

# **Two-stage biologically interpretable neural-network models for liver cancer prognosis prediction using histopathology and transcriptomic data**

Zhucheng Zhan<sup>1\*</sup>, Zheng Jing<sup>2\*</sup>, Bing He<sup>3</sup>, Noshad Hosseni<sup>3</sup>, Maria Westerhoff<sup>4</sup>, Eun-Young Choi<sup>4</sup>,  
Lana X. Garmire<sup>3§</sup>

<sup>1</sup>School of Science and Engineering, Chinese University of Hong Kong Shenzhen Campus, Shenzhen, P.R. China

<sup>2</sup>Department of Applied Statistics, University of Michigan, Ann Arbor, MI USA

<sup>3</sup>Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI USA

<sup>4</sup>Department of Pathology, University of Michigan, Ann Arbor, MI, USA

\*: These authors contributed equally to the work

§ corresponding author, email: [lgarmire@med.umich.edu](mailto:lgarmire@med.umich.edu)

## ABSTRACT

**Purpose:** Pathological images are easily accessible data with the potential as prognostic biomarkers. Moreover, integration of heterogeneous data types from multi-modality, such as pathological image and gene expression data, is invaluable to help predicting cancer patient survival. However, the analytical challenges are significant.

**Experimental Design:** Here we take the hepatocellular carcinoma (HCC) pathological image features extracted by CellProfiler, and apply them as the input for Cox-nnet, a neural network-based prognosis. We compare this model with conventional Cox-PH models, using C-index and log ranked p-values on HCC testing samples. Further, to integrate pathological image and gene expression data of the same patients, we innovatively construct a two-stage Cox-nnet model, and compare it with another complex neural network model PAGE-Net.

**Results:** pathological image based prognosis prediction using Cox-nnet (median C-index 0.74 and log-rank p-value  $4e-6$ ) is significantly more accurate than Cox-PH model (median C-index 0.72 and log-rank p-value of  $3e-4$ ). Moreover, the two-stage Cox-nnet complex model combining histopathology image and transcriptomics RNA-Seq data achieves better prognosis prediction, with a median C-index of 0.75 and log-rank p-value of  $6e-7$  in the testing datasets. The results are much more accurate than PAGE-Net, a CNN based complex model (median C-index of 0.67 and log-rank

p-value of 0.02). Imaging features present additional predictive information to gene expression features, as the combined model is much more accurate than the model with gene expression alone (median C-index 0.70). Pathological image features are modestly correlated with gene expression. Genes having correlations to top imaging features have known associations with HCC patient survival and morphogenesis of liver tissue.

**Conclusion:** This work provides two-stage Cox-nnet, a new class of biologically relevant and relatively interpretable models, to integrate multi-modal and multiple types of data for survival prediction.

**Key words:** prognosis, survival, prediction, neural network, modelling, Cox proportional hazards, pathology, image, gene expression, omics, RNA-Seq, data integration

## INTRODUCTION

Prognosis prediction is important for providing effective disease monitoring and management. Various biomaterials have been proposed as potential biomarkers to predict patient survival. Among them, hematoxylin and eosin (H&E) stained histopathological images, are very attractive materials to extract biomarker features. Compared to genomics materials, such as RNA-Seq transcriptomics, these images are much more easily accessible and cheaper to obtain, through processing archived formalin-fixed paraffin-embedded (FFPE) Blocks. In H&E staining, the hematoxylin is oxidized into phematein, a basic dye which stains acidic (basophilic) tissue components (ribosomes, nuclei, and

rough endoplasmic reticulum) into darker purple color. Whereas acidic eosin dye stains other protein structures of the tissue (stroma, cytoplasm, muscle fibers) into a pink color. As patients' survival information is retrospectively available in electronic medical record data and FFPE blocks are routinely collected clinically, the histopathology images can be generated and used for highly valuable and predictive prognosis models.

Previously, we developed a neural network model called Cox-nnet to predict patient survival, using transcriptomics data [1]. Cox-nnet is an alternative to the conventional methods, such as Cox proportional hazards (Cox-PH) methods with LASSO or ridge penalization. We have demonstrated that Cox-nnet is more optimized for survival prediction from high throughput gene expression data, with comparable or better performance than other conventional methods, including Cox-PH, Random Survival Forests [2] and Coxboost [3]. Moreover, Cox-nnet reveals much richer biological information, at both the pathway and gene levels, through analysing the survival related “surrogate features” represented as the hidden layer nodes in Cox-nnet. However, it remains to be explored whether other data types that are less biologically intuitive than genomics data, such as histopathology imaging data, are also suitable input features for Cox-nnet. Moreover other neural network based models have been proposed [15-19], and some of them were designed to handle multi-modal data. For example, PAGE-Net is a complex neural network model that has a convolutional neural network (CNN) layer followed by pooling and a genomics model involved in transformation of the gene layer to pathway layer. The genomics neural network portion is followed by two hidden layers, the latter of which is combined with the image neural network model to predict glioblastoma patient survival. Though PAGE-Net uses CNN, the resulting predictive C-index value based on imaging data appeared almost random (C-index=0.509), raising the concern of overfitting. It is therefore important to test if a model built upon Cox-nnet, using pre-extracted, biologically informative features, can combine multiple types of data, eg. imaging and genomics data, and if so, how well it performs relative to models such as PAGE-Net.

In this study, we extend Cox-nnet to take up pathological image features extracted from imaging processing tool *CellProfiler* [4], and compare the predictive performance of Cox-nnet relative to Cox proportional hazards, the standard method for survival analysis, which was also the second best method in the previous survival prediction study using pan-cancer datasets [1]. Moreover, we propose a new type of two-stage complex Cox-nnet model, which combines the hidden node features from multiple first-stage Cox-nnet models, and then use these combined features as the input nodes to train a second stage Cox-nnet model. We applied the models on TCGA hepatocellular carcinoma (HCC), which we had previously gained domain experience on [5-6]. Hepatocellular carcinoma (HCC) is the most prevalent type of liver cancer that accounts for 70%-90% of all liver cancer cases. It is a devastating disease with poor prognosis, where the 5-year survival rate is only 12% [20]. And the prognosis prediction becomes very challenging due to the high level of heterogeneity in HCC as well as the complex etiologic factors. Limited treatment strategies in HCC, relative to other cancers, also imposes an urgent need to develop tools for patient survival prediction. As comparison, we also evaluated the performance of another CNN based model called PAGE-Net, and showed that Cox-nnet achieves higher accuracy in testing data.

## **METHODS**

### **Datasets**

The histopathology images and their associated clinical information are downloaded from The Cancer Genome Atlas (TCGA). A total of 384 liver tumor images are collected. Among them 322 samples are clearly identified with tumor regions by pathology inspection. Among these samples, 290 have gene expression RNA-Seq data, and thus are selected for pathology-gene expression integrated prognosis prediction. The gene expression RNA-Seq dataset is also downloaded from TCGA, each feature was then normalized into RPKM using the function *ProcessRNASeqData* by TCGA-Assembler [21].

## Tumor Image Pre-processing

For each FFPE image stained with H&E, the tumor regions are labelled by pathologists at University of Michigan. The tumor regions are then extracted using Aperio software *ImageScope* [7]. To reduce computational complexities, each extracted tumor region is divided into non-overlapping 1000 by 1000 pixel tiles. The density of each tile is computed as the summation of red, green and blue values, and 10 tiles with the highest density are selected for further feature extraction, following the guideline of others [8]. To ensure that the quantitative features are measured under the same scale, the red, green and blue values are rescaled for each image. Image #128 with the standard background color (patient barcode: TCGA-DD-A73D) is selected as the reference image for the others to be compared with. The means of red, green and blue values of the reference image are computed and the rest of the images are normalized by the scaling factors of the means of red, green, blue values relative to those of the reference image.

## Feature extraction from the images

CellProfiler is used for feature extraction [14]. Images are first preprocessed by '*UnmixColors*' module to H&E stains for further analysis. '*IdentifyPrimaryObject*' module is used to detect unrelated tissue folds and then removed by '*MaskImage*' module to increase the accuracy for detection of tumour cells. Nuclei of tumour cells are then identified by '*IdentifyPrimaryObject*' module again with parameters set by Otsu algorithm. The identified nuclei objects are utilised by '*IdentifySecondaryObject*' module to detect the cell body objects and cytoplasm objects which surround the nuclei. Related biological features are computed from the detected objects, by a series of feature extraction modules, including '*MeasureGranularity*', '*MeasureObjectSizeShape*', '*MeasureObjectIntensity*', '*MeasureObjectIntensityDistribution*', '*MeasureTexture*', '*MeasureImageAreaOccupied*', '*MeasureCorrelation*', '*MeasureImageIntensity*' and '*MeasureObjectNeighbors*'. To aggregate the features from the primary and secondary objects, the related summary statistics (mean, median, standard deviation and quartiles) are then calculated to summarize data from object level to image

level, yielding 2429 features in total. Each patient is represented by 10 images, and the median of each feature is selected to represent the patient's image biological feature.

### **Survival prediction models**

**Cox-nnet model:** The Cox-nnet model is implemented in the Python package named Cox-nnet [1]. Current implementation of Cox-nnet is a fully connected, two-layer neural network model, with a hidden layer and an output layer for cox regression. The drop-out method is used to avoid overfitting. We used a hold-out method by randomly splitting the dataset to 80% training set and 20% testing set. We used grid search and 5-fold cross-validation to optimise the hyper-parameters for the deep learning model on the selected training set. The model is then trained under the optimised hyperparameter setting using the training set and further evaluated on the remaining testing set, the procedure is repeated 20 times to assess the average performance. More details about Cox-nnet is described earlier in Ching et al [1].

**Cox proportional hazards Model:** Since the number of features produced by *CellProfiler* exceeds the sample size, an elastic net Cox proportional hazard model is built to select features and compute the prognosis index (PI) [9]. Function *cv.glmnet* in the *Glmnet* R package is used to perform cross-validation to select the tuning parameter *lambda*. The parameter *alpha* that controls the trade-off between quadratic penalty and linear penalty is selected using grid search. Same hold-out setting is employed by training the model using 80% randomly selected data and evaluated on the remaining 20% testing set. The procedure is repeated 20 times to calculate the mean accuracy of the model.

**Two-stage Cox-nnet model:** The two-stage Cox-nnet model has two phases, as indicated in the name. For the first stage, we construct two separate Cox-nnet models in parallel, one for the image data and the other one for gene expression data. For each model, we optimize the hyper-parameters using grid search under 5-fold cross-validation, as described earlier. In the second stage, we extract and combine the nodes of the hidden layer from each Cox-nnet model as the new input features for a new Cox-nnet

model. We construct and evaluate the second stage Cox-nnet model with the same parameter-optimisation strategy as in the first-stage.

**PAGE-Net model:** it is another neural-network method that can combine imaging and genomics (eg. gene expression) information to predict patient survival [15]. The imaging prediction module is very complex, with a patch-wide pre-trained convolutional neural network (CNN) layer followed by pooling them together for another neural network. The genomics model involved transformation of gene layer to pathway layer, and then followed by two hidden layers, the latter of which is combined with the image NN to predict patient survival. Due to this issue, we only repeated the step of integrative layer training with different train-test splits, but we do not repeat the steps of CNN-pretrain or feature extraction due to the running time issue.

### Model evaluation

Similar to the previous studies [1,5,6], we also use concordant index (C-index) and log-rank p-value as the metrics to evaluate model accuracy. C-index signifies the fraction of all pairs of individuals whose predicted survival times are correctly ordered and is based on Harrell C statistics. The equation is as follows:

$$C = \frac{\# \text{ concordant pairs}}{\# \text{ concordant pairs} + \# \text{ discordant pairs}} = \frac{\sum_{i \neq j} 1\{\eta_i < \eta_j\} 1\{T_i > T_j\} d_j}{\sum_{i \neq j} 1\{T_i > T_j\} d_j} ,$$

where  $\eta$  is the predicted risk score,  $T$  is the “time-to-event” response,  $d$  is an auxiliary variable such that  $d=1$  if event is observed and  $d=0$  if patient is censored. A C-index of 1 means the model fits the survival data perfectly, whereas a score around 0.50 means randomness. In practice, a C-index around 0.70 indicates a good model. As both Cox-nnet and Cox-PH models quantify the patient's prognosis by log hazard ratios, we use the predicted median hazard ratios to stratify patients into two risk groups (high vs. low survival risk groups). We also compute the log-rank p-value to test if two Kaplan-Meier survival curves produced by the dichotomised patients are significantly different, similar to earlier reports [5-6, 9, 22-24].



## **Feature evaluation**

The input feature importance score is calculated by drop-out. The values of a variable are set to its mean and the log likelihood of the model is recalculated. The difference between the original log likelihood and the new log likelihood is considered as feature importance [13]. We select 100 features with the highest feature scores from Cox-nnet for association analysis between pathology image and gene expression features. We regress each selected image feature (y) over all the gene expression features (x) using LASSO penalization, where lamda.min is selected as the value that gives the minimum mean error over cross-validation. We use the R-square statistic as the correlation metric.

## **Code availability**

The code for two-stage Cox-nnet, including integration of hidden nodes, feature extraction, and feature analysis are all available at github: <https://github.com/lanagarmire/two-stage-cox-nnet>

## **RESULTS**

### **Overview of Cox-nnet model on pathological image data**

In this study, we tested if pathological images can be used to predict cancer patients. The initial task is to extract image features that can be used as the input for the predictive models. As described in the Methods, pathological images of 322 TCGA HCC patients are individually annotated with tumor contents by pathologists, before being subject to a series of processing steps. The tumor regions of these images then undergo segmentation, and the top 10 tiles (as described in section 2.2) out of 1000 by 1000 tiles are used to represent each patient. These tiles are next normalized for RGB coloring against a common reference sample, and 2429 image features of different categories are extracted by *CellProfiler*. Summary statistics (mean, median, standard deviation and quartiles) are calculated for

each image feature, and the median values of them over 10 tiles are used as the input imaging features for survival prediction.

We applied these imaging features on Cox-nnet, a neuron-network based prognosis prediction method previously developed by our group. The architecture of Cox-nnet is shown in **Figure 1A**. Briefly, Cox-nnet is composed of the input layer, one fully connected hidden layer and an output “proportional hazards” layer. We use 5-fold cross-validation (CV) to find the optimal regularization parameters. Based on the results on RNA-Seq transcriptomics previously, we use dropout as the regularization method. Additionally, to evaluate the results on pathology image data, we compare Cox-nnet with Cox-PH model, the previously 2nd-best prognosis model on RNA-Seq data.

### **Comparison of prognosis prediction between Cox-nnet and Cox-PH over pathology imaging data**

We use two accuracy metrics to evaluate the performance of models in comparison: C-index and log-rank p-values. C-index measures the fraction of all pairs of individuals whose predicted survival times are correctly ordered by the model. The higher C-index, the more accurate the prognosis model is. On the other hand, log-rank p-value tests if the two Kaplan-Meier survival curves based on the survival risk-stratification are significantly different (log-rank p-value  $<0.05$ ). In this study, we stratify the patients by the median score of predicted prognosis index (PI) from the model. As shown in **Figure 2**, on the testing datasets, the median C-index score from Cox-nnet (0.74) is significantly higher ( $p < 0.001$ ) than Cox-PH (0.72). Additionally, the discrimination power of Cox-nnet on patient Kaplan-Meier survival difference (**Figure 3 B and D**) is much better than Cox-PH model (**Figure 3 A and C**), using median PI based survival risk stratification. In the training dataset, Cox-nnet achieves a log-rank p-value of  $1e-12$ , compared to  $5e-9$  for Cox-PH; in the testing dataset, Cox-nnet gives a log-rank p-value of  $4e-6$ , whereas Cox-PH has a log-rank p-value of  $3e-4$ .

We next investigate the top 100 image features according to Cox-nnet ranking (**Supplementary Figure 1**). Interestingly, the most frequent features are those involved in textures of the image, accounting for 48% of raw input features. Intensity and Area/Shape parameters make up the 2nd and 3rd highest categories, with 18% and 15% features. Density, on the other hand, is less important (3%). It is also worth noting that among the 47 features selected by the conventional Cox-PH model, 70% (33) are also found in the top 100 features selected by Cox-nnet, showing the connections between the two models.

### **Two-stage Cox-nnet model to predict prognosis on combined histopathology imaging and gene expression RNA-Seq data**

Multi-modal and multi-type data integration is challenging, particularly for survival prediction. We next ask if we can utilize Cox-nnet framework for such purpose, exemplified by pathology imaging and gene expression RNA-Seq based survival prediction. Towards this, we propose a two-stage Cox-nnet complex model, inspired by other two-stage models in genomics fields [10-12]. The two-stage Cox-nnet model is depicted in **Figure 1B**. For the first stage, we construct two Cox-nnet models in parallel, using the image data and gene expression data of HCC, respectively. For each model, we optimize the hyper-parameters using grid search under 5-fold cross-validation. Then we extract and combine the nodes of the hidden layer from each Cox-nnet model as the new input features for the second-stage Cox-nnet model. We construct and evaluate the second-stage Cox-nnet model with the same parameter-optimisation strategy as in the first-stage.

As shown in **Figure 2**, the resulting two-stage Cox-nnet model yields very good performance, judging by the C-index values. The median C-index scores for the training and testing sets are 0.89 and 0.75, respectively. These C-index values are significantly improved, compared to the Cox-nnet models that are built on either imaging (described earlier) or gene expression RNA-Seq data alone. For example, on the testing datasets, the median C-index score from two-stage Cox-nnet (0.75) is significantly

higher ( $p < 0.0005$ ) than the Cox-nnet model built on gene expression data (0.70). It is also significantly higher ( $p < 0.005$ ) than the Cox-nnet model built on image data (0.74). The superior predictive performance of the two-stage Cox-nnet model is also confirmed by the log-rank p-values in the Kaplan-Meier survival curves (**Figure 4**). It achieves a log-rank p-value of  $6e-7$  in testing data (**Figure 4A**), higher than the Cox-nnet models based on pathological image data (**Figure 4B**) or gene expression RNA-Seq data (**Figure 4C**), which have log-rank p-values of  $4e-6$  and 0.01, respectively.

### **Comparing two-stage Cox-nnet model with other imaging and gene expression based prognosis models**

We compare two-stage Cox-nnet with PAGE-Net, another neural-network method that combines imaging and genomics information to predict patient survival [15]. The imaging prediction module of pagenet is very complex, with a patch-wide pre-trained convolutional neural network (CNN) layer followed by pooling them together for another neural network. The genomics model involves transformation of gene layer to pathway layer, and then followed by two hidden layers, the latter of which is combined with the image NN to predict patient survival. For a fair comparison, we use the same image inputs, tumor-selected images, for both PAGE-Net and Cox-nnet models. We also use the same pathway data used in PAGE-Net paper to construct the gene-pathway layer. We perform the same train-test splits for two models; 290 samples are split into 80% training and 20% testing data. For PAGE-Net, the training set is further split into 90% training and 10% validation, used for selecting hyperparameters and avoiding overfitting, following its code. We repeat the experiment 20 times with different train-test splits.

As shown in **Figure 5A**, on the testing datasets, the median C-index score of 0.75 from the two-stage Cox-nnet model is significantly higher ( $p\text{-value} < 3.7e-6$ ) than that of PAGE-Net (0.67). The C-index values from the PAGE-Net model are much more variable (less stable), compared to those from two-stage Cox-nnet. Moreover, PAGE-Net model appears to have an overfitting issue: the median

C-index score of PAGE-Net model on the training set is very high (0.97), however, its predictability on hold-out testing data is much poorer. Moreover, impractical running time is another concern for PAGE-Net. Even on Graphic Processing Unit (GPUs) of Nvidia V100-PCIE with 16GB of memory each, it takes over a week to pretrain CNN and extract image features from only 290 samples, prohibiting its practical use. Confirming the results in C-index metric, the Kaplan-Meier survival difference on testing data based on two-stage Cox-nnet prediction (**Figure 5D**) is also much better than that of PAGE-Net model (**Figure 5E**). Using median PI based survival risk stratification, Cox-nnet achieves a much better log-rank p-value of  $6e-7$  compared to 0.02 for PAGE-Net, despite that PAGE-Net has higher log-rank p-value of  $9e-24$  in training data (**Figure 5C**).

### **Relationship between histology and gene expression features in the two-stage Cox-nnet model**

We also investigate the correlations between the top imaging features with those RNA-Seq gene expression features. For this we performed two types of correlation analysis. We first examined the pairwise correlations between top histopathology features and top gene expression features. The bipartite graph shows the Pearson's correlations between top 10 histopathology features and top 50 gene expression features (**Figure 6**). Top genes that are associated with top 10 histopathology features include long intergenic non-protein coding RNA 1554 (LINC01554), MAP7 domain containing 2 (MAP7D2), homeobox D9 (HOXD9), mucin 6 (MUC6), keratin 17 (KRT17) and matrix metalloproteinase 7 (MMP7). Gene Set Enrichment (GSEA) analysis on top 1000 genes correlated to each image feature shows that image feature correlated genes are significantly (false discovery rate,  $FDR < 25\%$  as recommended) enriched in pathways-in-cancer (normalized enrichment score,  $NES=1.41$ ,  $FDR=0.18$ ) and focal-adhesion ( $NES=1.22$ ,  $FDR=0.21$ ) pathways, which are upregulated in HCC patients with poor prognosis. Pathways-in-cancer is a pan-pathway that covers multiple important cancer-related signalling pathways, such as PI3-AKT signaling, MAPK signaling and p53 signaling. Focal-adhesion pathway includes genes that involve cell-matrix adhesions, which play

essential roles in important biological processes including cell motility, cell proliferation, cell differentiation, and cell survival. These two pathways play important roles in the survival of HCC patients [25]. We also regress each selected image feature (y) over all the gene expression features (x) using LASSO penalization. Nine image features have R-squares  $> 0.10$  with the gene expression features (**Supplementary Table 1**). The one feature with the best fitted linear relationship is related to nuclei intensity (StDev\_Nuclei\_Intensity\_MeanIntensity\_MaskedHWWithoutOverlap) with R-square=0.19. This result shows that imaging features extracted using *CellProfiler* have modest correlations with the RNA-Seq gene expression features. Most of the image features can provide additional predictive values to prognosis, supporting the observed significant increase in C-index (**Figure 2**) and log-rank p-values (**Figure 4**), after adding RNA-Seq features to imaging features.

#### 4 CONCLUSIONS

Driven by the objective to build a uniform workframe to integrate multi-modal and multi-type data to predict patient survival, we extend Cox-nnet model, a neural-network based survival prediction method, on pathology imaging and transcriptomics data. Using TCGA HCC pathology images as the example, we demonstrate that Cox-nnet is more robust and accurate at predicting survival, compared to Cox-PH the standard method which was also the second-best method in the original RNA-Seq transcriptomic study [1]. Moreover, we propose a new two-stage complex Cox-nnet model to integrate imaging and RNA-Seq transcriptomic data, and showcase its superior accuracy on HCC patient survival prediction, compared to another neural network PAGE-Net. The two-stage Cox-nnet model combines the transformed, hidden node features from the first-stage of Cox-nnet models for imaging or gene expression RNA-Seq data respectively and uses these combined hidden features as the new inputs to train a second-stage Cox-nnet model.

Rather than using convolutional neural network (CNN) models that are more complex, such as PAGE-Net, we utilized a less complex but more biologically interpretable approach, where we extract

imaging features defined by the tool *CellProfiler*. These features are then used as input nodes in relatively simple, two-layer neural network models. Hidden features extracted from each Cox-nnet model can then be combined flexibly to build new Cox-nnet models. On the other hand, PAGE-Net uses a pretrained CNN for images and a gene-pathway layer to handle gene expression data. Despite great efforts, image features extracted by CNN in PAGE-Net are not easily interpretable, the model appears to be over-fit given the limited sample size, and requires very long training time. The significantly higher predictive performance of two stage Cox-nnet model argues for the advantages to use a relatively simple neural network model with input nodes of biological relevance, such as those extracted by imaging processing tools and gene expression input features.

Besides the interpretability of histopathology image features themselves, correlation analysis between top gene features and top image features identified genes known to be related to survival of HCC patients and/or morphology of the tissue, such as LINC01554, HOXD9, MUC6, and MMP7. LINC01554 is a long non-coding RNA that is down-regulated in HCC and its expression corresponded to good survival of HCC patients previously [26]. HOXD9 is a highly conserved transcription factor which was reported to promote the epithelial–mesenchymal transition [27] of HCC cells and associated with poor survival of HCC patients [28]. MUC6 is a mucin protein that participates in the remodeling of the ductal plate in the liver [29], which was also involved in the carcinogenesis of HCC [30]. MMP7, also known as matrilysin, is an enzyme that breaks down extracellular matrix by degrading macromolecules including casein, type I, II, IV, and V gelatins, fibronectin, and proteoglycan [31]. MMP7 participates in the remodeling of extracellular matrix [32] and impacts the morphology of liver tissue [33], which may explain its link to histopathology features. MMP7 expression was also associated with poor prognosis in patients with HCC [34].

In summary, we extend the previous Cox-nnet model to process pathological imaging data, and propose a new class of two-stage Cox-nnet neural network model that creatively addresses the general challenge of multi-modal data integration, for patient survival prediction. Using input imaging

features extracted from CellProfiler, Cox-nnet models are biologically interpretable. Some image features are also correlated with genes of known HCC relevance, enhancing their biological interpretability.

## **AUTHOR CONTRIBUTIONS**

LXG envisioned the project, supervised the study and wrote the majority of the manuscript. ZZ and ZJ performed modeling and data analysis, with help from BH and NH. MW and EYC assisted with the tumor section labeling of all the images.

## **CONFLICT OF INTEREST**

The authors declare no conflict of interest.

## **ACKNOWLEDGEMENTS**

We thank previous Garmire group members Dr. Fadhl Alakwaa, Dr. Olivier Poirion and Dr. Travers Ching for the discussions. LXG would like to thank the support by grants K01ES025434 awarded by NIEHS through funds provided by the trans-NIH Big Data to Knowledge (BD2K) initiative ([www.bd2k.nih.gov](http://www.bd2k.nih.gov)), R01 LM012373 and R01 LM012907 awarded by NLM, and R01 HD084633 awarded by NICHD to L.X. Garmire.

## **REFERENCES**

- [1] Ching T, Zhu X, Garmire LX. Cox-nnet: An artificial neural network method for prognosis prediction of high-throughput omics data. *PLoS Comput. Biol.*, vol. 14, no. 4, p. e1006076, Apr. 2018.
- [2] Ishwaran H, Lu M. Random Survival Forests. *Wiley StatsRef: Statistics Reference Online*. pp. 1–13, 2019.



- [3] Bin RD. Boosting in Cox regression: a comparison between the likelihood-based and the model-based approaches with focus on the R-packages CoxBoost and mboost. *Computational Statistics*, vol. 31, no. 2. pp. 513–531, 2016.
- [4] McQuin C, Goodman A, Chernyshev V, Kametsky L, Cimini BA, Karhohs KW, et al. CellProfiler 3.0: Next-generation image processing for biology. *PLoS Biol.*, vol. 16, no. 7, p. e2005970, Jul. 2018.
- [5] Chaudhary K, Poirion OB, Lu L, Garmire LX. Deep Learning-Based Multi-Omics Integration Robustly Predicts Survival in Liver Cancer. *Clin. Cancer Res.*, vol. 24, no. 6, pp. 1248–1259, Mar. 2018.
- [6] Chaudhary K, Poirion OB, Lu L, Huang S, Ching T, Garmire LX. Multi-modal meta-analysis of 1494 hepatocellular carcinoma samples reveals vast impacts of consensus driver genes on phenotypes. *Clin Cancer Res.* 2019;25(2):463-472. doi:10.1158/1078-0432.CCR-18-0088
- [7] Marinaccio C, Ribatti D. A simple method of image analysis to estimate CAM vascularization by APERIO ImageScope software. *Int. J. Dev. Biol.*, vol. 59, no. 4–6, pp. 217–219, 2015.
- [8] Yu KH, Zhang C, Berry GJ, Altman RB, Ré C, Rubin DL, et al. Predicting non-small cell lung cancer prognosis by fully automated microscopic pathology image features. *Nat. Commun.*, vol. 7, p. 12474, Aug. 2016.
- [9] Huang S, Yee C, Ching T, Yu H, Garmire LX. A novel model to combine clinical and pathway-based transcriptomic information for the prognosis prediction of breast cancer. *PLoS Comput. Biol.*, vol. 10, no. 9, p. e1003851, Sep. 2014.
- [10] Schulz-Streeck T, Ogutu JO, Piepho HP. Comparisons of single-stage and two-stage approaches to genomic selection. *Theor. Appl. Genet.*, vol. 126, no. 1, pp. 69–82, Jan. 2013.
- [11] Wei R, Vivo ID, Huang S, Zhu X, Risch H, Moore JH, et al. Meta-dimensional data integration identifies critical pathways for susceptibility, tumorigenesis and progression of endometrial cancer. *Oncotarget*, vol. 7, no. 34, pp. 55249–55263, Aug. 2016.
- [12] Pinu FR, Beale DJ, Paten AM, Kouremenos K, Swarup S, Schirra HJ, et al. Systems Biology

and Multi-Omics Integration: Viewpoints from the Metabolomics Research Community. *Metabolites*. 2019;9(4):76. Published 2019 Apr 18.

[13] Bengio Y, Boulanger-Lewandowski N, Pascanu R. Advances in optimizing recurrent networks. *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*; 2013: IEEE.

[14] Kametsky L, Jones TR, Fraser A, Bray MA, Logan DJ, Madden KL, et al. Improved structure, function and compatibility for CellProfiler: modular high-throughput image analysis software. *Bioinformatics* 27, 1179–1180 (2011).

[15] Hao J, Kosaraju SC, Tsaku NZ, Song DH, Kang M. PAGE-Net: Interpretable and Integrative Deep Learning for Survival Analysis Using Histopathological Images and Genomic Data. *Pac Symp Biocomput*. 2020;25:355-366.

[16] Yao J, Zhu X, Zhu F, Huang J. Deep Correlational Learning for Survival Prediction from Multi-modality Data. In: Descoteaux M., Maier-Hein L., Franz A., Jannin P., Collins D., Duchesne S, editors. *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2017*. MICCAI 2017. *Lecture Notes in Computer Science*, vol 10434. Springer, Cham

[17] Hao J, Kim Y, Mallavarapu T, Oh JH, Kang M. Cox-PASNet: Pathway-based Sparse Deep Neural Network for Survival Analysis. *2018 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, Madrid, Spain, 2018, pp. 381-386, doi: 10.1109/BIBM.2018.8621345.

[18] Hao J, Masum M, Oh JH, Kang M. Gene- and Pathway-Based Deep Neural Network for Multi-omics Data Integration to Predict Cancer Survival Outcomes. In: Cai Z, Skums P, Li M, editors. *Bioinformatics Research and Applications. ISBRA 2019*. *Lecture Notes in Computer Science*, vol 11490. Springer, Cham, doi:10.1007/978-3-030-20242-2\_10.

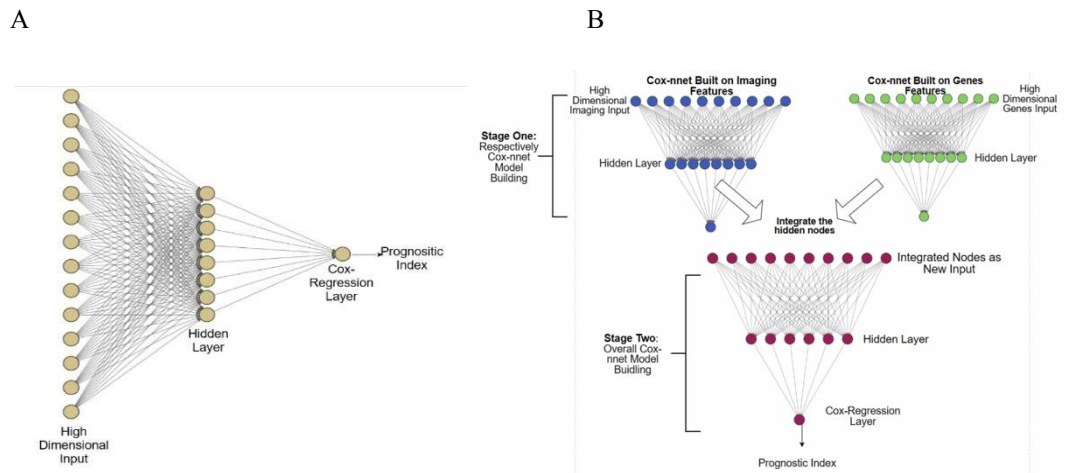
[19] Paul R, Hawkins SH, Hall LO, Goldgof DB, Gillies RJ. Combining deep neural network and traditional image features to improve survival prediction accuracy for lung cancer patients from diagnostic CT. *2016 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, Budapest, 2016, pp. 002570-002575, doi: 10.1109/SMC.2016.7844626.

- [20] Khalaf N, Ying J, Mittal S, Temple S, Kanwal F, Davila J, et al. Natural History of Untreated Hepatocellular Carcinoma in a US Cohort and the Role of Cancer Surveillance. *Clin Gastroenterol Hepatol.* 2017;15(2):273-281.e1. doi:10.1016/j.cgh.2016.07.033
- [21] Zhu Y, Qiu P, Ji Y. TCGA-Assembler: open-source software for retrieving and processing TCGA data. *Nat Methods* 11, 599–600 (2014). doi:10.1038/nmeth.2956
- [22] Huang S, Chong N, Lewis NE, Jia W, Xie G, Garmire LX. Novel personalized pathway-based metabolomics models reveal key metabolic pathways for breast cancer diagnosis. *Genome Med* 8, 34 (2016). doi:10.1186/s13073-016-0289-9
- [23] Fang X, Liu Y, Ren Z, Du Y, Huang Q, Garmire LX. Lilikoi V2.0: a deep-learning enabled, personalized pathway-based R package for diagnosis and prognosis predictions using metabolomics data. *bioRxiv*, 2020. doi:10.1101/2020.07.09.195677.
- [24] Chaudhary K, Huang S, Garmire LX. Multi-omics-based pan-cancer prognosis prediction using an ensemble of deep-learning and machine-learning models. *medRxiv*, 2019 doi:10.1101/19010082.
- [25] Ding X, Zhu Q, Zhang S, Guan L, Li T, Zhang L, et al. Precision medicine for hepatocellular carcinoma: driver mutations and targeted therapy. *Oncotarget.* 2017;8(33):55715-55730. Published 2017 Jun 6. doi:10.18632/oncotarget.18382
- [26] Ding Y, Sun Z, Zhang S, Chen Y, Zhou B, Li G, et al. Down-regulation of Long Non-coding RNA LINC01554 in Hepatocellular Cancer and its Clinical Significance. *J Cancer.* 2020;11(11):3369-3374. Published 2020 Mar 5. doi:10.7150/jca.40512
- [27] Lv X, Li L, Lv L, Qu X, Jin S, Li K, et al. HOXD9 promotes epithelial-mesenchymal transition and cancer metastasis by ZEB1 regulation in hepatocellular carcinoma. *J Exp Clin Cancer Res.* 2015;34:133. Published 2015 Oct 29. doi:10.1186/s13046-015-0245-3
- [28] Long J, Zhang L, Wan X, Lin J, Bai Y, Xu W, et al. A four-gene-based prognostic model predicts overall survival in patients with hepatocellular carcinoma. *J Cell Mol Med.* 2018;22(12):5928-5938. doi:10.1111/jcmm.13863

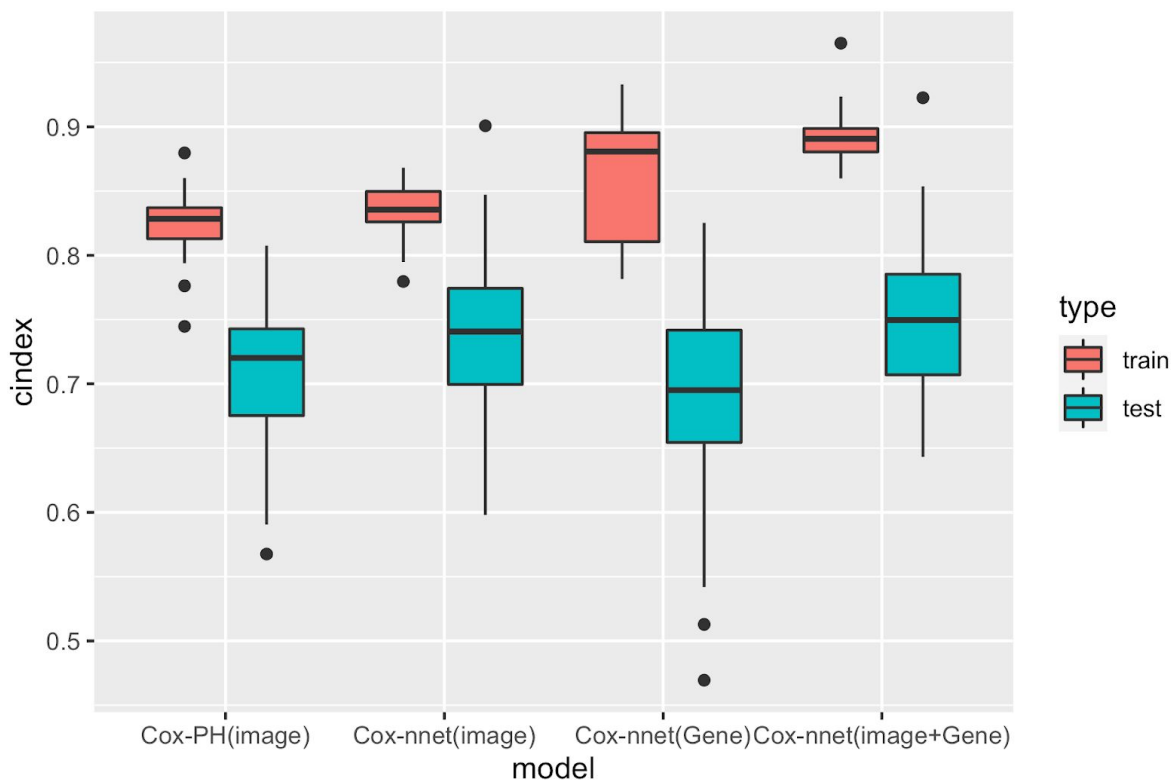
- [29] Terada T. Human fetal ductal plate revisited: II. MUC1, MUC5AC, and MUC6 are expressed in human fetal ductal plate and MUC1 is expressed also in remodeling ductal plate, remodeled ductal plate and mature bile ducts of human fetal livers. *Int J Clin Exp Pathol.* 2013;6(4):571-585.
- [30] Kasprzak A, Adamek A. Mucins: the Old, the New and the Promising Factors in Hepatobiliary Carcinogenesis. *Int J Mol Sci.* 2019;20(6):1288. Published 2019 Mar 14. doi:10.3390/ijms20061288
- [31] Yokoyama Y, Grünebach F, Schmidt SM, Heine A, Häntschel M, Stevanovic S, et al. Matrilysin (MMP-7) is a novel broadly expressed tumor antigen recognized by antigen-specific T cells. *Clin Cancer Res.* 2008;14(17):5503-5511. doi:10.1158/1078-0432.CCR-07-4041
- [32] Benyon RC, Arthur MJ. Extracellular matrix degradation and the role of hepatic stellate cells. *Semin Liver Dis.* 2001;21(3):373-384. doi:10.1055/s-2001-17552
- [33] Huang CC, Chuang JH, Chou MH, Wu CL, Chen CM, Wang CC, et al. Matrilysin (MMP-7) is a major matrix metalloproteinase upregulated in biliary atresia-associated liver fibrosis. *Mod Pathol.* 2005;18(7):941-950. doi:10.1038/modpathol.3800374
- [34] Rong W, Zhang Y, Yang L, Feng L, Wei B, Wu F, et al. Post-surgical resection prognostic value of combined OPN, MMP7, and PSG9 plasma biomarkers in hepatocellular carcinoma. *Front Med.* 2019;13(2):250-258. doi:10.1007/s11684-018-0632-1

## FIGURES AND LEGENDS

**Figure 1: the architectures of Cox-nnet model and two-stage Cox-nnet model:** A. The sketch of Cox-nnet model for prognosis prediction, based on a single data type. B. the architectures of two-stage Cox-nnet complex model for prognosis prediction, which integrates multiple data types (eg. pathology image and gene expression).

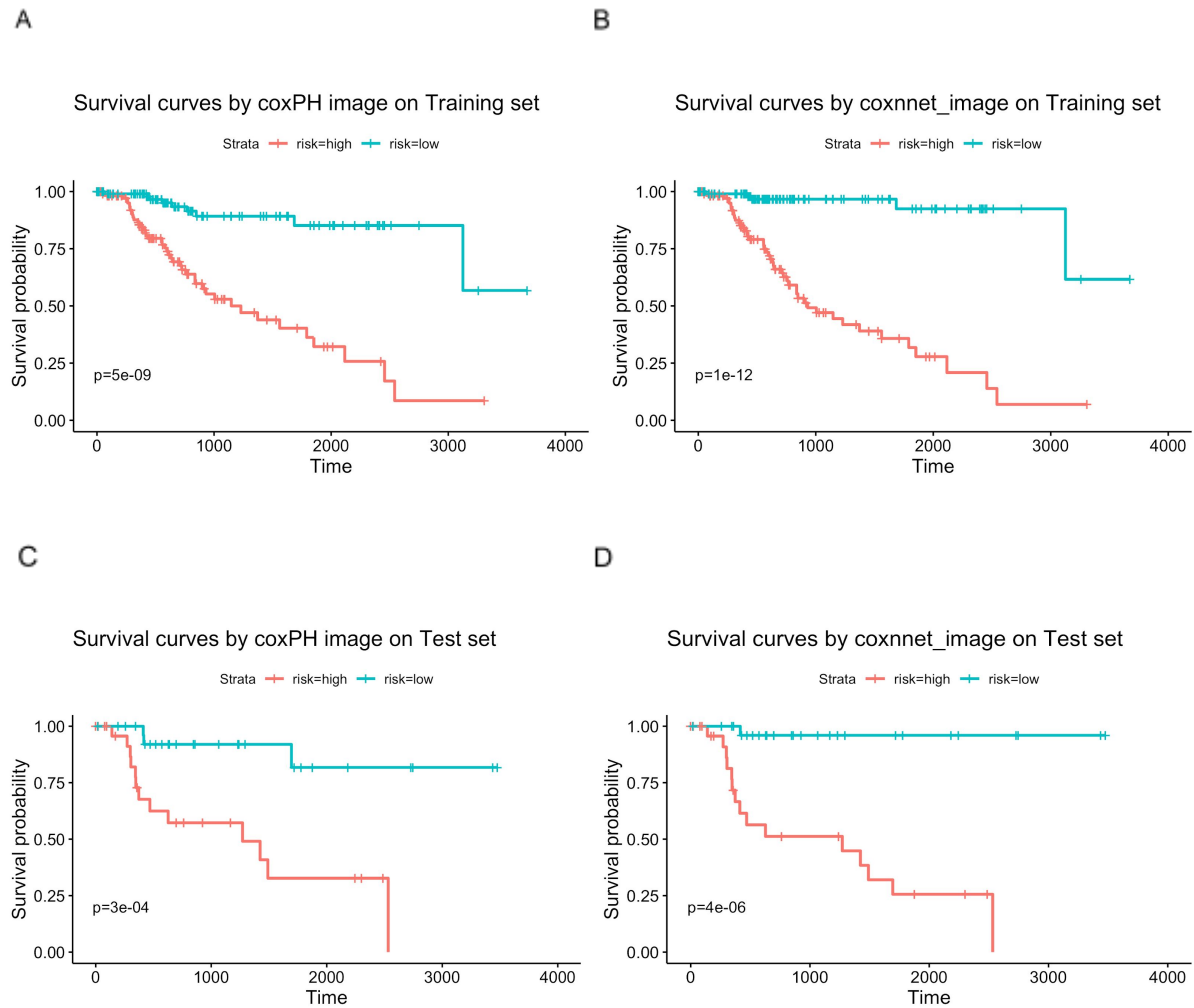


**Figure 2: Comparison of prognosis prediction with different models and data types.** The boxplots shown are on training (red) and testing (blue) data, on 20 repetitions.

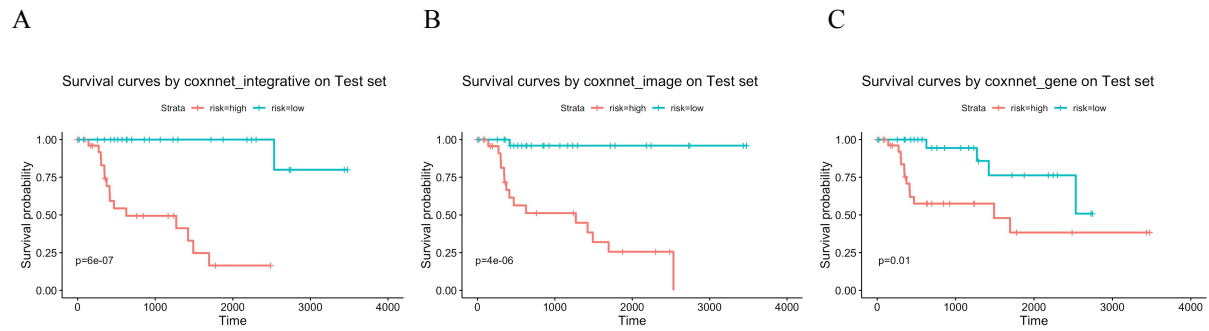


**Figure 3: Comparison of Kaplan-Meier survival curves resulting from Cox-PH and Cox-nnet models, based on pathological images.**

A. coxph image - training   B. Cox-nnet image - training   C. coxph image - test   D. Cox-nnet image - test

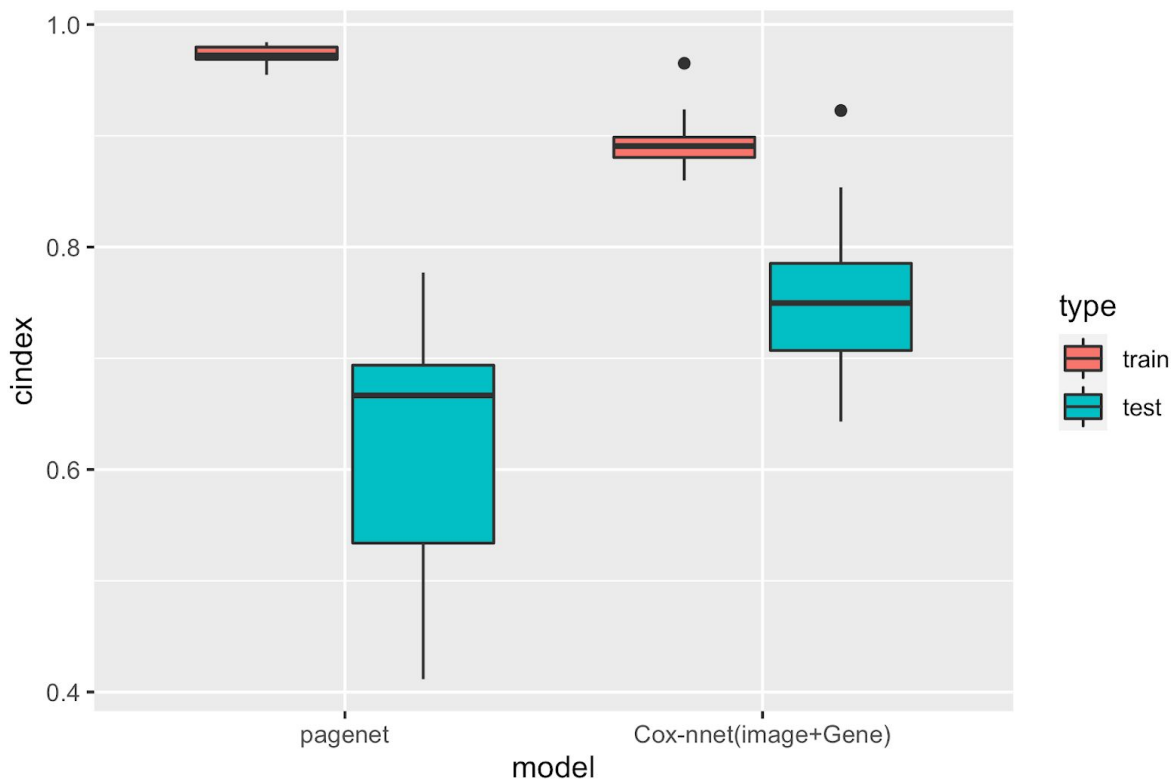


**Figure 4: Comparison of Kaplan-Meier survival curves on testing datasets, among two-stage Cox-nnet model combining pathological images and gene expression RNA-Seq data, the Cox-nnet model on image only, and the Cox-nnet model on gene expression RNA-Seq data only. (A) two-stage Cox-nnet model combining images and gene expression data. (B) Cox-nnet model on imaging data only. (C) Cox-nnet model on gene expression data only.**

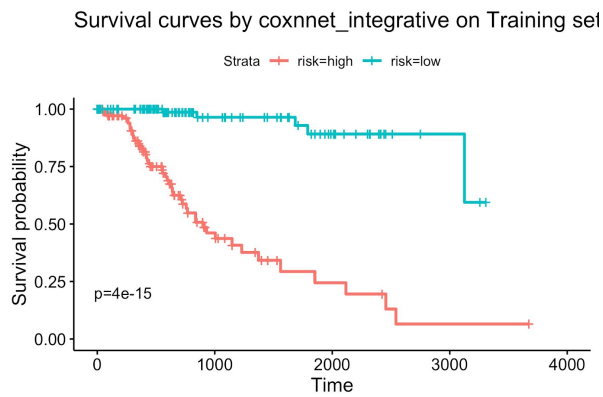


**Figure 5: Comparison of two-stage Cox-nnet and PAGE-Net, based on combined pathological images and gene expression.** (A) C-index of the two methods on training (red) and testing (blue) datasets, on 20 repetitions. (B-E) Kaplan-Meier survival curves resulting from the Cox-nnet (B, D) and PAGE-NET model (C, E) using training and testing datasets, respectively.

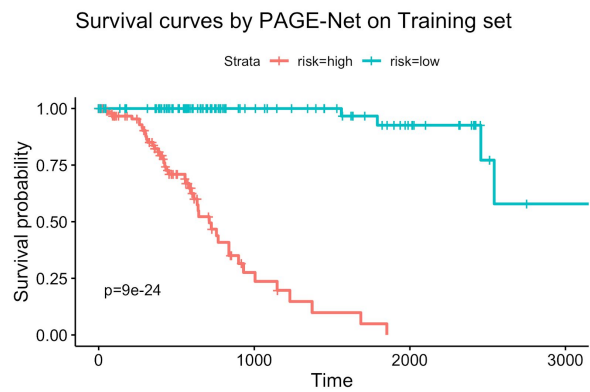
A.



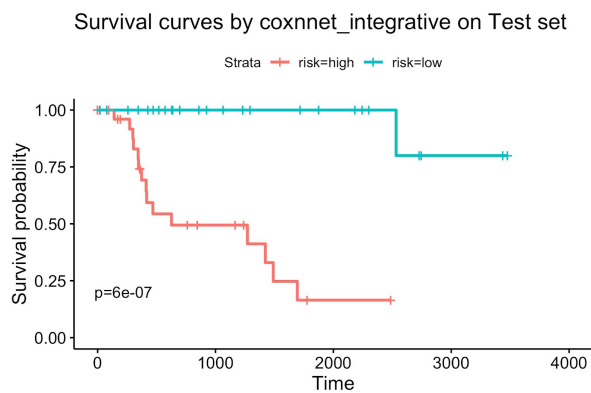
B.



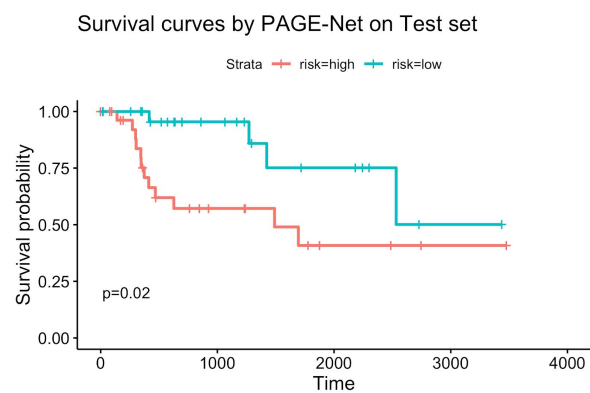
C.



D.

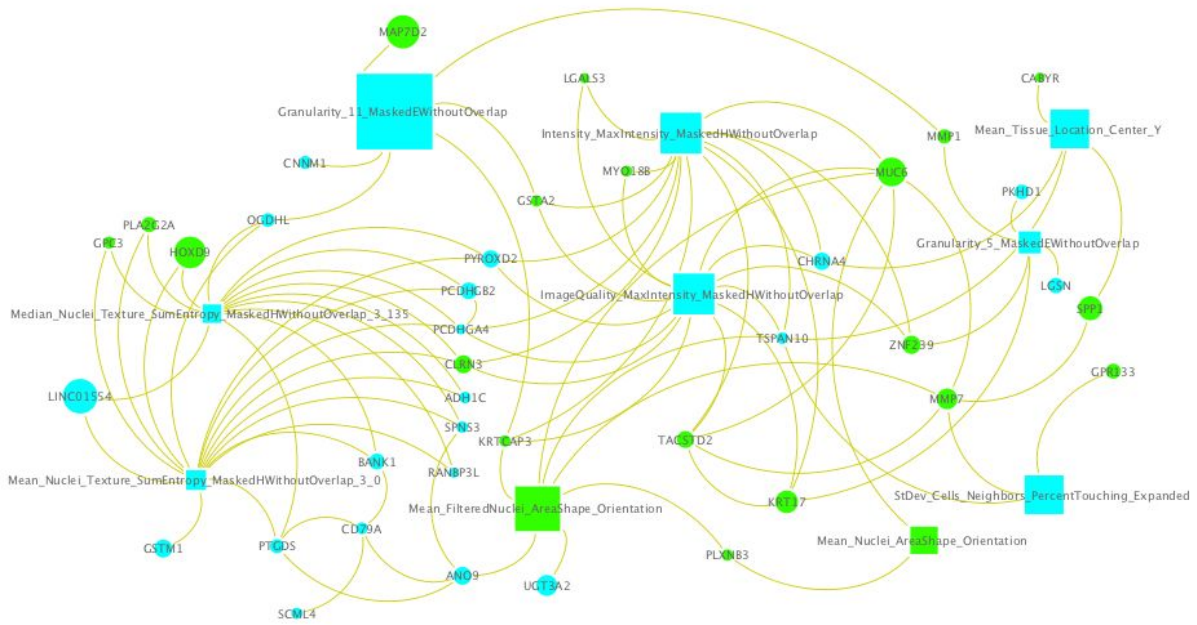


E.



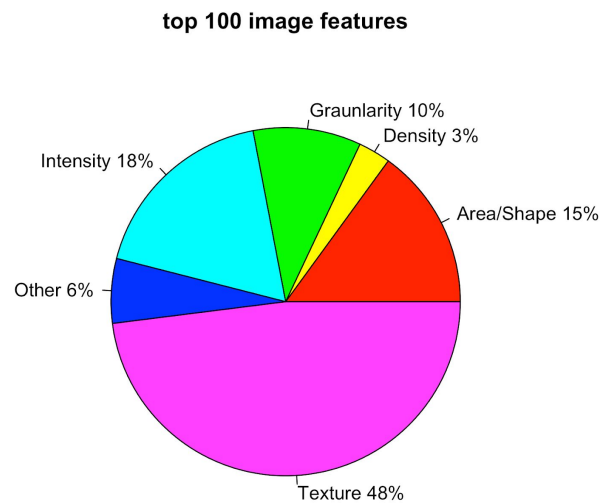
**Figure 6: Relationship between top imaging and gene features.** Rectangle nodes are image features and circle nodes are gene features. Node sizes are proportional to importance scores from Cox-nnet. Two gene nodes are connected only if their correlation is greater than 0.5; an image node and a gene node are connected only if their correlation is greater than 0.1. Green nodes represent features with positive coefficients (hazard ratio) in univariate Cox-PH regression, indicating worse prognosis. Blue nodes represent features with negative coefficients (hazard ratio) in univariate Cox-PH regression, indicating protection against bad prognosis.





## SUPPLEMENTARY MATERIALS

**Supplementary Figure 1: Categories of the top 100 most important image features in Cox-nnet.**



**Supplementary Table 1** : R-square correlations of top image features (  $R^2 > 0.10$  ), by regressing each of them on all gene features

<b>feature</b>	<b>R-square</b>
<b>StDev_Nuclei_Intensity_MeanIntensity_MaskedHWithoutOverlap</b>	<b>0.1938</b>
<b>Median_Tissue_Location_Center_Y</b>	<b>0.1525</b>
<b>Texture_SumAverage_MasedHWithoutOverlap_3_45</b>	<b>0.1476</b>
<b>Texture_SumAverage_MasedHWithoutOverlap_3_135</b>	<b>0.1473</b>
<b>StDev_Cells_Neighbors_PercentTouching_Expanded</b>	<b>0.1472</b>
<b>Median_Nuclei_Texture_SumEntropy_MaskedHWithoutOverlap_3_45</b>	<b>0.1343</b>
<b>Median_Nuclei_Texture_SumEntropy_MaskedHWithoutOverlap_3_135</b>	<b>0.1267</b>
<b>Mean_Nuclei_Texture_SumEntropy_MaskedHWithoutOverlap_3_0</b>	<b>0.1103</b>
<b>Mean_Nuclei_Texture_SumEntropy_MaskedHWithoutOverlap_3_90</b>	<b>0.1074</b>