



# **MGCT: Mutual-Guided Cross-Modality Transformer for Survival Outcome Prediction using Integrative Histopathology-Genomic Features**

Mingxin Liu<sup>1</sup>, Yunzan Liu<sup>1</sup>, Hui Cui<sup>2</sup>, Chunquan Li<sup>3,#</sup>, Jiquan Ma<sup>1,#</sup>

*Department of Computer Science and Technology, Heilongjiang University  
Department of Computer Science and Information Technology, La Trobe University  
The First Affiliated Hospital, Hengyang Medical School, University of South China*

# Background

- I. Cancer ranks the **leading cause of death** worldwide and has become one of **the five most common** diseases in China and developing or developed countries.
- II. In China, **55** people die of cancer in every **10 minutes**.
- III. There were an estimated **19,292,789** new cases **and 9,958,133** cancer deaths worldwide in 2020. (excluding nonmelanoma, skin cancer, and basal cell carcinoma)
- IV. In 2023, **1,958,310** new cancer cases and **609,820** cancer deaths are projected to occur in the United States.
- V. **Accurately diagnosing and prognosis** the cancer is of **paramount clinical importance**.

[1] C. Xia et al. Cancer statistics in China and United States, 2022: profiles, trends, and determinants. Chinese Medical Journal, 2022.

[2] F. Bray et al. The ever-increasing importance of cancer as a leading cause of premature death worldwide. Cancer, 2021.

[3] H. Sung et al. Global Cancer Statistics 2020: GLOBOCAN Estimates of Incidence and Mortality Worldwide for 36 Cancers in 185 Countries. CA: A Cancer Journal for Clinicians, 2021.

[4] R. L. Siegel et al. Cancer statistics, 2023. CA: A Cancer Journal for Clinicians, 2023.

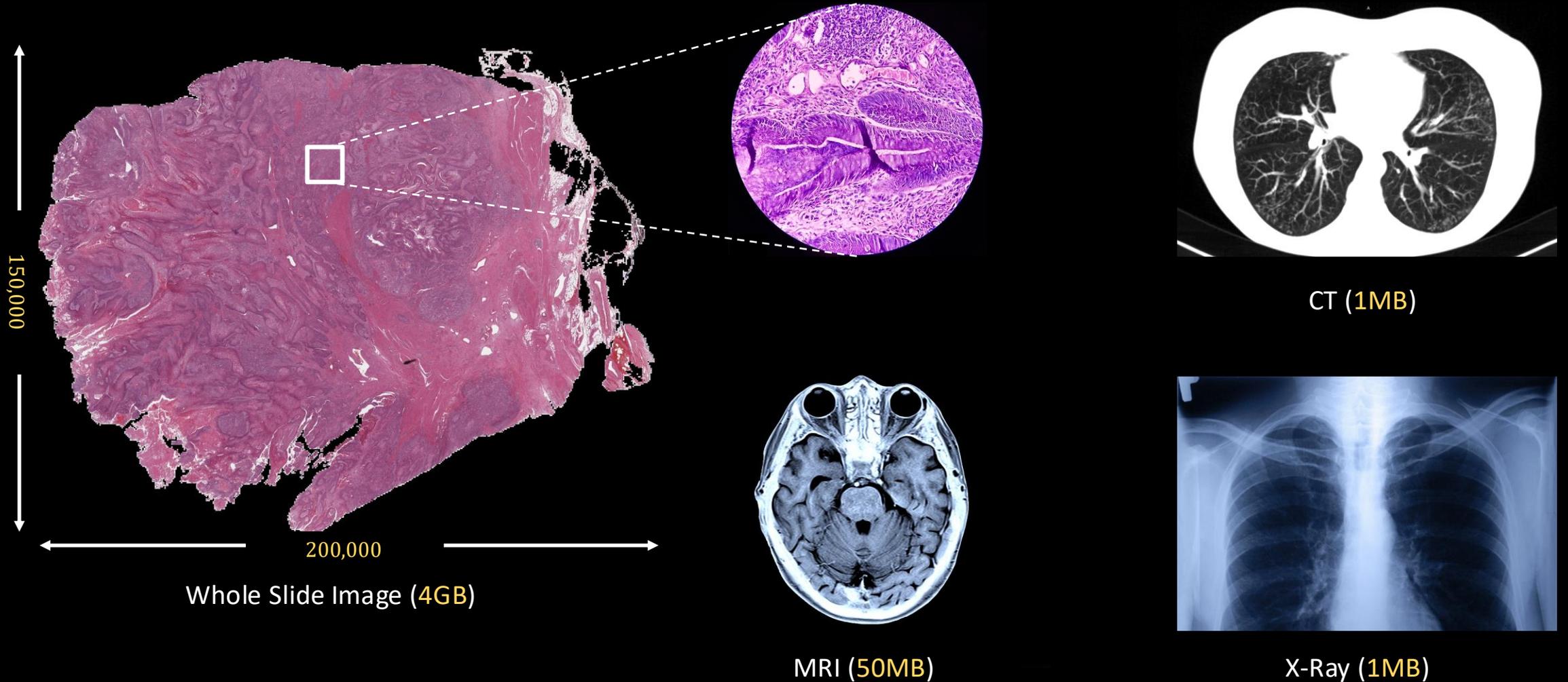
# Background

---

- I. Survival analysis is a crucial topic in clinical research, which aims to predict the time elapsed from a known origin to an event of interest, such as death, relapse of disease, and development of an adverse reaction.
- II. Traditionally, survival analysis relies on short term clinical indicators and long-term follow-up reports which are time-consuming and impractical in clinical applications.
- III. Recently, deep learning based medical image analysis is unfolding its infinity power.
- IV. While current deep learning-based survival outcome prediction techniques are single-modality, pathology or genomics alone, which inevitably reduce their potential to accurately predict patient prognosis.

# Challenges

— The enormous heterogeneity of gigapixel WSIs



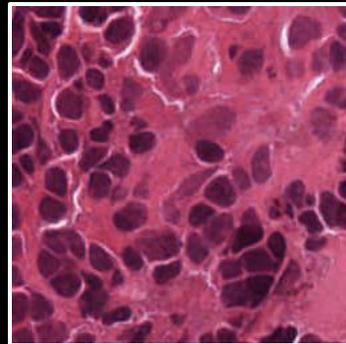
TCGA: The Cancer Genome Atlas. [\[GDC Data Portal\]](#)

# Challenges

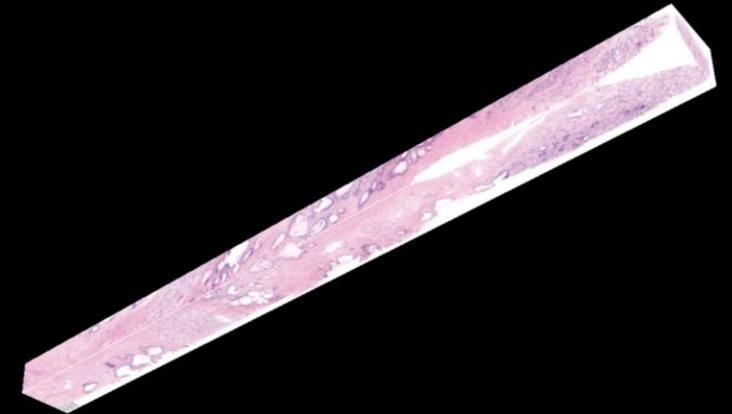
- The absence of spatially corresponding relationship

```
>chr1:90006571-90007309
CTGAAGGAAATAATTTTGCAAATAATTGAATATATTATAA
>chr1:230843894-230855664
CTCTAGGGTGGGACATGAGAAAGGACAGAAAAAGAAAGA
>chr1:239632290-239678180
GATCCTTTTCTAGTTGAGCTATTTCCCTTGAAAAGGGGGA
>chr1:182489205-182492965
AGAGCAGTCACTGTAATTTTTTTGACCTTTACATATGGGC
>chr1:230236424-230236654
GTACTGACACTGATTTATCCCTGTGTCTGGCTCTTCTCT
>chr1:51463798-51465258
CTACAATAAAAAACAAAATCACAAGTAAATAATAAGAGA
>chr1:68047367-68050547
AAGAAAGTGAGACGGTGAGAGATGGTAGTGGTAACCTACA
>chr1:35895269-35901471
GTGAGATTGCCAAGTAATGGCTGGGGAATAGGCATTGTAT
```

Genomics



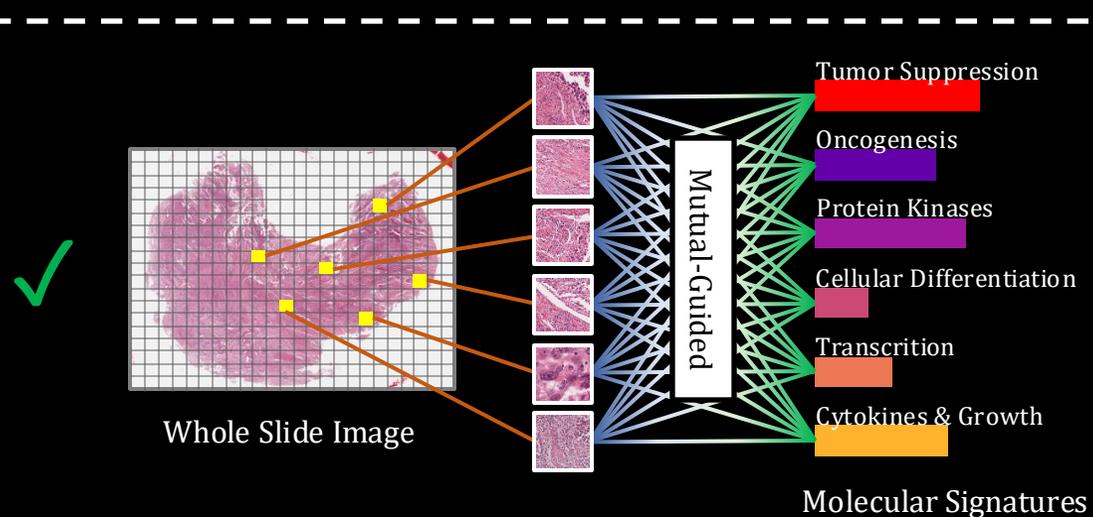
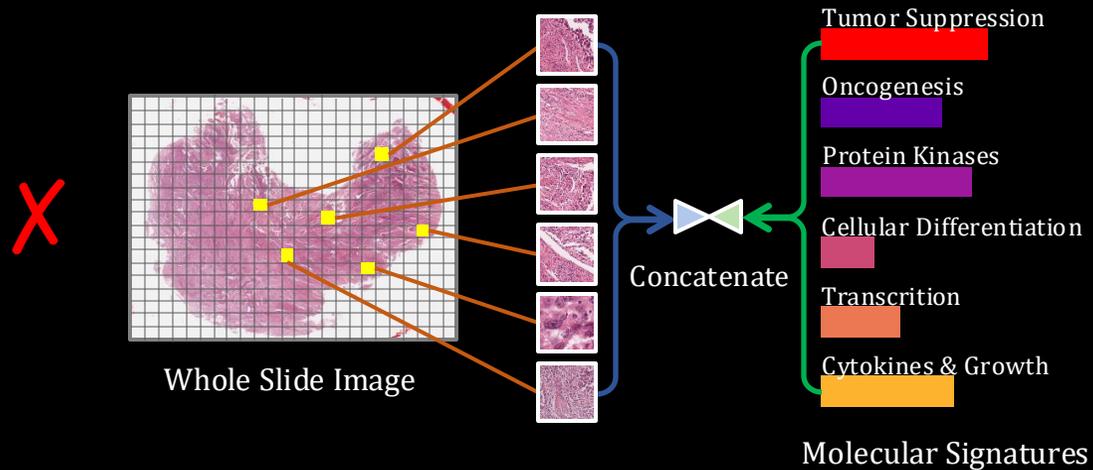
Pathological Image



3D Pathology

# Challenges

— SOTAs struggle to capture the explicit interactions



- Current SOTA methods are almost using **early, late, intermediate** multimodal feature fusion strategies which **cannot fully exploit the crucial interactions** between histopathology feature and genomic data.
  - Some guided-fusion based approaches are **solely using the genomic data as the guidance** to integrate multimodal pathomic features.
- ↓
- However, the gigapixel WSIs encompass **abundant crucial information** including **cell appearance, tumor microenvironment (TME), geometrical characteristics**.
  - Therefore, we designed a novel framework to capture the **genotype-phenotype interactions** by making these two modalities' data **guide each other mutually**.



# Methodology

## — Problem Formulation

- We denote the **input WSI** as  $\mathbb{X}_i$ , the feature vector of **genomic attributes with the WSI** as  $\mathbb{G}_i$ , the **overall survival time** (in months) as  $t_i \in \mathbb{R}^+$ , and the **right uncensorship status** (death observed) as  $c_i \in \{0,1\}$ .
- Therefore, we can represent the observations for all patient samples as a quadruple  $\{\mathbb{X}_i, \mathbb{G}_i, t_i, c_i\}_{i=1}^N$ .
- The main objective is to develop and optimize  $\mathcal{T}(\cdot)$  for integrating  $\mathbb{X}_i$  and  $\mathbb{G}_i$  to estimate the hazard function:

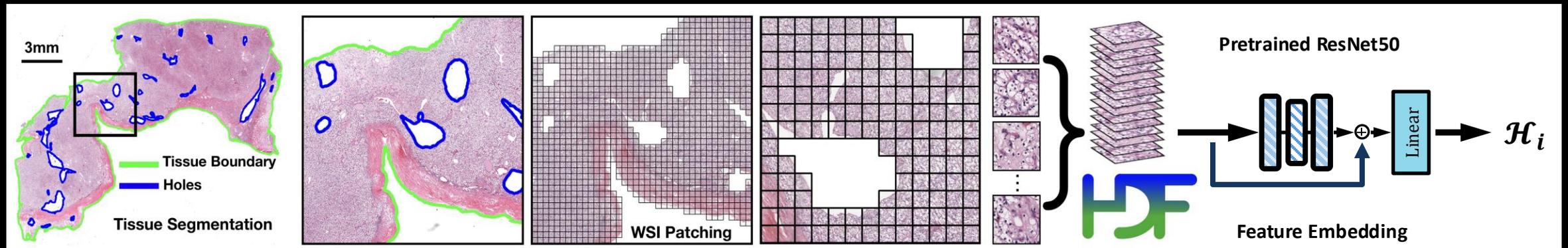
$$\tilde{t}_i = \mathcal{T}(\mathbb{X}_i, \mathbb{G}_i) = \emptyset \left( \xi \left( \rho([f(x_1), f(x_2), \dots, f(x_{N_i})]), \mathbb{G}_i \right) \right)$$

- $f(\cdot)$  is an instance-level encoder that processes features for each instance independently
- $\rho(\cdot)$  is the method for multimodal pathomic features integration
- $\xi(\cdot)$  is a **permutation-invariant** instance aggregator which aggregate and pools the features to a bag-level embedding
- $\emptyset(\cdot)$  is a bag-level classifier to make final survival outcome predictions

# Methodology

## — Histopathology Feature Extraction

- For input WSI  $\mathbb{X}_i$ , CLAM repository is employed for **automated tissue segmentation**.
- Then we extract  $256 \times 256$  image patches  $\{x_k\}_{k=1}^{N_i}$  **without spatial overlapping** at the **20 $\times$**  magnification.
- We further utilize an **ImageNet-pretrained** ResNet-50 to generate a 1024-dim feature embedding  $\mathbf{h}_k \in \mathbb{R}^{1024}$ .
- Finally, we assemble the feature embeddings into a **WSI-level bag representation**  $\mathcal{H}_i \in \mathbb{R}^{1024 \times N_i}$ .

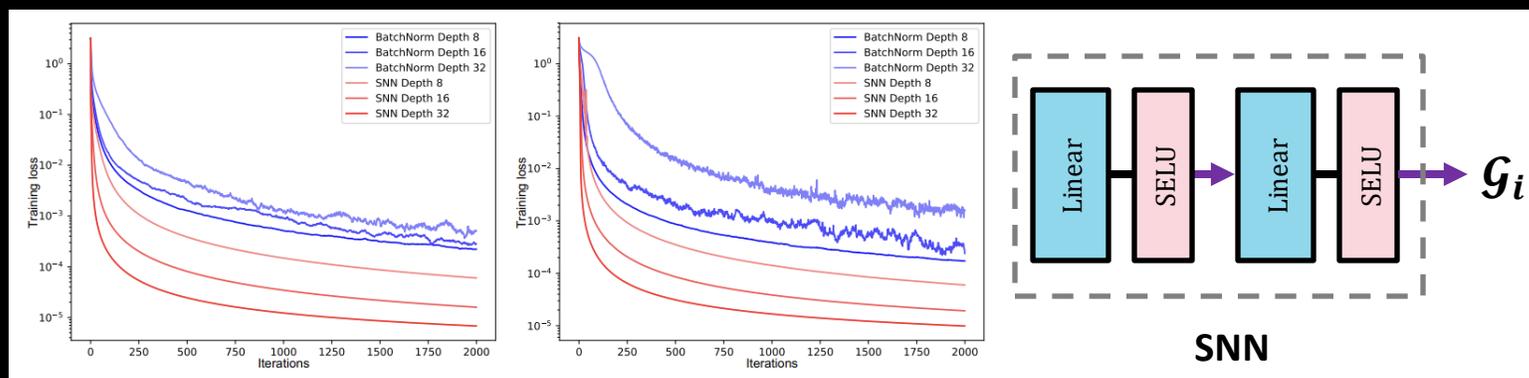


MY Lu et al. Data-efficient and weakly supervised computational pathology on whole-slide images. Nature biomedical engineering, 2021.

# Methodology

## — Genomic Feature Embedding

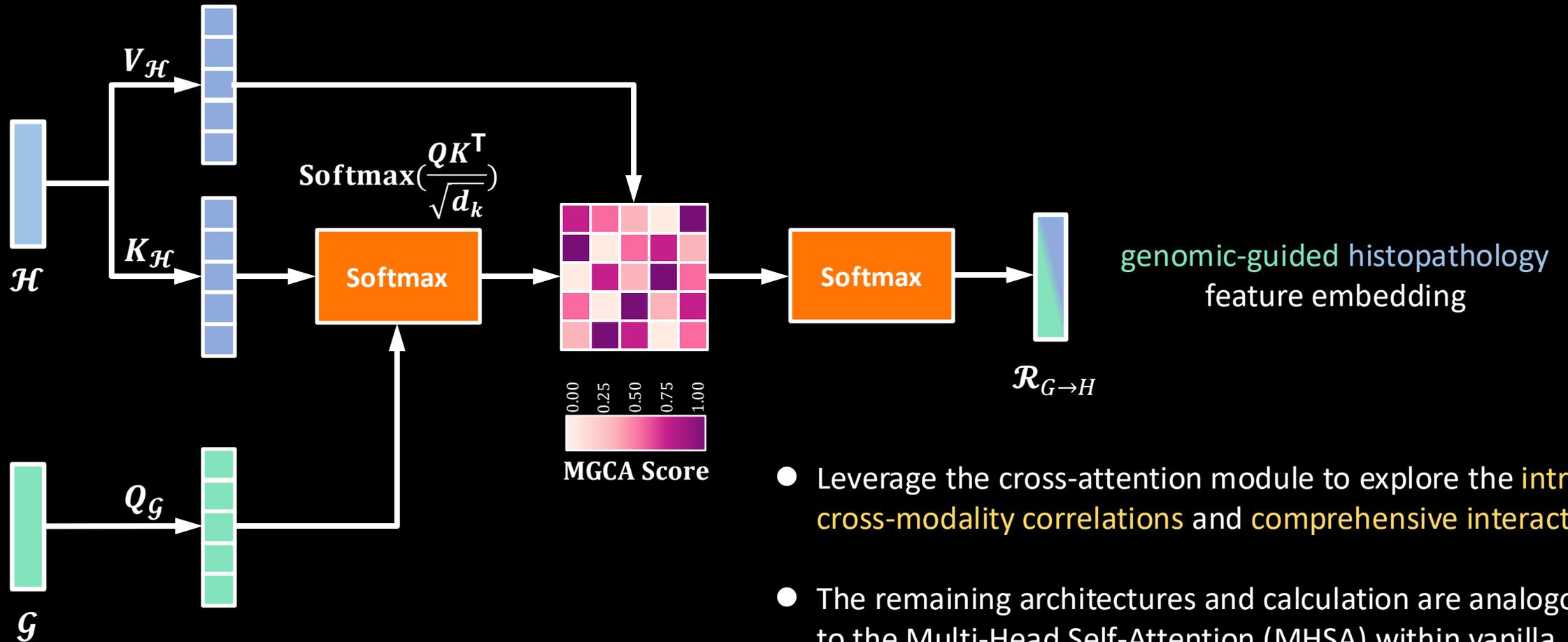
- We select **transcript abundance (bulk RNA-Seq)**, **gene mutation status**, **copy number variation** as the input genomics.
- These  $1 \times 1$  measurements exhibit a **high-dimensional low-sample (HDLSS)** nature which leads to **overfitting** problem.
- Therefore, we leverage the **Self-Normalizing Neural Network (SNN)** to formulate the genomic feature embedding.
- We further aggregate and structure the genomic embeddings based on **S** related **biological functional impacts**.
- Finally, we can generate the **bag-level genomic feature embedding** as  $\mathcal{G}_i \in \mathbb{R}^{1024 \times S}$ .



G. Klambauer et al. Self-normalizing neural networks. NeurIPS, 2017.

# Methodology

## — Mutual-Guided Cross-Modality Attention



- Leverage the cross-attention module to explore the **intrinsic cross-modality correlations** and **comprehensive interactions**.
- The remaining architectures and calculation are analogous to the Multi-Head Self-Attention (MHSA) within vanilla transformer encoder layer.

A. Vaswani et al. Attention is all you need. NeurIPS, 2017.

# Methodology

## — Mutual-Guided Cross-Modality Transformer

- The procedure for MGCT layer calculation:

$$\begin{aligned} \mathbf{MGCA}(\mathcal{G}_i, \mathcal{H}_i, \mathcal{H}_i) &= \mathbf{Softmax}\left(\frac{\mathbf{Q} \cdot \mathbf{K}^\top}{\sqrt{d_k}}\right) \\ &= \mathbf{Softmax}\left(\frac{\mathbf{W}_q \cdot \mathcal{G}_i \cdot \mathcal{H}_i^\top \cdot \mathbf{W}_k^\top}{\sqrt{d_k}}\right) \cdot \mathbf{W}_v \cdot \mathcal{H}_i \rightarrow \mathcal{R}_{G \rightarrow H} \end{aligned}$$

$$\begin{aligned} \mathcal{R}'_{G \rightarrow H} &= \mathbf{AttnPool}\left(\sum_{i=1}^N \alpha_i\right) \cdot \mathcal{R}_{G \rightarrow H} \quad \text{where} \\ \alpha_i &= \frac{\exp\{\mathbf{W}(\mathbf{tanh}(V \cdot \mathcal{R}_i^\top) \odot \mathbf{sigm}(U \cdot \mathcal{R}_i^\top))\}}{\sum_{j=1}^N \exp\{\mathbf{W}(\mathbf{tanh}(V \cdot \mathcal{R}_j^\top) \odot \mathbf{sigm}(U \cdot \mathcal{R}_j^\top))\}} \end{aligned}$$

$$\mathcal{R}''_{G \rightarrow H} = \zeta(\mathbf{MLP}(\mathcal{R}'_{G \rightarrow H})\mathbf{W}_{\mathbf{MLP}}) \cdot \mathbf{W}_\zeta$$

$\zeta(\cdot)$  is a permutation-invariant instance aggregator

---

**Algorithm 1:** The proposed MGCT framework

---

**Input:**

- I. WSI bag embedding  $\mathcal{H}_i \in \mathbb{R}^{1024 \times N_i}$ .
- II. Genomic bag embedding  $\mathcal{G}_i \in \mathbb{R}^{1024 \times S}$ .
- III. # MGCT layers in two multimodal feature integration stages,  $S_1$  and  $S_2$ .

1: **for**  $s_1 = 1$  to  $S_1$  **do**

- 2:  $\mathcal{R}''_{G \rightarrow H} \leftarrow \mathbf{MGCT-Layer}(\mathcal{G}_i, \mathcal{H}_i, \mathcal{H}_i)$
- $\mathcal{R}''_{H \rightarrow G} \leftarrow \mathbf{MGCT-Layer}(\mathcal{H}_i, \mathcal{G}_i, \mathcal{G}_i)$

3: **end for**

- 4:  $\mathcal{R}_{F_1} \leftarrow \mathbf{Concatenate}(\mathcal{R}''_{G \rightarrow H}, \mathcal{R}''_{H \rightarrow G})$

5: **for**  $s_2 = 1$  to  $S_2$  **do**

- 6:  $\mathcal{R}''_{F_1 \rightarrow H} \leftarrow \mathbf{MGCT-Layer}(\mathcal{R}_{F_1}, \mathcal{H}_i, \mathcal{H}_i)$
- $\mathcal{R}''_{H \rightarrow F_1} \leftarrow \mathbf{MGCT-Layer}(\mathcal{H}_i, \mathcal{R}_{F_1}, \mathcal{R}_{F_1})$

7: **end for**

- 8:  $\mathcal{R}_{\text{Final}} \leftarrow \mathbf{Concatenate}(\mathcal{R}''_{F_1 \rightarrow H}, \mathcal{R}''_{H \rightarrow F_1})$

**Return** final multimodal feature embedding  $\mathcal{R}_{\text{Final}}$

---

# Experiments

## — Datasets

- Five benchmarks were used for model evaluation
- BLCA: Bladder Urothelial Carcinoma
- BRCA: Breast Invasive Carcinoma
- LUAD: Lung Adenocarcinoma.
- GBMLGG: Glioblastoma Multiforme & Brain Lower Grade Glioma
- UCEC: Uterine Corpus Endometrial Carcinoma

Cancer	# Cases	# WSIs	# Patches	Censorship
BLCA	373	437	7,648,953	0.547
BRCA	957	1,023	12,306,155	0.860
LUAD	453	516	6,717,757	0.651
GBMLGG	569	1,042	12,742,037	0.766
UCEC	480	539	9,136,545	0.844
<b>Overall</b>	<b>2,832</b>	<b>3,557</b>	<b>48,551,447</b>	<b>0.734</b>

Genomic profiles are grouped by:

- Tumor Suppression
- Oncogenesis
- Protein Kinases
- Cellular Differentiation
- Transcription
- Cytokines and Growth

View Gene Families

The following table provides a functional overview of the MSigDB gene sets by categorizing their genes into a small number of carefully chosen "gene families". To categorize the genes in a gene set, use the gene set page or the Investigate Gene Sets page.

Click on a gene family or gene family intersection to retrieve annotations for those genes.

	cytokines and growth factors	transcription factors	homeodomain proteins	cell differentiation markers	protein kinases	translocated cancer genes	oncogenes	tumor suppressors
tumor suppressors	1	14	2	3	6	1	0	82
oncogenes							328	
translocated cancer genes	APC ASXL1 ATM BLM BMPRIA BRCA1 BRCA2 BRIP1 BUB1B CBLB							
protein kinases	CBLC CDC73 CDH1 CDKN2A CDKN2C CHEK2 CYLD DDB2 DICER1 EP300							
cell differentiation markers	ERCC2 ERCC3 ERCC4 ERCC5 EXT1 EXT2 FAM123B FANCA FANCC FANCD2							
homeodomain proteins	FANCE FANCF FANCG FAS FBXW7 FH GATA3 HNF1A KDM5C KDM6A							
transcription factors	KLF6 MAP2K4 MEN1 MLH1 MSH2 MSH6 MUTYH NBN NF1 NF2							
cytokines and growth factors	PALB2 PHOX2B PIK3R1 PMS1 PMS2 PRF1 PTCH1 PTEN RB1 RECQL4							
transcription factors	SBDS SDHAF2 SDHB SDHC SDHD SETD2 SMAD4 SMARCA4 SMARCB1 SOCS1							
transcription factors	STK11 SUFU TET2 TNFAIP3 TP53 TSC1 TSC2 VHL WRN WT1							
transcription factors	XPA XPC							
transcription factors	0	1537						
cytokines and growth factors	452							

A. Subramanian et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. PNAS, 2005.

# Experiments

---

## — Evaluation Metrics

- 5-fold Monte Carlo cross-validation for each cancer type was used for model evaluation.
- Concordance index (C-index) values were employed to measure the predictive ability of the model.
- Kaplan-Meier curves (KM curve) were leveraged to visually represent the quality of patient stratification.
- Log-rank test was introduced to determine the statistical significance of patient stratification.

## — Implementation Details

- MGCT is trained on a workstation equipped with an NVIDIA Quadro GV100 GPU for 20 epochs (about 7.5 hours).
- Adam optimization with learning rate of  $2e-4$  and weight decay of  $1e-5$ .
- Batch size is 1 (due to samples having varying bag sizes) and 32 gradient accumulation steps.
- Our related models and scripts will be publicly made available ASAP at <https://github.com/lmxmercy/MGCT>.

# Experiments

---

## — Baselines

- Unimodal Baselines:

- a) Genomic Only: MLP, SNN, DeepSurv, CoxRegression

- b) Pathology Only: Deep Sets, Attention MIL, CLAM, DeepAttnMISL, Patch-GCN

- Multimodal Baselines:

- c) Enhanced MILs with concatenation and bilinear pooling as multimodal baselines

- d) Current State-of-the-Art methods: PORPOISE, MCAT

# Experiments

## — Concordance index Comparison

	Methods	BLCA	BRCA	LUAD	GBMLGG	UCEC	Overall
Genomic	SNN	$0.541 \pm 0.016$	$0.466 \pm 0.058$	$0.539 \pm 0.069$	$0.598 \pm 0.054$	$0.493 \pm 0.096$	0.527
	DeepSurv	$0.567 \pm 0.049$	$0.598 \pm 0.054$	$0.608 \pm 0.026$	$0.810 \pm 0.020$	$0.577 \pm 0.058$	0.632
	CoxRegression	$0.591 \pm 0.041$	$0.568 \pm 0.077$	$0.574 \pm 0.042$	$0.705 \pm 0.014$	$0.464 \pm 0.099$	0.580
Pathology	Deep Sets	$0.500 \pm 0.000$	$0.500 \pm 0.000$	$0.496 \pm 0.008$	$0.498 \pm 0.014$	$0.500 \pm 0.000$	0.499
	CLAM	$0.565 \pm 0.027$	$0.578 \pm 0.032$	$0.582 \pm 0.072$	$0.776 \pm 0.034$	$0.609 \pm 0.082$	0.622
	DeepAttnMISL	$0.504 \pm 0.042$	$0.524 \pm 0.043$	$0.548 \pm 0.050$	$0.734 \pm 0.029$	$0.597 \pm 0.059$	0.581
	Patch-GCN	$0.560 \pm 0.034$	$0.580 \pm 0.025$	$0.585 \pm 0.012$	$0.824 \pm 0.024$	$0.629 \pm 0.052$	0.636
Multimodal	Attention MIL (Concat)	$0.605 \pm 0.045$	$0.551 \pm 0.077$	$0.563 \pm 0.050$	$0.816 \pm 0.011$	$0.614 \pm 0.052$	0.630
	DeepAttnMISL (Concat)	$0.611 \pm 0.049$	$0.545 \pm 0.071$	$0.595 \pm 0.061$	$0.805 \pm 0.014$	$0.615 \pm 0.020$	0.634
	PORPOISE	$0.613 \pm 0.021$	$0.563 \pm 0.056$	<b><math>0.621 \pm 0.045</math></b>	$0.818 \pm 0.011$	$0.622 \pm 0.061$	0.647
	MCAT	$0.624 \pm 0.034$	$0.580 \pm 0.069$	$0.620 \pm 0.032$	$0.817 \pm 0.021$	$0.622 \pm 0.019$	0.653
	<b>MGCT (Ours)</b>	<b><math>0.640 \pm 0.039</math></b>	<b><math>0.608 \pm 0.026</math></b>	$0.596 \pm 0.078$	<b><math>0.827 \pm 0.024</math></b>	<b><math>0.645 \pm 0.039</math></b>	<b>0.663</b>

# Experiments

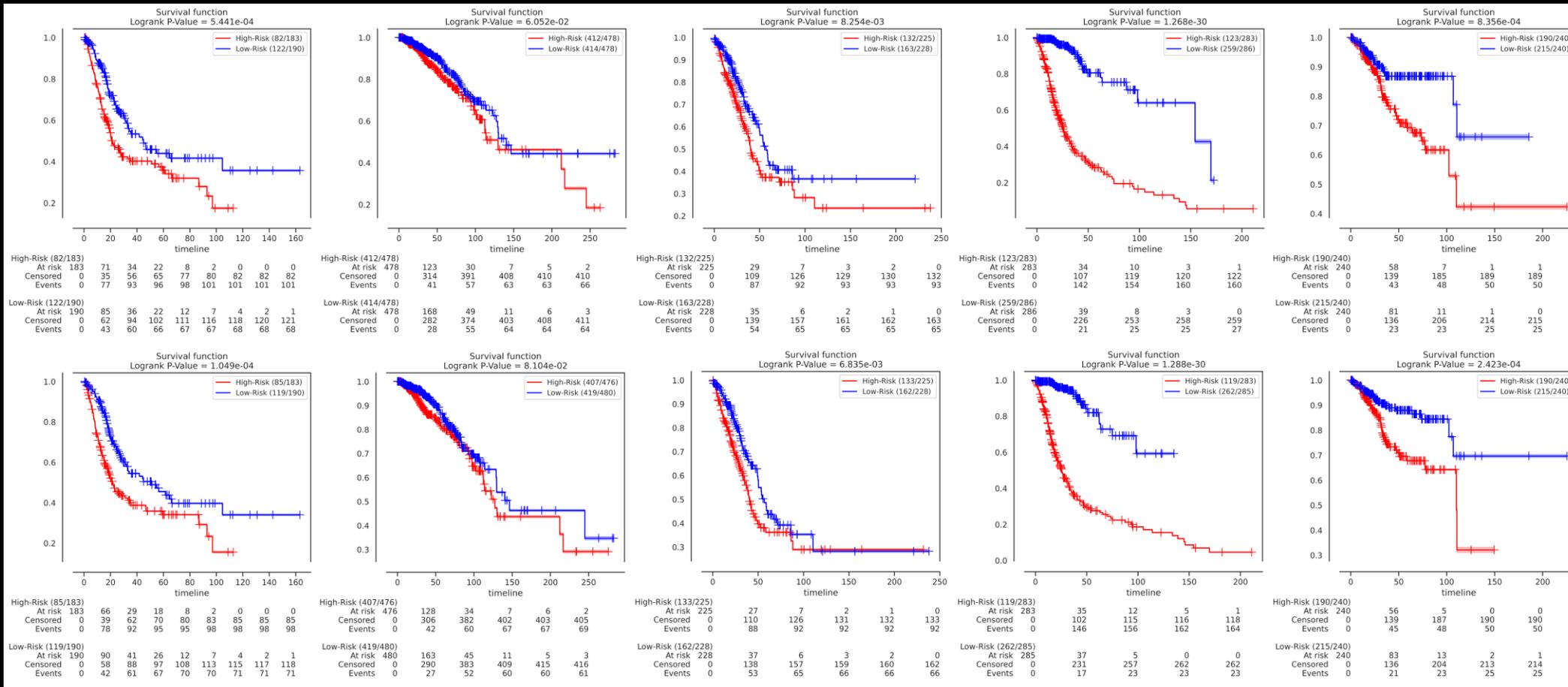
## — Patient Stratification: Kaplan-Meier Survival Curves

MCAT

MGCT (Ours)

Cumulative Proportion Surviving

Cumulative Proportion Surviving



# Experiments

— Ablation study on designed components

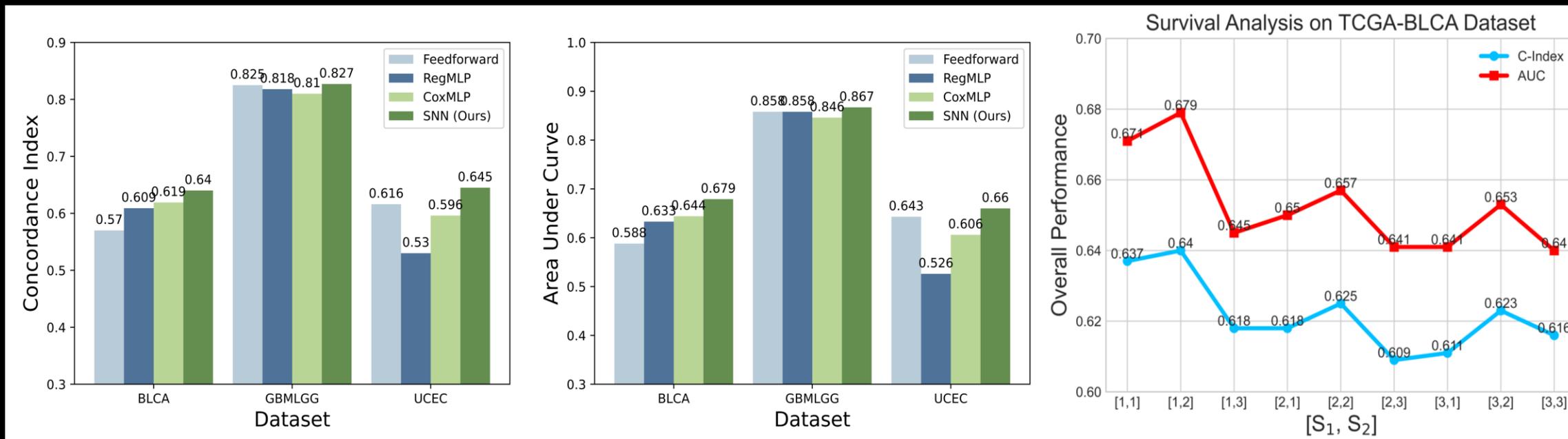
- Test on TCGA-BLCA and TCGA-UCEC two benchmarks.
- Deep Fusion: stack two parallel MGCT layers in depth.
- MGCA: mutual-guided cross-modality attention.
- GAP: gated-attention pooling operation in MGCT layer.
- Feedforward: position-wise feed-forward network in MGCT layer.

Model	Designs in MGCT				TCGA-BLCA		TCGA-UCEC	
	Deep Fusion	MGCA	GAP	Feedforward	C-index $\uparrow$	AUC $\uparrow$	C-index $\uparrow$	AUC $\uparrow$
A					$0.499 \pm 0.002$	$0.499 \pm 0.002$	$0.499 \pm 0.002$	$0.499 \pm 0.002$
B	✓				$0.535 \pm 0.038$	$0.532 \pm 0.045$	$0.541 \pm 0.063$	$0.558 \pm 0.034$
C	✓	✓			$0.590 \pm 0.045$	$0.621 \pm 0.072$	$0.608 \pm 0.062$	$0.627 \pm 0.071$
D	✓	✓	✓		$0.601 \pm 0.047$	$0.621 \pm 0.072$	$0.608 \pm 0.062$	$0.627 \pm 0.071$
E	✓	✓	✓	✓	<b><math>0.640 \pm 0.039</math></b>	<b><math>0.679 \pm 0.039</math></b>	<b><math>0.645 \pm 0.039</math></b>	<b><math>0.660 \pm 0.039</math></b>

# Experiments

— Ablation study on genomic feature embedding method & #MGCT layers

- Test on TCGA-BLCA, TCGA-GBMLGG, and TCGA-UCEC three benchmarks.



# References

- Xia C, Dong X, Li H, et al. Cancer statistics in China and United States, 2022: profiles, trends, and determinants. Chinese medical journal, 2022, 135(05): 584-590.
- Bray F, Laversanne M, Weiderpass E, et al. The ever-increasing importance of cancer as a leading cause of premature death worldwide. Cancer, 2021, 127(16): 3029-3030.
- Sung H, Ferlay J, Siegel R L, et al. Global cancer statistics 2020: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. CA: a cancer journal for clinicians, 2021, 71(3): 209-249.
- Siegel R L, Miller K D, Wagle N S, et al. Cancer statistics, 2023. CA: a cancer journal for clinicians, 2023, 73(1): 17-48.
- Zhou F, Chen H. Cross-Modal Translation and Alignment for Survival Analysis. Proceedings of the IEEE/CVF International Conference on Computer Vision. 2023: 21485-21494.
- Song A H, Williams M, Williamson D F K, et al. Weakly Supervised AI for Efficient Analysis of 3D Pathology Samples. arXiv preprint arXiv:2307.14907, 2023.
- Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need. Advances in neural information processing systems, 2017, 30.
- Subramanian A, Tamayo P, Mootha V K, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. Proceedings of the National Academy of Sciences, 2005, 102(43): 15545-15550.
- Klambauer G, Unterthiner T, Mayr A, et al. Self-normalizing neural networks. Advances in neural information processing systems, 2017, 30.

# References

- Katzman J L, Shaham U, Cloninger A, et al. DeepSurv: personalized treatment recommender system using a Cox proportional hazards deep neural network. *BMC medical research methodology*, 2018, 18(1): 1-12.
- Kvamme H, Borgan Ø, Scheel I. Time-to-event prediction with neural networks and Cox regression. *arXiv preprint arXiv:1907.00825*, 2019.
- Ilse M, Tomczak J, Welling M. Attention-based deep multiple instance learning. *International conference on machine learning*. PMLR, 2018: 2127-2136.
- Lu M Y, Williamson D F K, Chen T Y, et al. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature biomedical engineering*, 2021, 5(6): 555-570.
- Yao J, Zhu X, Jonnagaddala J, et al. Whole slide images based cancer survival prediction using attention guided deep multiple instance learning networks. *Medical Image Analysis*, 2020, 65: 101789.
- Chen R J, Lu M Y, Shaban M, et al. Whole slide images are 2d point clouds: Context-aware survival prediction using patch-based graph convolutional networks. *Medical Image Computing and Computer Assisted Intervention—MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part VIII 24*. Springer International Publishing, 2021: 339-349.
- Chen R J, Lu M Y, Williamson D F K, et al. Pan-cancer integrative histology-genomic analysis via multimodal deep learning. *Cancer Cell*, 2022, 40(8): 865-878. e6.
- Chen R J, Lu M Y, Weng W H, et al. Multimodal co-attention transformer for survival prediction in gigapixel whole slide images. *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021: 4015-4025.



*Thank You!*

*E-mail: mxliu.mercy@gmail.com*

22