learning-based survival prediction for multiple cancer types using histopathology images. *PLoS ONE* **15**, 1–18 (2020).
- Jiang, D. et al. A machine learning-based prognostic predictor for stage III colon cancer. *Sci. Rep.* **10**, 1–9 (2020).
- Dimitriou, N., Arandjelović, O., Harrison, D. J. & Caie, P. D. A principled machine learning framework improves accuracy of stage II colorectal cancer prognosis. *NPJ Digit. Med.* **1**, 1–9 (2018).
- Yener, B. Cell-graphs: image-driven modeling of structure-function relationship. *Commun. ACM* **60**, 74–84 (2016).
- Wang, M. et al. High-dimensional analyses reveal a distinct role of t-cell subsets in the immune microenvironment of gastric cancer. *Clin. Transl. Immunol.* **9**, e1127 (2020).
- Huang, Y.-K. et al. Macrophage spatial heterogeneity in gastric cancer defined by multiplex immunohistochemistry. *Nat. Commun.* **10**, 1–15 (2019).
- Carstens, J. L. et al. Spatial computation of intratumoral T cells correlates with survival of patients with pancreatic cancer. *Nat. Commun.* **8**, 1–13 (2017).
- Berry, S. et al. Analysis of multispectral imaging with the AstroPath platform informs efficacy of PD-1 blockade. *Science* **372**, eaba2609 (2021).
- Lu, M. Y., Sater, H. A. & Mahmood, F. Multiplex computational pathology for treatment response prediction. *Cancer Cell* **39**, 1053–1055 (2021).
- Barua, S. et al. Spatial interaction of tumor cells and regulatory T cells correlates with survival in non-small cell lung cancer. *Lung Cancer* **117**, 73–79 (2018).
- Gunduz, C., Yener, B. & Gultekin, S. H. The cell graphs of cancer. *Bioinformatics* **20**, i145–i151 (2004).
- Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M. & Monfardini, G. The graph neural network model. *IEEE Trans. Neural Netw.* **20**, 61–80 (2008).
- Gao, H. & Ji, S. Graph U-Nets. In *ICML*, 2083–2092 (2019).
- Xu, K., Hu, W., Leskovec, J. & Jegelka, S. How powerful are graph neural networks? In *ICLR* <https://arxiv.org/abs/1810.00826> (2019).

- Wang, J. et al. scGNN is a novel graph neural network framework for single-cell RNA-Seq analyses. *Nat. Commun.* **12**, 1–11 (2021).
- Zhao, T., Hu, Y., Valsdottir, L. R., Zang, T. & Peng, J. Identifying drug–target interactions based on graph convolutional network and deep neural network. *Brief. Bioinform.* **22**, 2141–2150 (2021).
- Harrell, F. E., Califf, R. M., Pryor, D. B., Lee, K. L. & Rosati, R. A. Evaluating the yield of medical tests. *JAMA* <https://doi.org/10.1001/jama.1982.03320430047030> (1982).
- Bergstra, J., Yamini, D. & Cox, D. D. Hyperopt: a python library for optimizing the hyperparameters of machine learning algorithms. In *Proceedings of the 12th Python in Science Conference*, 13, 20 (Citeseer, 2013).
- Nishihara, R. et al. Real-time machine learning: the missing pieces. In *Workshop on Relational Representation Learning (R2L) at NIPS*, 106–110 (2018).
- Kipf, T. N. & Welling, M. Semi-supervised classification with graph convolutional networks. In *ICLR* <https://arxiv.org/abs/1609.02907> (2017).
- Cangea, C., Veličković, P., Jovanović, N., Kipf, T. & Liò, P. Towards sparse hierarchical graph classifiers. In *Workshop on Relational Representation Learning (R2L) at NIPS* <https://arxiv.org/abs/1811.01287> (2018).
- Knyazev, B., Taylor, G. W. & Amer, M. R. Understanding attention in graph neural networks. In *NeurIPS* <https://arxiv.org/abs/1905.02850v3> (2019).
- Lee, J., Lee, I. & Kang, J. Self-attention graph pooling. In *ICML*, 3734–3743 (2019).
- Dong, B. Sparse representation on graphs by tight wavelet frames and applications. *Appl. Comput. Harmon. Anal.* **42**, 452–479 (2017).
- Zheng, X., Zhou, B., Wang, Y. G. & Zhuang, X. Decimated framelet system on graphs and fast G-framelet transforms. *J. Mach. Learn. Res.* **23**, 1–68 (2022).
- Wang, Y. G. & Zhuang, X. Tight framelets on graphs for multiscale data analysis. In *Wavelets and Sparsity XVIII*, vol. 11138, 111380B (International Society for Optics and Photonics, 2019).
- Wang, Y. G. & Zhuang, X. Tight framelets and fast framelet filter bank transforms on manifolds. *Appl. Comput. Harmon. Anal.* **48**, 64–95 (2020).
- Zheng, X. et al. How framelets enhance graph neural networks. <https://arxiv.org/abs/2102.06986> (2021).
- Voduc, D., Kenney, C. & Nielsen, T. O. Tissue microarrays in clinical oncology. *Semin. Radiat. Oncol.* **18**, 89–97 (2008).
- Agarap, A. F. Deep learning using rectified linear units (ReLU). <https://arxiv.org/abs/1803.08375> (2019).
- Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet classification with deep convolutional neural networks. In *NeurIPS*, 1097–1105 (2012).
- Wu, Z. et al. A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems* <https://arxiv.org/abs/1901.00596> (2020).
- Fey, M. & Lenssen, J. E. Fast graph representation learning with PyTorch Geometric. In *ICLR 2019 Workshop on Representation Learning on Graphs and Manifolds* <https://arxiv.org/abs/1903.02428> (2019).

## ACKNOWLEDGEMENTS

The authors would like to thank the Ruijin Hospital affiliated with Shanghai Jiao Tong University School of Medicine for providing support for this project. We would also like to thank all the collaborators and colleagues for the enlightening discussions and feedback. This work was supported by the Major Inter-Disciplinary Research (IDR) Grant awarded by Monash University and project grants from the National Natural Science Foundation of China (NSFC) (81871274, 82071811, and 31670905).

## AUTHOR CONTRIBUTIONS

Y.W. and Y.G.W. contributed equally to this work. Y.W. and Y.G.W. conceived and conducted the experiments, analyzed the results, and wrote the first draft. Y.W., Y.G.W., C.H., M.L., P.L., G.I.W., and J.S. were responsible for the methodology and experiment design. Y.F., N.O., T.K., R.J.D., J.Z., A.B., and G.M. helped to analyze the results. I.S., Q.G., Y.H., and D.X. processed and curated the mIHC data. P.L., D.X., G.I.W., and J.S. supervised this study. All authors reviewed or revised the manuscript.

## COMPETING INTERESTS

The authors declare no competing interests.

## ETHICS DECLARATION

This study was approved by the Ethics Committee of Ruijin Hospital, Shanghai Jiao Tong University School of Medicine (ID: 2021-194). All researchers were blinded to the patient private data during the experimental analysis.

**ADDITIONAL INFORMATION**

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41698-022-00285-5>.

**Correspondence** and requests for materials should be addressed to Pietro Liò, Dakang Xu, Geoffrey I. Webb or Jiangning Song.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, 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 license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license 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 license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022