

# Dynamic Entity-Masked Graph Diffusion Model for histopathological image Representation Learning

Zhenfeng Zhuang<sup>1\*</sup>, Min Cen<sup>2\*</sup>, Yanfeng Li<sup>1\*</sup>, Fangyu Zhou<sup>1</sup>, Lequan Yu<sup>3</sup>, Baptiste Magnier<sup>4,5</sup>, Liansheng Wang<sup>1†</sup>

<sup>1</sup>Department of Computer Science at School of Informatics, Xiamen University, Xiamen, China

<sup>2</sup>School of Artificial Intelligence and Data Science, University of Science and Technology of China, Hefei, Anhui, China

<sup>3</sup>Department of Statistics and Actuarial Science, The University of Hong Kong, Hong Kong SAR, China

<sup>4</sup>EuroMov Digital Health in Motion, Univ Montpellier, IMT Mines Ales, Ales, France

<sup>5</sup>Service de Médecine Nucléaire, Centre Hospitalier Universitaire de Nîmes, Université de Montpellier, Nîmes, France  
{zhuangzhenfeng, liyanfeng, lswang}@stu.xmu.edu.cn, cenmin0127@mail.ustc.edu.cn

## Abstract

Significant disparities between the features of natural images and those inherent to histopathological images make it challenging to directly apply and transfer pre-trained models from natural images to histopathology tasks. Moreover, the frequent lack of annotations in histopathology patch images has driven researchers to explore self-supervised learning methods like mask reconstruction for learning representations from large amounts of unlabeled data. Crucially, previous mask-based efforts in self-supervised learning have often overlooked the spatial interactions among entities, which are essential for constructing accurate representations of pathological entities. To address these challenges, constructing graphs of entities is a promising approach. In addition, the diffusion reconstruction strategy has recently shown superior performance through its random intensity noise addition technique to enhance the robust learned representation. Therefore, we introduce **H-MGDM**, a novel self-supervised Histopathology image representation learning method through the Dynamic Entity-Masked Graph Diffusion Model. Specifically, we propose to use complementary subgraphs as latent diffusion conditions and self-supervised targets respectively during pre-training. We note that the graph can embed entities' topological relationships and enhance representation. Dynamic conditions and targets can improve pathological fine reconstruction. Our model has conducted pretraining experiments on three large histopathological datasets. The advanced predictive performance and interpretability of H-MGDM are clearly evaluated on comprehensive downstream tasks such as classification and survival analysis on six datasets. Our code will be publicly available at <https://github.com/centurion-crawler/H-MGDM>.

## Introduction

Achieving concise and informative histopathology patch image representation is the cornerstone for solving many tasks in the field of computational histopathology analysis, such as cancer diagnosis, grading, segmentation, and

\*These authors contributed equally.

†Corresponding author.

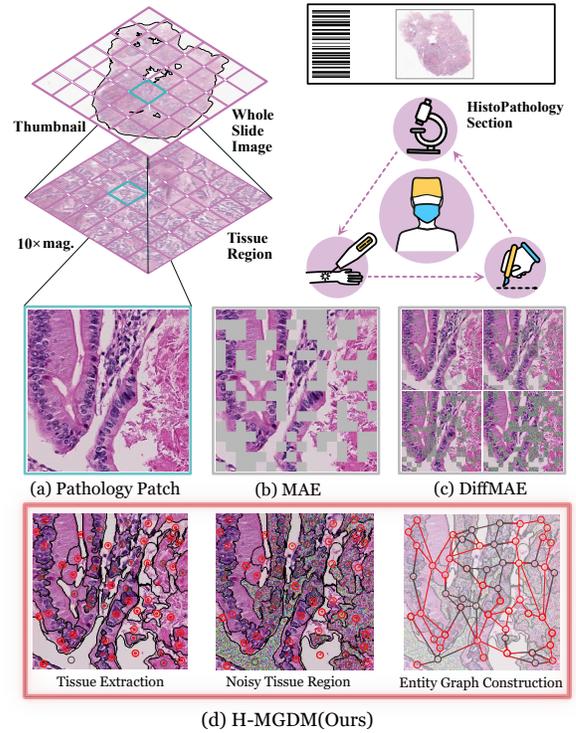


Figure 1: Pathological slide inspection process from the overall view to details. Unlike comparison methods, H-MGDM focuses on masked pathological tissue regions rather than grid tiles in patches, constructing the masked subgraph with varying intensities and noise for reconstruction by complementary conditional subgraph.

prognosis tasks within whole Slide Image (WSI) (Panayides et al. 2020). There are many extensive investigations on pre-training features for pathology images, which can help alleviate the time-consuming and tedious task of manual slide inspection by pathologists (Song et al. 2023). Today, the representation of pathology image patches relies heavily on transfer learning methods (Dosovitskiy et al. 2020; Sharmay

et al. 2021; Li et al. 2022) and the supervised pathology classification models such as KimiaNet (Riasatian et al. 2021). The above approaches exist problems like the domain gap and category bias (Guan and Liu 2021), and scarce and high-cost annotations limit retraining. Therefore, unlabeled self-supervised learning (SSL) has emerged as a solution to alleviate these limitations by learning salient representations without using labels.

Previous SSL methods like MAE (He et al. 2022) have shown that using masks and reconstruction tasks in SSL effectively enables models to learn from unlabeled pathological data. In the context of pathological images, the topological connections among pathological tissues, including cellular interactions and their surrounding environment, are crucial for various tasks. Recent approaches have underscored the importance of structure function relationships by linking the spatial organization of cells within tissues through cell graphs. These methods enable the extraction of biomarker-based pathological features (Jaume et al. 2021b,a), capturing complex semantic associations that extend beyond pixel-level data to encompass tissues and cells. These approaches align more closely with pathological diagnostic procedures. Entity-based topological analysis provides enhanced control over tissue modeling and facilitates the integration of pathological priors into task-specific histopathological entity representations. This indicates that it is crucial to recognize the interactions among "Image-to-Graph" based pathological tissue entities. However, when SSL is applied to pathological images, recent studies often focus on using pathological grid tiles in patches as masks, while neglecting the impact of entities (e.g., cellular or tissue regions) mask strategy on the overall semantic representation. Therefore, we introduce a new method to convert images into graphics to capture the structure of pathological entities, such as tissues, in self-supervised learning.

On the other hand, diffusion strategies have recently improved robust learning representations as conditions through their technique of adding noise with varying intensities (Purma et al. 2023; Wei et al. 2023; Yang and Wang 2023). We propose using an entity-masked graph as the input for the diffusion process, with encoded features from different layers of the graph serving as conditions to maintain strong performance. This approach captures powerful and complex entity information, thereby enhancing the representation of pathological images.

In summary, the motivation of our paper is to better learn the knowledge of entity graphs under self-supervised reconstruction progress. We propose a novel approach that converts histopathological images into entity graphs with dynamic mask and noise for diffusion pre-training to obtain better pathological image representations in the pathology inspection process (see Fig.1). Our contributions are:

- We propose a novel framework the **H-MGDM**, a **novel self-supervised histopathological image representation learning method through Dynamic Entity-Masked Graph Diffusion Model**. A strategy for partially visible entities as conditioning to prompt masked noisy entities to graph diffusion. Random masks and dynamic intensity noises can enhance representations in

histopathological images.

- In H-MGDM, we **convert pathology images into entity graphs of latent space to incorporate structural information of pathological entities**. This allows for more comprehensive spatial and semantic priors.
- We conducted pretraining experiments on three large histopathological datasets. The advanced predictive performance and interpretability of H-MGDM are clearly evaluated on comprehensive downstream tasks on six datasets. All performance **across several downstream tasks is averagely improved by 5.18%**. This demonstrates the effectiveness of the H-MGDM framework for pre-training in histopathological image analysis.

## Related Works

### Graph Representation for Digital Pathology

Recently, graph neural networks (GNN) have been employed to represent patches as graphs for pathology tasks. CGC-Net (Zhou et al. 2019) introduces a cell-graph convolutional neural network that converts large histology images into graphs, where vertices represent nuclei and edges denote cellular interactions based on vertex similarity. HACT (Pati et al. 2022) develops a hierarchical cell-to-tissue graph representation to jointly model both low-level and high-level graphs, incorporating intra- and inter-level coupling based on the topological distributions and interactions among entities. SHGNN (Hou et al. 2022a) proposes a novel spatial hierarchical GNN framework, equipped with a dynamic structure learning module, to capture entity location attributes and semantic representations, thereby extracting the characteristics of different entities in images. However, these methods focus supervision on global representation. Our proposed method introduces entities subgraph self-supervised targets. This enables entities to capture the contextual implications of local information.

### Denosing Diffusion Models

Diffusion models (Rombach et al. 2022; Ho, Jain, and Abbeel 2020) are known for their ability to generate sophisticated images under conditional control. Diffusion models also exhibit variable masking capabilities and have also been used to enhance representation learning in the self-supervised learning domain, in learning paradigms such as DiffAE (Preechakul et al. 2022). Also, GenSelfDiff-HIS (Purma et al. 2023) proposes a diffusion-based generative pre-training process for self-supervision to learn efficient histopathological image representations. DiffMAE (Wei et al. 2023) integrates diffusion's nuanced detail reconstruction capabilities with MAE's comprehensive semantic representation capability, symbolizing a convergence of methodologies in pursuit of enhanced representation.

### Mask for Self-Supervised Representation

Since the introduction of BERT (Devlin et al. 2018) in language models, masking prediction has reattracted attention. Lots of self-supervised masking methods have been

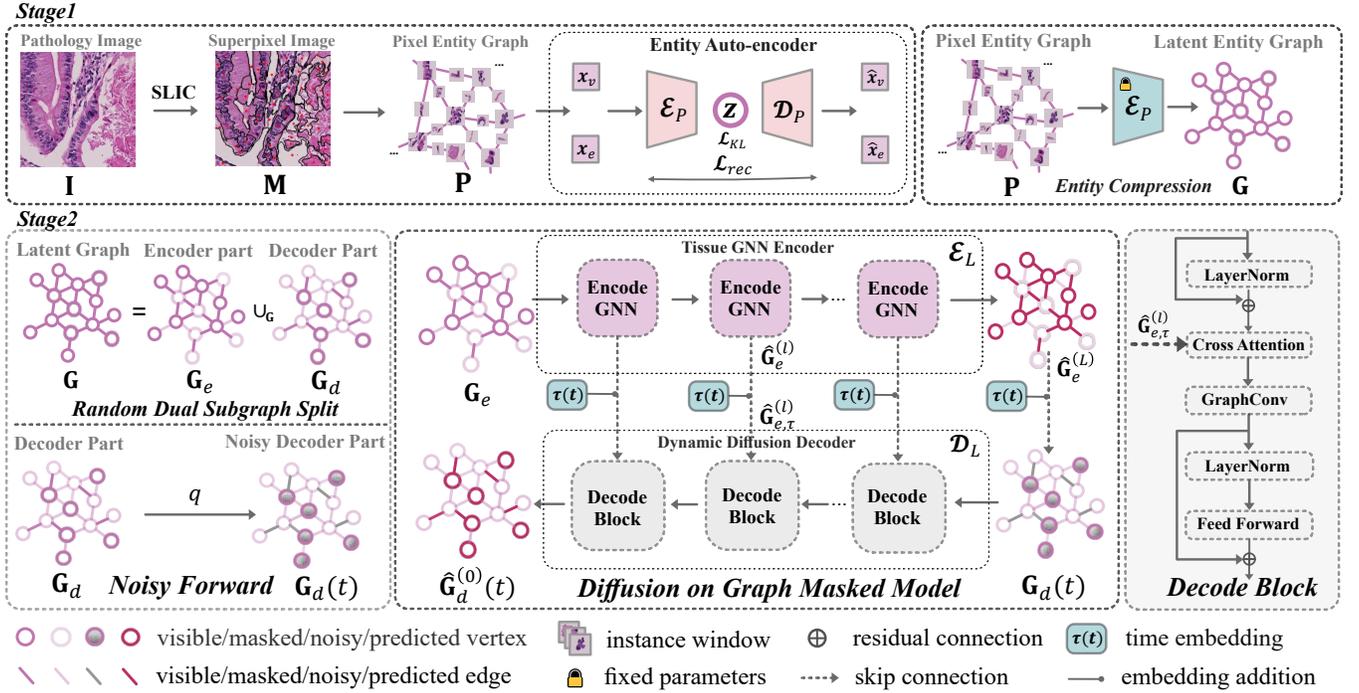


Figure 2: Overview of the H-MGDM pretraining stages. Conditional diffusion reverse process in the decoder.  $G_e$  and  $G_d$  are two complementary subgraphs of  $G$ .  $G_d(t)$  are from the diffusion forward process  $q_L$  of  $G_d$ . The target is to denoise  $G_d(t)$  to  $\hat{G}_d^{(0)}(t)$  close to  $G_d$  at sampling time  $t$ .

proposed in vision tasks, delineating enhanced representation techniques rooted in various theoretical frameworks: MAE (He et al. 2022) devised an asymmetric architecture tailored for pixel-level reconstruction, while MixMIM (Liu et al. 2022) endeavors to narrow the chasm between pre-training and fine-tuning through the stochastic amalgamation of masks. Additionally, GraphMAE (Hou et al. 2022b) and GraphMAE2 (Hou et al. 2023) advocate employing masked graph convolution to facilitate feature reconstruction, but the mentioned methods are confined to fixed-intensity masking approaches. The variation mechanism forces the features to adapt to each situation.

## Methodology

Figure 2 illustrates the framework of H-MGDM. First, the histopathological entity graph is constructed using the superpixel algorithm SLIC (Achanta et al. 2012) to extract the topological relations among tissues in the image and compress entities to latent space. Then, the GNN encoder is used to encode the visible subgraph as a condition. The latent graph diffusion model is introduced to reconstruct the dynamic self-supervised target of the masked subgraph to obtain robust representations of patches.

### Pathological Entity Graph Construction

To strengthen the concept of entity within limited structural constraints, we utilize priori pathological tissue superpixels as tissue entities when constructing the graph. First, a pathological image  $I \in \mathbb{N}^{h_I \times w_I \times 3}$  with the height  $h_I$  and the

width  $w_I$  is partitioned via SLIC algorithms (Achanta et al. 2012) resulting in a set of superpixels. For each superpixel,  $s$ , a window of size  $a \times a$  centered at  $s$  is considered as a vertex  $v$  of the pathological entity graph  $P$  in pixel space. Pixels in the window that do not belong to  $s$  are assigned the background color. Subsequently, edges will be established between every two vertices with adjacent boundaries, considering the local interactions between adjacent vertices more. The edge  $e$  originates from the region after the expansion operation along the boundary between  $s^i$  and another neighboring superpixel  $s^j$ . Thus, the image  $I$  is transformed into a pathological entity graph  $P(V_P, E_P, A, D)$ , where  $V_P = \{s_i\}_{i \in [0, N_V]}$ ,  $E_P = \{e_j\}_{j \in [0, N_E]}$  are sets of vertices and edges, respectively. And the adjacency matrix is  $A$  and the degree matrix of  $A$  is  $D$ .  $N_V$  and  $N_E$  are the numbers of vertices and edges, respectively.

### Entity Compression into Latent Space

Our compression model, based on previous work (Kingma and Welling 2013; Esser, Rombach, and Ommer 2021), utilizes an auto-encoder in stage 1. Given a pathological entity graph  $P$ , the encoder  $\mathcal{E}_P$  transforms each entity  $x \in (X_v \cup X_e) \subset \mathbb{N}^{a \times a \times 3}$  in  $P$  into a latent representation  $z = \mathcal{E}_P(x)$ , where  $z \in \mathbb{R}^{l \times l \times c}$ . The encoder learns an approximate posterior distribution  $q(z|x) = \mathcal{N}(z; \mu(x), \sigma^2(x)I)$ , with  $\mu(x)$  and  $\sigma(x)$  being learned mean and standard deviation of  $x$ . And the decoder  $\mathcal{D}_P$  then reconstructs the image from this latent space,  $\hat{x} = \mathcal{D}_P(z) = \mathcal{D}_P(\mathcal{E}_P(x))$ . The downsampling factor of the image is  $f = a/l$  and we ex-

plore various downsampling factors  $f$ . After the first stage of training is completed, we infer  $\mathcal{E}_P$  to encode all entities into the latent space, resulting in sets  $\mathbf{V}_G$  and  $\mathbf{E}_G$ . Those compose the latent space entity graph  $\mathbf{G}(\mathbf{V}_G, \mathbf{E}_G, \mathbf{A}, \mathbf{D})$ .

### Latent Graph Diffusion Model

The forward process of the latent diffusion model can add noise to the graph entities, describing the degraded sequence caused by Gaussian noise on the latent space, which does not contain much semantics. Given a well-defined latent diffusion diffusion (Ho, Jain, and Abbeel 2020) forward process  $q_L : \{\mathbf{G}_d(t)\}_{[0,T]}$  with variance time dependence, and a noise schedule  $\{\beta(t)\}_{[0,T]}$  where the integer time  $t \in [0, T]$ , based on Markov chain and diffusion characteristics (Huang et al. 2023; Wei et al. 2023), we have:

$$\begin{cases} q_L(\mathbf{G}_d(t)|\mathbf{G}_d(t-1)) = \mathcal{N}(\mathbf{G}_d(t)|\sqrt{1-\beta(t)}\mathbf{G}_d(t-1), \beta(t)\mathbf{I}) \\ q_L(\mathbf{G}_d(t)|\mathbf{G}_d(0)) = \mathcal{N}(\mathbf{G}_d(t)|\sqrt{\bar{\alpha}(t)}\mathbf{G}_d(0), (1-\bar{\alpha}(t))\mathbf{I}) \end{cases} \quad (1)$$

$\mathbf{G}_d(t)$  is reparameterized as  $\mathbf{G}_d(t) = \sqrt{\bar{\alpha}(t)}\mathbf{G}_d(0) + \sqrt{1-\bar{\alpha}(t)}\epsilon$ , noise follows a normal distribution:  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , where  $\alpha(t) = 1 - \beta(t)$ ,  $\bar{\alpha}(t) = \prod_{i=1}^t \alpha(i)$ , signal-to-noise ratio  $\{\frac{\alpha(t)}{\beta(t)}\}_{[0,T]}$  are chosen to noise gradually that  $q_L(\mathbf{G}_d(T)) \approx \mathcal{N}(\mathbf{0}, \mathbf{I})$ . The conditional diffusion reverse above by modeling the reverse distribution  $p_L$  which implies masked part conditioned on visible graph  $\hat{\mathbf{G}}_e$  from encoder:

$$p_L(\mathbf{G}_d(t-1)|\mathbf{G}_d(t), \hat{\mathbf{G}}_e) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (2)$$

Then, a reverse diffusion network  $\mathcal{D}_L$  with graph conditioning on  $\{\beta(t)\}_{[0,T]}$  is introduced. Considering the graph structure, we can apply the continuous diffusion process to  $\mathbf{V}_d(t)$  and  $\mathbf{E}_d(t)$  respectively to facilitate the restoration of noisy latent within subgraph pathological entities.

### Dynamic Diffusion on Masked Graph Model

In stage 2 illustrated in Fig. 2, an asymmetric auto-encoder mode is employed, utilizing GNN layers (Kipf and Welling 2016) as the encoder and ViT variants (Dosovitskiy et al. 2020) as the decoder. From the LDM perspective, the encoder also provides encoding conditions for the decoder’s denoising process. For a graph input  $G$ , is dynamically randomly divided into two complementary subgraphs  $\mathbf{G}_e(\mathbf{V}_e, \mathbf{E}_e, \mathbf{A}_e, \mathbf{D}_e)$  for encoder and  $\mathbf{G}_d(\mathbf{V}_d, \mathbf{E}_d, \mathbf{A}_d, \mathbf{D}_d)$  for decoder according to the given masking ratio  $r_m = \frac{N_{V_d}}{N_V} = \frac{N_{E_d}}{N_E}$ . The noise-added  $\mathbf{G}_d(t)$  serves as the initial input to the decoder. Furthermore, the topological information  $\mathbf{A}_* \mathbf{D}_* (* = e, d)$  is maintained during training.

**Tissue GNN Encoder** The latent encoder  $\mathcal{E}_L$  employs GNN that integrates tissue vertices and edges for  $L$  layers. In our pathological graph-based construction, the latent domains of graph vertices and edges are identical. Therefore, the vertex-based message passing can be used to forward edge latent for  $\hat{\mathbf{G}}_e^{(l)}(V_e^{(l)}, E_e^{(l)}, A_e, D_e)$  in the  $l$ -th layer:

$$\mathcal{E}_L^{(l)} : \begin{cases} \mathbf{V}_e^{(l+1)} = \sigma(\tilde{\mathbf{A}}\mathbf{V}_e^{(l)}\mathbf{W}_V^{(l)}) \\ \mathbf{E}_e^{(l+1)} = \sigma(\tilde{\mathbf{A}}^*\mathbf{E}_e^{(l)}\mathbf{W}_E^{(l)}) \end{cases}, \quad (3)$$

where  $\tilde{\mathbf{A}} = \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$ .  $\tilde{\mathbf{A}}$  and  $\tilde{\mathbf{A}}^*$  are the normalized symmetric adjacency matrices of the graph  $\mathbf{G}_e$  and the dual graph  $\mathbf{G}_e^*$ , respectively.  $\mathbf{V}_e^{(l)}, \mathbf{E}_e^{(l)}$  are inputs to the  $l$ -th layer,  $\mathbf{W}_V^{(l)}, \mathbf{W}_E^{(l)}$  are the vertices and edges weight matrices of graph convolution  $\mathcal{E}_L^{(l)}$ , and  $\sigma(\cdot)$  is a non-linear activation.

**Dynamic Diffusion Decoder** The decoder  $\mathcal{D}_L$  utilizes the conditional latent graph diffusion model. The noise level  $t$  serves as the forward sampling time during pre-training to generate  $\mathbf{G}_d(t)$ . Similar to the Transformer architecture (Vaswani et al. 2017), in the  $l$ -th decode block, cross-attention  $CA(\cdot, \cdot)$  utilizes the visible latent  $\hat{\mathbf{G}}_e^{(l)}$  from the  $l$ -th encoder layers as the conditional control after adding the time embedding  $\tau(t)$ :  $\hat{\mathbf{G}}_{e,\tau}^{(l)} = \hat{\mathbf{G}}_e^{(l)} + \tau(t)$ , aiding in denoising  $\hat{\mathbf{G}}_d^{(l)}(t)$  during decoding. And the graph convolution of decoder  $C_d^{(l)}$  is to perform the message passing of the fused latent according to  $A_d$  next. The Feed Forward  $\overline{FF}(\cdot)$  with layer normalized residual block is employed as the final layer within the Decoder Block to induce feature activation, resulting in  $\hat{\mathbf{G}}_d^{(l-1)}(t)$ . Then, it proceeds to the subsequent blocks to predict  $\hat{\mathbf{G}}_d^{(0)}(t)$ :

$$\begin{cases} \hat{\mathbf{G}}_d^{(l-1)}(t) = \overline{FF}(C_{dec}(CA(\hat{\mathbf{G}}_d^{(l)}(t), \hat{\mathbf{G}}_{e,\tau}^{(l)}))) \\ \hat{\mathbf{G}}_d^{(0)}(t) = \mathcal{D}_L(\mathbf{G}_d(t), t, \{\hat{\mathbf{G}}_e^{(l)}\}_{l \in [0,L]}) \end{cases} \quad (4)$$

The skip connection from  $\hat{\mathbf{G}}_{e,\tau}^{(l)}$  forms a U-shaped configuration, typically advantageous for graph representation across different levels and noisy latent restoration.

### Training Strategy

**Objectives** The optimization of the model in pre-taining has two stages: In order to avoid high-variance latent space in the first stage, the auto-encoder is optimized by a combination of a reconstruction loss and KL Divergence  $D_{KL}$  like VAE (Kingma and Welling 2013). The pixel-space reconstruction constraint  $\mathcal{L}_{rec}$  enforces local realism avoids blurriness and ensures that the reconstructions are confined to the image manifold. Here we use the MSE form:

$$\mathcal{L}_{rec} = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \|\mathbf{x} - \hat{\mathbf{x}}\|_2, \mathcal{L}_{VAE} = \mathcal{L}_{rec} + \lambda D_{KL}, \quad (5)$$

where  $\mathbb{E}_{q(\mathbf{z}|\mathbf{x})}$  represents the expectation of the distribution of the latent variable  $z$ ,  $\lambda$  is the loss weight. And  $q(\mathbf{z}|\mathbf{x})$  is the posterior distribution of the latent  $\mathbf{z}$  given  $\mathbf{x}$ .

The second stage is guided by minimizing the following objective function to optimize parameters  $\theta$  of the model. Notably,  $\mathbf{x}_0$ -mode is used dynamic denoising on graph masking of vertices and edges:

$$\mathcal{L}_{Simple} = \mathbb{E}_{t,\theta,\epsilon} [\|\mathbf{V}_d(0) - \hat{\mathbf{V}}_d^{(0)}(t)\| + \|\mathbf{E}_d(0) - \hat{\mathbf{E}}_d^{(0)}(t)\|]. \quad (6)$$

Here,  $\epsilon$  represents sampled noise. Leveraging variational inference, the well-known Mean Square Error (MSE) objective, derived from the evidence lower bound, is utilized to predict the denoised  $\hat{\mathbf{V}}_d^{(0)}$  and  $\hat{\mathbf{E}}_d^{(0)}$  as reconstruction targets.

Cancer Subtype / Tissue Classification									
Strategies			Datasets	Komura et al.		PANDA		IBD	
Graph	Mask	Diffusion		ACC(%)	F1(%)	ACC(%)	F1(%)	ACC(%)	F1(%)
✗	✗	✗	SimCLR (Chen et al. 2020)	69.24±2.12	62.45±1.72	61.10±2.19	58.48±1.25	71.44±1.54	69.47±2.43
✗	✗	✗	KimiaNet* (Riasatian et al. 2021)	72.66±1.69	65.93±1.75	67.68±1.88	55.40±1.19	76.42±0.98	70.75±1.37
✗	✓	✗	Dino (Caron et al. 2021)	78.05±1.36	70.64±2.01	69.92±1.28	64.11±1.31	82.45±0.82	75.23±1.41
✗	✓	✗	MAE (He et al. 2022)	77.37±1.51	69.44±0.89	70.51±1.78	62.73±0.89	81.06±0.88	76.44±1.02
✓	✓	✗	GraphMAE (Hou et al. 2022b)	76.69±2.10	67.60±1.87	72.22±2.19	60.34±2.55	75.85±1.25	72.78±0.81
✓	✓	✗	GraphMAE2 (Hou et al. 2023)	78.86±3.05	65.97±1.05	72.56±1.80	63.49±2.64	78.12±1.42	72.32±0.79
✗	✗	✓	DiffAE (Preechakul et al. 2022)	79.11±1.92	68.18±1.74	71.54±2.12	61.62±1.11	81.24±1.48	74.51±2.08
✗	✓	✓	DiffMAE (Wei et al. 2023)	78.14±2.10	70.23±2.14	71.92±1.22	65.82±1.82	84.58±1.72	74.74±2.15
✓	✓	✓	<b>H-MGDM (Ours)</b>	<b>82.06±1.17</b>	<b>72.41±1.36</b>	<b>74.51±1.13</b>	<b>67.32±1.40</b>	<b>86.23±2.31</b>	<b>78.92±1.79</b>
Ablation study			w/o edge latents	80.91±1.92	70.85±1.58	72.39±1.29	65.89±1.82	84.11±1.72	77.41±1.54
			w/o skip connection	79.29±1.27	68.43±1.45	70.98±1.75	62.11±2.16	82.35±2.32	74.61±2.14
			noise intensity fixed	79.70±1.14	67.14±1.24	72.52±1.65	63.73±2.31	83.58±1.69	77.24±2.53

Table 1: Comparing classification performance across methods with ablations. The best results are marked in **bold**. “\*” indicates label-supervised of the method. The results are reported in: *mean±std* (where *std* stands for standard deviation).

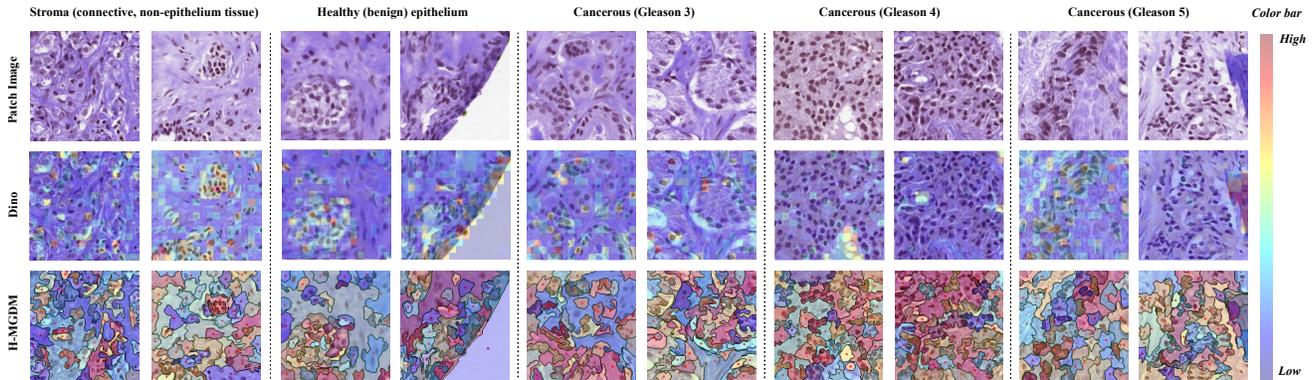


Figure 3: Original images and their attention heatmaps of five different categories of the PANDA dataset, showing the interpretability of our method under the pathological entity graph construction comparison with Dino.

**Downstream Tasks** For downstream tasks, we deploy two stages’ encoders  $\mathcal{E}_P$   $\mathcal{E}_L$  to inference to obtain the global graph representation  $\mathbf{O}_G$  by readout  $r_G$ :

$$\mathbf{G}_o = \hat{\mathbf{G}}^{(L)} = \mathcal{E}_L(\mathcal{E}_P(\mathbf{G}_P)), \mathbf{O}_G = r_G(\mathbf{G}_o). \quad (7)$$

- **Classification.** For the downstream tuning, the cross-entropy loss  $\mathcal{L}_{CE}$  is adopted for patch classification tasks. The model is optimized through the classification layer  $\text{MLP}_c$  to obtain the predicted  $\hat{Y}$  by global representation  $\mathbf{O}_G$  as probabilities of the classes and to supervise it with the classification labels  $Y$ .
- **Regression.** We introduce the Cox proportional hazards model, a semi-parametric regression model. Using survival events and survival time as labels. Here, pathological image features are used to predict risks and analyze the impact on survival. Considering a problem involving two explanatory variables as predictors of survival time  $t_i$  and  $t_j$  of patients  $i, j$ ,  $\delta$  represents the termination event (1: death, recurrence, 0: not occurred)  $R(t_i)$  represents the condition  $t_j > t_i$ , for neg log partial likelihood loss:

$$\mathcal{L}_{Cox} = - \sum_{i:\delta_i=1} [\hat{h}_i - \log \sum_{j \in R(t_i)} e^{\hat{h}_j}], \quad (8)$$

where we use  $\text{MLP}_h$  linear map the graph representation  $\mathbf{O}_{G_i}$  to the final hazards prediction  $h_i$ .

## Experiments

### Experimental Settings

**Datasets** The proposed framework underwent pretraining and classification using three large extensive histopathology datasets: **Komura et al.** (Komura et al. 2022) (1.6M images, 32 cancer types), **Prostate Annotation Data Archive (PANDA)** dataset (Bulten et al. 2022) (11,000 digitized H&E-stained WSIs, obtaining 12.5M images with 6 level Gleason region annotations), our in-house colorectal cancer data **IBD** (23M images, with 360K patches annotated into 9 common tissue types). For the survival analysis task, we compare the performances on two public cohorts: **TCGA-KIRC** (512 cases) and **TCGA-ESCA** (155 cases) and one privately collected primary-metastatic pathology colorectal cancer dataset **CRC-PM** (388 cases). We use different methods pre-trained on the pancancer dataset Komura et al. as feature extractors for patches in WSI <sup>1</sup>.

<sup>1</sup>More descriptions in Appendix B.2 about datasets

Survival Analysis Regression										
Methods	Datasets	TCGA-KIRC			TCGA-ESCA			CRC-PM		
		DeepSurv	AB-MIL	PatchGCN	DeepSurv	AB-MIL	PatchGCN	DeepSurv	AB-MIL	PatchGCN
SimCLR (Chen et al. 2020)		61.91±4.71	62.09±4.33	62.46±4.26	59.26±4.35	58.06±4.60	62.59±4.82	56.30±2.47	58.76±6.51	59.87±4.06
KimiaNet* (Riasatian et al. 2021)		62.16±4.72	64.12±5.28	65.76±3.14	60.53±6.69	59.00±2.97	58.16±5.75	59.06±9.04	59.95±6.11	57.67±8.89
Dino (Caron et al. 2021)		67.92±5.61	66.94±4.42	66.64±3.01	57.88±3.91	58.85±4.14	56.58±5.11	58.02±6.72	61.01±9.91	63.20±8.50
MAE (He et al. 2022)		57.78±2.42	61.30±2.54	65.08±3.69	59.71±5.02	57.84±6.29	60.62±5.21	60.92±7.63	59.44±8.98	62.80±8.82
GraphMAE2 (Hou et al. 2023)		64.31±6.73	66.76±4.21	65.84±2.04	57.26±5.55	59.71±4.90	60.35±5.55	58.16±6.29	59.38±6.94	60.42±7.08
DiffAE (Preechakul et al. 2022)		63.26±3.88	67.31±1.98	68.29±3.84	60.74±5.54	60.23±5.74	63.61±5.49	60.81±6.30	60.61±5.80	62.87±9.52
DiffMAE (Wei et al. 2023)		62.41±3.24	65.92±5.04	67.32±4.18	60.60±2.57	59.49±4.83	63.89±6.32	61.16±5.27	60.64±6.05	61.29±7.12
<b>H-MGDM (Ours)</b>		<b>66.99±4.56</b>	<b>69.88±3.97</b>	<b>71.17±4.51</b>	<b>62.68±3.04</b>	<b>62.55±3.28</b>	<b>64.82±3.45</b>	<b>62.27±5.09</b>	<b>63.89±6.94</b>	<b>66.05±7.81</b>
w/o edge latent		65.91±5.74	67.88±4.73	70.53±5.01	59.23±4.29	60.62±4.29	64.10±4.83	<b>62.32±6.12</b>	61.39±7.12	65.12±7.28
w/o skip connection		63.64±5.32	67.00±3.10	69.07±5.62	58.66±3.85	60.34±3.24	59.40±2.36	59.78±2.91	60.09±4.57	60.97±8.36
noise intensity $t$ fixed		63.16±3.62	69.44±4.50	69.99±4.45	59.17±6.49	58.97±3.97	61.17±3.80	60.34±7.38	62.05±8.69	62.33±5.29

Table 2: Survival Analysis performance across SOTAs features on external public validation data by conducting the 5-fold evaluation procedure with 5 runs. The experimental results of CI are reported by *mean±std*.

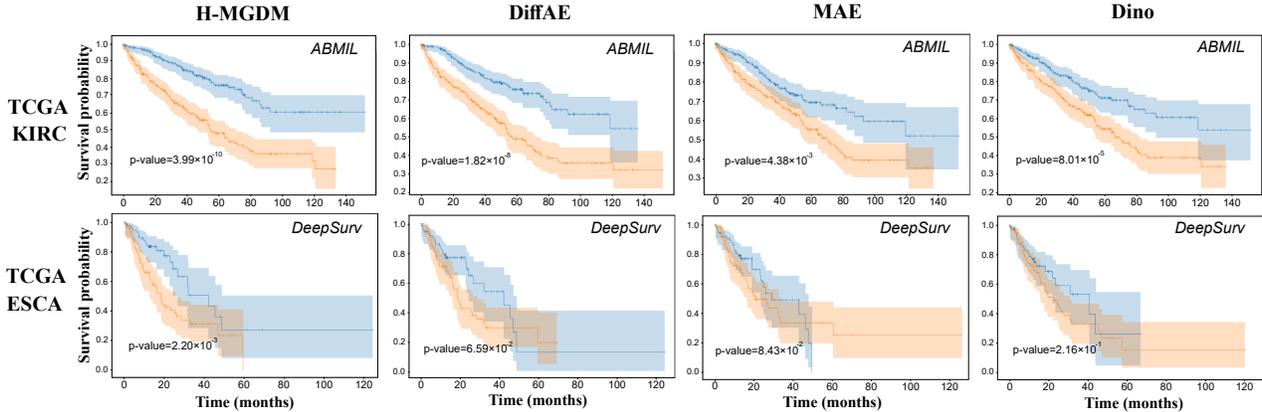


Figure 4: Kaplan-Meier Analysis of comparison methods and our framework. All patients from the five tests were pooled and analyzed. Each cohort is split into a high-risk (orange) and a low-risk group (blue) according to the median output of the cohort.

**Implementation Details** Experiments involving H-MGDM and comparative methods are conducted with a batch size set to 64 max. SLIC with the initial region number of 500, window size  $a$  is 64. Latent downsampling factors  $f = 2$ . Also, pre-training is optimized by Adam (Kingma and Ba 2014) with an initial learning rate of  $3e-4$  and the plateau scheduler with a minimum learning rate of  $1e-5$  for 250 epochs. The noise timesteps  $T$  is 1000, and the sigmoid schedule  $\{\beta(t)\}_{[0,T]}$  from  $1e-7$  to  $2e-3$  is used.

**Evaluation Metrics** Accuracy (ACC) and Marco F1 are used to evaluate the methods in classification comparison. To evaluate the effectiveness of entity latent diffusion, the Root Mean Square Error (RMSE) for graph entities is calculated as a qualitative evaluation metric. Harrell’s concordance index (C-index or CI) measures the survival model’s ability to correctly provide a reliable ranking of the survival times based on the individual risk scores. It ranges from 0 to 1, with higher values indicating better performance.

### Comparison with Baseline Methods

**Results** We conduct comparative experiments among our H-MGDM and other baseline pre-training models: *i.e.* SimCLR (Chen et al. 2020), Dino, MAE, GraphMAE, Graph-

MAE2, DiffAE, DiffMAE. We also mark the use of strategies (graph construction, mask strategy, diffusion-guided) by the comparison methods in Table 1. For comparison of the features from various pre-training methods, we use three backbones: DeepSurv (Katzman et al. 2018), AB-MIL (Ilse, Tomczak, and Welling 2018) and PatchGCN (Chen et al. 2021) for the survival prediction task. To the best of our knowledge, H-MGDM is the first time these three mechanisms have been introduced simultaneously into histopathology pre-training. Our method achieved outstanding performance due to incorporating structural information of pathological entities and enhancing mask learning during the diffusion procedure. The results are presented in Table 1, 2. For all baselines, our method achieves an average improvement of 5.99%, 5.43%, and 4.146% on the ACC, F1, and CI.

**Ablation Study** To further explore the proposed components with the effectiveness of our model. We perform ablation experiments in Table 1, 2. The study investigates the influence of edges in graph data, skip connections from encoder representations as diffusion conditioning, and time-varying intensity noise on model performance. The results show these components can enhance the representational capacity of downstream tasks with improvement of 2.50%,

3.17%, and 2.46% on the ACC, F1, and CI metrics.

**Kaplan-Meier Analysis and Significance** Based on the median survival risk output by our model, each cancer cohort is divided into high-risk and low-risk groups. If the survival predictions are consistent, the Kaplan-Meier (KM) curves of these groups should show significant differences. As expected, Fig. 4 shows our method’s log-rank p-values less than 0.05 for two different cancer cohorts, indicating the effectiveness of in predicting patient survival.

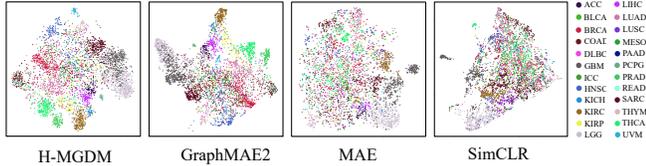


Figure 5: T-SNE plots of pan-cancer samples’ readout representations learning with H-MGDM and baseline methods.

**T-SNE Visualization of pan-cancer representations** To further confirm the performances of pre-training models, we randomly select 6.4K samples from the trained pre-training model from the Komura et al. (Komura et al. 2022) data test set and visualize the distribution of t-SNE (Hinton and Roweis 2002) as shown in Fig. 5. Our H-MGDM model can better distinguish HNSC (Head and Neck Squamous Cell Carcinoma), LUAD (Lung Adenocarcinoma) and UVM (Uveal Melanoma) etc. cancer types. They have high intra-class aggregation and inter-class separation.

### Analysis of Our Framework

**Interpretability** The graphical representation of our approach provides scalable interpretability. Unlike posterior explanations of previous methods like Grad-CAM (Selvaraju et al. 2017) and Shap (Štrumbelj and Kononenko 2014), H-MGDM can directly generate high-attention regions during inference readout. We visualize the attention heatmap in Fig. 3 of tissue regions sampling from normal to pathological classes. We also visualize the multi-head attention of the Dino method for comparison, highlighting its weaker resistance to interference (some Gleason 4 and 5 samples are focused on blank and outlined regions).

**Masking ratio Investigation** We examine rates of the dynamic masking strategy, which correspond to the ratio of the target subgraph division. As shown in Table 3, the entity graph masking rate  $r_m$  is about 50%-70% for reconstruction, and the rest is used as condition learning. Such pre-training features can improve the classification performance of the three datasets. Lower masking rates may hinder the accurate

Mask Ratio	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
Komura et.al.	65.87	69.90	74.44	77.07	78.38	<b>82.06</b>	79.19	76.24	73.91
PANDA	63.12	63.15	67.88	65.37	70.51	72.49	<b>74.51</b>	70.47	66.24
IBD	70.54	75.23	79.54	84.51	<b>86.23</b>	82.72	79.15	74.15	64.31

Table 3: Hyperparameter Investigation of mask ratios

learning of target entity semantics within the visible parts, while higher rates may result in excessive difficulty.

Strategy	NtoN	EtoE	NtoE	EtoN	NtoN & EtoE	NtoE & EtoN
ACC	73.74	71.59	71.98	68.59	<b>74.51</b>	73.82
RMSE	0.136	0.155	0.169	0.184	0.141	<b>0.127</b>

Table 4: Effectiveness of various decoder conditioning strategies on latent restoration (RMSE) and classification (ACC).

**Decode Strategy Studies** The decoder’s conditioning strategies need empirical investigation. We test six cross-attention block strategies by aligning different graph attributes between the encoder and decoder (N for vertex, E for edge): NtoN, EtoE, NtoE, EtoN, NtoN & EtoE, and NtoE & EtoN. Graph vertices generally contain more information than edges. PANDA results in Table 4 show that graph attribute alignment enhances representation for classification, while attribute heterogeneity improves reconstruction.

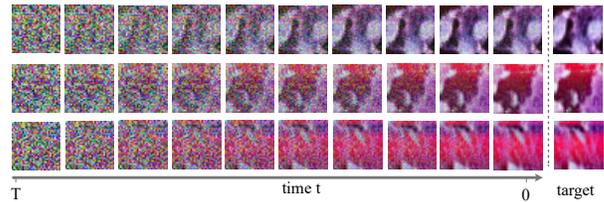


Figure 6: Visualization of diffusion process over time t.

**Entity Latent Diffusion Process** Fig. 6 represents the entity restoration as a pre-training proxy task in the latent space. We try to reconstruct fine-grained latent targets throughout the sampling process: as the sampling time  $t$  iterates, the latent similarity between diffuse entities and entities increases during the reverse of the process. It Demonstrates excellent representation for latent recovery by the conditioning encoder  $\mathcal{E}_L$  in the diffusion decoder  $\mathcal{D}_L$  of H-MGDM.

### Conclusion

Our proposed novel framework, the dynamic entity mask on graph diffusion model for histopathology (H-MGDM), addresses the challenge in representation learning and enhances pre-training histopathology representation by incorporating tissue structural information and a suitable masking technique to guide diffusion models during reconstruction. The results of our experiments on six histopathology datasets covering common cancer types demonstrate the superior classification and regression performance of H-MGDM compared to existing methodologies. We anticipate the widespread applicability of H-MGDM in diverse downstream tasks such as prognosis and generation. Additionally, we are committed to exploring the potential of the newer entity extraction methods, the deeper interpretability, and the larger-scale experiments of our framework in our future research endeavors.

## Acknowledgments

This work was supported by National Natural Science Foundation of China (Grant No. 62371409) and Fujian Provincial Natural Science Foundation of China (Grant No. 2023J01005).

## References

- Achanta, R.; Shaji, A.; Smith, K.; Lucchi, A.; Fua, P.; and Süsstrunk, S. 2012. SLIC superpixels compared to state-of-the-art superpixel methods. *IEEE transactions on pattern analysis and machine intelligence*, 34(11): 2274–2282.
- Brodersen, K. H.; Ong, C. S.; Stephan, K. E.; and Buhmann, J. M. 2010. The balanced accuracy and its posterior distribution. In *2010 20th international conference on pattern recognition*, 3121–3124. IEEE.
- Bulten, W.; Kartasalo, K.; Chen, P.-H. C.; Ström, P.; Pinckaers, H.; Nagpal, K.; Cai, Y.; Steiner, D. F.; van Boven, H.; Vink, R.; et al. 2022. Artificial intelligence for diagnosis and Gleason grading of prostate cancer: the PANDA challenge. *Nature medicine*, 28(1): 154–163.
- Caron, M.; Touvron, H.; Misra, I.; Jégou, H.; Mairal, J.; Bojanowski, P.; and Joulin, A. 2021. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, 9650–9660.
- Chen, R. J.; Lu, M. Y.; Shaban, M.; Chen, C.; Chen, T. Y.; Williamson, D. F.; and Mahmood, F. 2021. Whole slide images are 2d point clouds: Context-aware survival prediction using patch-based graph convolutional networks. In *Medical Image Computing and Computer Assisted Intervention—MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part VIII 24*, 339–349. Springer.
- Chen, T.; Kornblith, S.; Norouzi, M.; and Hinton, G. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, 1597–1607. PMLR.
- Cui, M.; and Zhang, D. Y. 2021. Artificial intelligence and computational pathology. *Laboratory Investigation*, 101(4): 412–422.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Esser, P.; Rombach, R.; and Ommer, B. 2021. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 12873–12883.
- Fey, M.; and Lenssen, J. E. 2019. Fast Graph Representation Learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*.
- Guan, H.; and Liu, M. 2021. Domain adaptation for medical image analysis: a survey. *IEEE Transactions on Biomedical Engineering*, 69(3): 1173–1185.
- He, K.; Chen, X.; Xie, S.; Li, Y.; Dollár, P.; and Girshick, R. 2022. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 16000–16009.
- Hinton, G. E.; and Roweis, S. 2002. Stochastic neighbor embedding. *Advances in neural information processing systems*, 15.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33: 6840–6851.
- Hou, W.; Huang, H.; Peng, Q.; Yu, R.; Yu, L.; and Wang, L. 2022a. Spatial-hierarchical graph neural network with dynamic structure learning for histological image classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 181–191. Springer.
- Hou, Z.; He, Y.; Cen, Y.; Liu, X.; Dong, Y.; Kharlamov, E.; and Tang, J. 2023. GraphMAE2: A Decoding-Enhanced Masked Self-Supervised Graph Learner. In *Proceedings of the ACM Web Conference 2023*, 737–746.
- Hou, Z.; Liu, X.; Cen, Y.; Dong, Y.; Yang, H.; Wang, C.; and Tang, J. 2022b. GraphMAE: Self-Supervised Masked Graph Autoencoders. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 594–604.
- Huang, H.; Sun, L.; Du, B.; and Lv, W. 2023. Conditional diffusion based on discrete graph structures for molecular graph generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 4302–4311.
- Humphrey, P. A. 2004. Gleason grading and prognostic factors in carcinoma of the prostate. *Modern pathology*, 17(3): 292–306.
- Ilse, M.; Tomczak, J.; and Welling, M. 2018. Attention-based deep multiple instance learning. In *International conference on machine learning*, 2127–2136. PMLR.
- Jaume, G.; Pati, P.; Anklin, V.; Foncubierta, A.; and Gabrani, M. 2021a. Histocartography: A toolkit for graph analytics in digital pathology. In *MICCAI Workshop on Computational Pathology*, 117–128. PMLR.
- Jaume, G.; Pati, P.; Bozorgtabar, B.; Foncubierta, A.; Annciello, A. M.; Feroce, F.; Rau, T.; Thiran, J.-P.; Gabrani, M.; and Goksel, O. 2021b. Quantifying explainers of graph neural networks in computational pathology. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 8106–8116.
- Kang, M.; Song, H.; Park, S.; Yoo, D.; and Pereira, S. 2023. Benchmarking Self-Supervised Learning on Diverse Pathology Datasets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3344–3354.
- Katzman, J. L.; Shaham, U.; Cloninger, A.; Bates, J.; Jiang, T.; and Kluger, Y. 2018. DeepSurv: personalized treatment recommender system using a Cox proportional hazards deep neural network. *BMC medical research methodology*, 18: 1–12.

- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kingma, D. P.; and Welling, M. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- Kipf, T. N.; and Welling, M. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*.
- Komura, D.; Kawabe, A.; Fukuta, K.; Sano, K.; Umezaki, T.; Koda, H.; Suzuki, R.; Tominaga, K.; Ochi, M.; Konishi, H.; et al. 2022. Universal encoding of pan-cancer histology by deep texture representations. *Cell Reports*, 38(9).
- Li, X.; Cen, M.; Xu, J.; Zhang, H.; and Xu, X. S. 2022. Improving feature extraction from histopathological images through a fine-tuning ImageNet model. *Journal of Pathology Informatics*, 13: 100115.
- Liu, C.; Fan, W.; Liu, Y.; Li, J.; Li, H.; Liu, H.; Tang, J.; and Li, Q. 2023. Generative diffusion models on graphs: Methods and applications. *arXiv preprint arXiv:2302.02591*.
- Liu, J.; Huang, X.; Liu, Y.; and Li, H. 2022. Mixmim: Mixed and masked image modeling for efficient visual representation learning. *arXiv preprint arXiv:2205.13137*.
- Liu, J.; Lichtenberg, T.; Hoadley, K. A.; Poisson, L. M.; Lazar, A. J.; Cherniack, A. D.; Kovatich, A. J.; Benz, C. C.; Levine, D. A.; Lee, A. V.; et al. 2018. An integrated TCGA pan-cancer clinical data resource to drive high-quality survival outcome analytics. *Cell*, 173(2): 400–416.
- Lu, M. Y.; Williamson, D. F.; Chen, T. Y.; Chen, R. J.; Barbieri, M.; and Mahmood, F. 2021. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature Biomedical Engineering*, 5(6): 555–570.
- Panayides, A. S.; Amini, A.; Filipovic, N. D.; Sharma, A.; Tsaftaris, S. A.; Young, A.; Foran, D.; Do, N.; Golemati, S.; Kurc, T.; et al. 2020. AI in medical imaging informatics: current challenges and future directions. *IEEE journal of biomedical and health informatics*, 24(7): 1837–1857.
- Pati, P.; Jaume, G.; Foncubierta-Rodriguez, A.; Feroce, F.; Anniciello, A. M.; Scognamiglio, G.; Brancati, N.; Fiche, M.; Dubruc, E.; Riccio, D.; et al. 2022. Hierarchical graph representations in digital pathology. *Medical image analysis*, 75: 102264.
- Preechakul, K.; Chatthee, N.; Wizadwongsa, S.; and Suwanajakorn, S. 2022. Diffusion Autoencoders: Toward a Meaningful and Decodable Representation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Purma, V.; Srinath, S.; Srirangarajan, S.; Kakkar, A.; et al. 2023. GenSelfDiff-HIS: Generative Self-Supervision Using Diffusion for Histopathological Image Segmentation. *arXiv preprint arXiv:2309.01487*.
- Riasatian, A.; Babaie, M.; Maleki, D.; Kalra, S.; Valipour, M.; Hemati, S.; Zaveri, M.; Safarpour, A.; Shafiei, S.; Afshari, M.; Rasoolijaberi, M.; Sikaroudi, M.; Adnan, M.; Shah, S.; Choi, C.; Damaskinos, S.; Campbell, C. J.; Diamandis, P.; Pantanowitz, L.; Kashani, H.; Ghodsi, A.; and Tizhoosh, H. 2021. Fine-Tuning and training of densenet for histopathology image representation using TCGA diagnostic slides. *Medical Image Analysis*, 70: 102032.
- Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695.
- Selvaraju, R. R.; Cogswell, M.; Das, A.; Vedantam, R.; Parikh, D.; and Batra, D. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, 618–626.
- Sharmay, Y.; Ehsany, L.; Syed, S.; and Brown, D. E. 2021. HistoTransfer: Understanding Transfer Learning for Histopathology. In *2021 IEEE EMBS International Conference on Biomedical and Health Informatics (BHI)*, 1–4.
- Song, A. H.; Jaume, G.; Williamson, D. F.; Lu, M. Y.; Vaidya, A.; Miller, T. R.; and Mahmood, F. 2023. Artificial intelligence for digital and computational pathology. *Nature Reviews Bioengineering*, 1(12): 930–949.
- Štrumbelj, E.; and Kononenko, I. 2014. Explaining prediction models and individual predictions with feature contributions. *Knowledge and information systems*, 41: 647–665.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Wei, C.; Mangalam, K.; Huang, P.-Y.; Li, Y.; Fan, H.; Xu, H.; Wang, H.; Xie, C.; Yuille, A.; and Feichtenhofer, C. 2023. Diffusion Models as Masked Autoencoder. In *ICCV*.
- Yang, X.; and Wang, X. 2023. Diffusion model as representation learner. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 18938–18949.
- Zhou, Y.; Graham, S.; Koohbanani, N. A.; Shaban, M.; Heng, P.-A.; and Rajpoot, N. 2019. CGC-Net: Cell Graph Convolutional Network for Grading of Colorectal Cancer Histology Images. In *The IEEE International Conference on Computer Vision (ICCV) Workshops*.

## Supplementary Material

We firstly provide the H-MGDM inference algorithm and introduce downstream tasks in Appendix A. The comprehensive introduction to the comparison methods is in Appendix B. We also provide additional interpretability and Kaplan-Meier analysis results in Appendix C. The limitations of our approach and areas for future work are shown in Appendix D.

### A. Preliminaries and Details of Inference

In this part, we make more additions to the graph latent space diffusion model, and model inference for downstream tasks.

#### A.1 Preliminaries of Latent Diffusion on Graph

Given a latent space instance graph  $\mathbf{G}(\mathbf{V}, \mathbf{E}, \mathbf{A}, \mathbf{D})$ . DDPM algorithm can generate graphs refer to (Liu et al. 2023):

$$\begin{aligned} q(\mathbf{G}(t)|\mathbf{G}(t-1)) &= (\mathbf{V}(t-1)\mathbf{Q}_t^{\mathbf{V}}, \mathbf{E}(t-1)\mathbf{Q}(t)^{\mathbf{E}}) \\ q(\mathbf{G}(t)|\mathbf{G}) &= (\mathbf{V}\tilde{\mathbf{Q}}(t)^{\mathbf{V}}, \mathbf{E}\tilde{\mathbf{Q}}(t)^{\mathbf{E}}) \end{aligned} \quad (9)$$

where  $\mathbf{G}(t) = (\mathbf{V}(t), \mathbf{E}(t))$  refers to the noisy graph composed of the node feature matrix  $\mathbf{V}_t$  and the edge attribute tensor  $\mathbf{E}_t$  at step  $t$ .  $\mathbf{Q}_t^{\mathbf{V}}$  and  $\mathbf{Q}_t^{\mathbf{E}}$  refer to the noise added to the node and edge, respectively. This Markov formulation allows sampling directly at an arbitrary time step without computing the previous steps. And  $\mathbf{G}_t$  can denoise with the conditional graph  $\mathbf{G}_c$  as follows:

$$\begin{aligned} p(\mathbf{G}(t-1)|\mathbf{G}(t), \mathbf{G}_c, t) \\ = (\mathbf{V}_t \mathbf{Q}_t^{CA(\mathbf{V}, \mathbf{V}_c)}, \mathbf{E}_t \mathbf{Q}_t^{CA(\mathbf{E}, \mathbf{E}_c)}) \end{aligned} \quad (10)$$

we introduce a domain-specific encoder to project  $G_c$  into an intermediate representation which is then mapped to the intermediate layers of the UNet structure via a cross-attention  $CA(\cdot, \cdot) \text{ softmax}(\frac{Q_* K_*^T}{V_*})$  with :

$$\begin{cases} Q_{\mathbf{V}} = W_{Q_{\mathbf{V}}} \cdot \mathbf{V}, K_{\mathbf{V}} = W_{K_{\mathbf{V}}} \cdot \mathbf{V}_c, V_{\mathbf{V}} = W_{V_{\mathbf{V}}} \cdot \mathbf{V}_c \\ Q_{\mathbf{E}} = W_{Q_{\mathbf{E}}} \cdot \mathbf{E}, K_{\mathbf{E}} = W_{K_{\mathbf{E}}} \cdot \mathbf{E}_c, V_{\mathbf{E}} = W_{V_{\mathbf{E}}} \cdot \mathbf{E}_c \end{cases}, \quad (11)$$

here,  $W_{Q_{\mathbf{V}}}, W_{K_{\mathbf{V}}}, W_{V_{\mathbf{V}}}, W_{Q_{\mathbf{E}}}, W_{K_{\mathbf{E}}}$  and  $W_{V_{\mathbf{E}}}$  are learnable projection matrices. Then we learn the conditional LDM on the graph optimized with:

$$\mathcal{L} = \mathbb{E}_{\mathcal{E}(\mathbf{G}), \mathbf{G}_c, \epsilon \sim \mathcal{N}(0,1), t} [|\epsilon - \epsilon_{\theta}(\mathbf{z}_t, t, \tau_{\theta}(\mathbf{G}_c))|], \quad (12)$$

where both  $\tau_{\theta}$  and  $\epsilon_{\theta}$  are jointly optimized via Eq.12. This conditioning mechanism is flexible with the same domain experts. In our method,  $\mathbf{G}_c$  are conditioning prompts in the dynamic diffusion decoder.

#### A.2 Detailed description of Inference and Downstream tasks

**Inference description** As mentioned in Section **Training Strategy**, the original image first undergoes two pre-training stages of encoders' inference to obtain the readout graph representation  $\mathbf{G}_o$  as the final representation of the image  $\mathbf{I}$  in algorithm 1. In Graph Neural Networks (GNNs), the readout operation creates a global representation of the graph

---

#### Algorithm 1: Inference of H-MGDM framework

---

**Require:** Pathology Image  $\mathbf{I}$ ; Optimized entity compressed encoder  $\mathcal{E}_P$  and tissue GNN encoder  $\mathcal{E}_L$  with  $L$  layers graph convolution, the Readout layer  $r_{\mathbf{G}}$ ;  
**Ensure:** Pooled and task-specific readout global graph representation  $\mathbf{O}_{\mathbf{G}}$ ;  
*# stage one: entity compression*  
1: Superpixel Image  $\mathbf{M} \leftarrow \mathbf{I}$  using SLIC;  
2: Adjacency algorithm constructs tissue entity graph  $\mathbf{P}$  from  $\mathbf{M}$  with vertex tiles  $\mathbf{V}_{\mathbf{P}}$ , edge tiles  $\mathbf{E}_{\mathbf{P}}$ , Adjacency matrix  $\mathbf{A}$  and Degree matrix  $\mathbf{D}$ ;  
3: **for** entity  $\mathbf{x}_{\mathbf{V}}$  in  $\mathbf{V}_{\mathbf{P}}$  **do**  
4:    $\mathbf{z}_{\mathbf{V}} \leftarrow \mathcal{E}_P(\mathbf{x}_{\mathbf{V}})$   
5: **end for**  
6: **for** entity  $\mathbf{x}_{\mathbf{E}}$  in  $\mathbf{E}_{\mathbf{P}}$  **do**  
7:    $\mathbf{z}_{\mathbf{E}} \leftarrow \mathcal{E}_P(\mathbf{x}_{\mathbf{E}})$   
8: **end for**  
9: Construct graph  $\mathbf{G}$ 's vertex and edge sets in latent:  
 $\mathbf{V}_{\mathbf{G}} \leftarrow \{\mathbf{z}_{\mathbf{V}}\}; \mathbf{E}_{\mathbf{G}} \leftarrow \{\mathbf{z}_{\mathbf{E}}\}$ ;  
10: Constructing latent graph  $\mathbf{G}(\mathbf{V}_{\mathbf{G}}, \mathbf{E}_{\mathbf{G}}, \mathbf{A}, \mathbf{D})$ ;  
*# stage two: graph encoding*  
11: **for**  $i$ -th graph convolution  $C_{enc}^{(l)}$  in  $\mathcal{E}_L$  **do**  
12:   **if**  $i \neq 0$  **then**  
13:      $\hat{\mathbf{G}}^{(l)} \leftarrow C_{enc}^{(l)}(\hat{\mathbf{G}}^{(l-1)})$   
14:   **else**  
15:      $\hat{\mathbf{G}}^{(1)} \leftarrow C_{enc}^{(1)}(\mathbf{G})$   
16:   **end if**  
17: **end for**  
18:  $\mathbf{G}_o = \hat{\mathbf{G}}_e^{(L)} \leftarrow \mathcal{E}_P(\mathbf{G}_{\mathbf{P}})$   
19: Readout representation  $\mathbf{O}_{\mathbf{G}} \leftarrow r_{\mathbf{G}}(\mathbf{G}_o)$

---

from node or edge embeddings. That includes *Mean Pooling*, *Max Pooling*, *Sum Pooling* and *Attention-based Readout*, etc. This is essential for tasks like graph classification and regression, as it combines local information into a global context, allowing the model to understand the entire graph. And then we use the final global  $\mathbf{O}_{\mathbf{G}}$  to perform downstream tasks.

#### Downstream tasks

- Histopathological image classification involves the analysis of microscopic images of tissue samples to identify and classify various diseases, including cancer and other medical conditions. We use the global representation  $\mathbf{O}_{\mathbf{G}}$  obtained from the model and optimize the classification layer  $\text{MLP}_c$  to predict the probabilities of the classes and to supervise it with the true labels. The criterion cross-entropy loss is defined as:

$$\mathcal{L}_{CE} = -\frac{1}{B} \sum_{i=1}^B \sum_{j=1}^C y_{ij} \log(\hat{y}_{ij}) \quad (13)$$

where  $Y = \{y_{ij}\}, \hat{Y} = \{\hat{y}_{ij}\}, \hat{y}_i = \text{MLP}_c(\mathbf{O}_{\mathbf{G}})$ ,  $B$  is the number of one batch,  $C$  is the number of classes.  $\mathbf{Y}$  is the one-hot label matrix for optimization.

- The Cox proportional hazards (PH) regression model is a class of survival models in statistics. Survival models relate the time that passes, before some event occurs, to one or more covariates that may be associated with that quantity of time. In a proportional hazards model, the unique effect of a unit increase in a covariate is multiplicative with respect to the hazard rate. As consisting of two parts: (1) the underlying baseline hazard function, often denoted  $\lambda_0(t)$  describing how the risk of event per time unit changes over time at baseline levels of covariates; (2) and the effect parameters, describing how the hazard varies in response to explanatory covariates.
- Kaplan-Meier survival curves and the Log-Rank test. To estimate survival probability at a given time, the Kaplan-Meier method utilizes the risk set, incorporating data on censored individuals up to the point of censorship rather than discarding it.

## B. Introduction of Comparison

### B.1 WSI and Patch Preprocessing

In this part, we follow CLAM (Lu et al. 2021) to conduct WSI processing with three stages: segmentation, patching, and label allocation for patches.

#### Segmentation

The pipeline begins with the automated segmentation of foreground tissue for each digitized slide. The whole slide image (WSI) is initially read at a reduced resolution and converted from the RGB to the HSV color space. A binary mask delineating tissue areas is generated by applying a threshold to the saturation channel after performing median blurring, followed by morphological closing to bridge any gaps. Contours detected in this process are filtered based on their area and then stored for subsequent analysis. The segmentation mask can be visually reviewed, and a text file is produced that lists the processed slides and key parameters used during segmentation.

**Patching** After segmentation, the algorithm systematically extracts patches of size  $h_I \times w_I$  from the identified tissue regions at the desired magnification. These patches, along with their spatial coordinates and related slide metadata, are saved in the HDF5 hierarchical data format. Depending on the dimensions of the WSI and the selected magnification, the number of patches per slide can range from a few hundred to several hundred thousand.

**Allocation** The pathology images derived from patching are either sent to pathologists for annotation or are assigned task-specific labels based on pre-existing annotations of the WSI. A script is employed to select homogeneous tissue areas, with an emphasis on avoiding non-structured regions as much as possible.

### B.2 More descriptions about datasets

The proposed framework undergoes pretraining and classification using three large extensive histopathology datasets. For the survival analysis task, we used the pre-trained

weights from the pan-cancer dataset Komura et al. and compared the performances of different models on two public TCGA cohorts and one privately collected primary-metastatic pathology colorectal dataset.

**Komura et al.** (Komura et al. 2022) comprises 1.6M image patches of H&E stained histological samples representing 32 solid cancer types sourced from the GDC legacy database. These images, extracted from 8,736 diagnostic slides of 7,175 patients, were randomly cropped at 6 magnification levels, the size ranging from  $128 \times 128$  to  $256 \times 256 \mu m$ , within annotated regions by trained pathologists (Komura et al. 2022).

**Prostate cANcer graDe Assessment (PANDA)** dataset (Bulten et al. 2022) includes approximately 11,000 whole-slide images of digitized H&E-stained biopsies. This dataset incorporates annotations for five grades of Gleason scores (Humphrey 2004) provided by the Radboud University Medical Center and the Karolinska Institute. We sampled patches according to Radboud’s Gleason area annotation. Valid annotation values are 0: background (nontissue) or unknown; 1: stroma (connective tissue, non-epithelium tissue); 2: healthy (benign) epithelium; 3: cancerous epithelium (Gleason 3); 4: cancerous epithelium (Gleason 4); 5: cancerous epithelium (Gleason 5). After processing, 12.5M images are generated.

**IBD** On the privately collected pathological colorectal cancer data set of 23M patches at 10x magnification. And 360K patches are annotated for 9 common tissue types(0: background; 1:tumor; 2:necrotic; 3:mucous; 4:fat; 5:smooth muscle; 6:thrombus; 7:neplasm; 8:normal colon mucosa) for downstream classification tasks. Note that non-pathological section patches will be removed. All data are pre-trained with data augmentation such as random scaling and rotation.

**TCGA-KIRC** The Cancer Genome Atlas Kidney Renal Clear Cell Carcinoma (TCGA-KIRC, 512 histopathology cases) data collection is part of a larger effort to build a research community focused on connecting cancer phenotypes to genotypes by providing clinical images matched to subjects from The Cancer Genome Atlas (TCGA).

**TCGA-ESCA** The Cancer Genome Atlas Esophageal Carcinoma (TCGA-ESCA, 155 histopathology cases) is part of the TCGA project and is dedicated to studying the genomic characteristics of esophageal cancer.

**CRC-PM** CRC-PM is a private dataset consisting of three batches of pathological images and survival information of rectal cancer and colon cancer (388 histopathology cases), with a censorship rate of  $\geq 80\%$ .

All pathological sections WSI are normalized at 10x, and divided into patches. Features are extracted from different pre-trained models for survival analysis experiments.

### B.3 Baseline Methods

**SimCLR** SimCLR (A Simple Framework for Contrastive Learning of Visual Representations) is an unsupervised method that applies random augmentations to each image to create two views encoded into feature vectors by a neural

network. These vectors are then transformed into projection vectors. The Normalized Temperature-scaled Cross Entropy (NT-Xent) loss is used to bring similar images closer and separate dissimilar ones. For a positive pair  $(i, j)$ , the loss function is:

$$l_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}. \quad (14)$$

where  $\text{sim}(\mathbf{z}_i, \mathbf{z}_j)$  is the cosine similarity,  $\tau$  is the temperature parameter,  $N$  is the batch size, and  $\mathbb{1}_{[k \neq i]}$  ensures the positive pair is excluded from the denominator.

**KimiaNet** KimiaNet utilizes the DenseNet-121 architecture with four dense blocks, which have been fine-tuned and trained using histopathology images across various configurations.

**Dino** DINO is a label-free knowledge distillation method where the student network learns from the teacher’s output using cross-entropy loss  $\mathcal{L}_{CE}$ . The teacher’s parameters  $\theta_t$  are updated using the student’s parameters  $\theta_s$  with momentum:

$$\theta_t \leftarrow m\theta_t + (1 - m)\theta_s. \quad (15)$$

To avoid collapse, the method includes prediction centering and sharpening of image  $\mathbf{I}$ :

$$\tilde{p}_t = p_t - c, p_t(\mathbf{I}) = \frac{p_t(i)^{1/\tau}}{\sum p_t^{1/\tau}}. \quad (16)$$

DINO works with both CNNs and ViTs without changing the original architecture.

**MAE** MAE (Masked Autoencoder) is a simple, effective, and scalable approach for visual representation learning. MAE applies the concept of masked autoencoding, inspired by the success of BERT in NLP, where random patches of an image are masked and then reconstructed in pixel space. The method features an asymmetric encoder-decoder design, where the encoder processes only a subset of visible patches, and a lightweight decoder reconstructs the input from the latent representation and mask tokens. This design allows for a high masking rate, significantly reducing redundancy and computational costs, while still learning robust representations.

**GraphMAE** GraphMAE is a graph neural network model based on an autoencoder. It masks node features during encoding and uses a GNN to create embeddings. In decoding, it remasks the same nodes to improve feature learning. The model learns meaningful information through unsupervised feature reconstruction. The node feature  $\tilde{\mathbf{x}}_i$  for  $\mathbf{v}_i \in \mathcal{V}$  in the masked feature matrix  $\tilde{\mathbf{x}}$  can be defined as:

$$\tilde{\mathbf{x}}_i = \begin{cases} \mathbf{x}_{[M]} & \mathbf{v}_i \in \tilde{\mathcal{V}}, \\ \mathbf{x}_i & \mathbf{v}_i \notin \tilde{\mathcal{V}}. \end{cases} \quad (17)$$

The objective of GraphMAE is to reconstruct the masked features of nodes in the subset  $\mathcal{V}$  given the partially observed node signals  $\tilde{\mathbf{X}}$  and the input adjacency matrix.

**GraphMAE2** GraphMAE2 addresses inaccuracies in node feature semantics by enhancing the masked prediction approach. It introduces regularization in the decoding stage by repeatedly re-masking the encoded representations, forcing the decoder to reconstruct features from corrupted data. Additionally, it predicts masked node representations in an embedding space distinct from the input feature space to reduce direct input influence. It employs the scaled cosine error to measure the reconstruction error and sum over the errors of the multi views  $K$  for training:

$$\mathcal{L}_{input} = \frac{1}{|\tilde{\mathcal{V}}|} \sum_{j=1}^K \sum_{\mathbf{v}_i \in \tilde{\mathcal{V}}} \left(1 - \frac{\mathbf{x}_i^T \mathbf{z}_i^{(j)}}{\|\mathbf{x}_i\| \cdot \|\mathbf{z}_i^{(j)}\|}\right)^\gamma. \quad (18)$$

where  $\mathbf{z}_i^{(j)}$  is the  $i$ -th row of predicted feature, and  $\gamma \geq 1$  is the scaled coefficient. It also learns the parameters  $\theta$  of the encoder and projector by minimizing the following scaled cosine error with gradient descent:

$$\mathcal{L}_{latent} = \left(1 - \frac{\tilde{\mathbf{z}}_i^T \tilde{\mathbf{x}}_i}{\|\tilde{\mathbf{z}}\| \cdot \|\tilde{\mathbf{x}}\|}\right)^\gamma. \quad (19)$$

And the parameters of the target generator  $\zeta$  are updated via an exponential moving average of  $\theta$  using weight decay  $\tau$ :  $\zeta = \tau\zeta + (1 - \tau)\theta$ . This approach prevents overfitting to input features and enhances the model’s generalization capabilities.

**DiffAE** DiffAE is a novel approach that leverages diffusion probabilistic models (DPM) for image attribute editing by combining the strengths of DPMs with autoencoder architectures. The key idea behind DiffAE is to use DPMs as decoders within an autoencoder framework, where the autoencoder’s encoder outputs a semantically meaningful latent code, denoted as  $\mathbf{z}_{sem}$ . This latent code, analogous to those in GANs or VAEs, captures high-level semantic features of the input image.

**DiffMAE** DiffMAE is an innovative method that rethinks the role of generative models in pre-training visual representations by integrating principles from denoising diffusion models. In DiffMAE, diffusion models are restructured into a masked autoencoder framework, where the input data is masked. It generates masked regions by sampling from  $p(\mathbf{x}_0^m | \mathbf{x}_0^v)$ , which is approximated by recursively sampling from  $p(\mathbf{x}_{t-1}^m | \mathbf{x}_t^m, \mathbf{x}_0^v)$  can also considered Gaussian distributed. We optimize the simple objective proposed by DDPM:

$$\mathcal{L} = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \|\mathbf{x}_0^m - D_\theta(\mathbf{x}_t^m, t, E_\phi(\mathbf{x}_0^v))\|^2. \quad (20)$$

## B.4 Downstream Backbones

**DeepSurv** DeepSurv uses deep learning to enhance the nonlinear Cox proportional hazards model through a deep linear feed-forward network. It predicts how patient covariates affect their hazard rate by learning the weights  $\beta$  in the linear risk function  $\hat{h}_\beta(\mathbf{x}) = \beta^T \mathbf{x}$ . The model optimizes  $\beta$  using the Cox partial likelihood:

$$\mathcal{L}_{cox}(\beta) = \prod_{i: E_i=1} \frac{\exp(\hat{h}_\beta(\mathbf{x}_i))}{\sum_{j \in \mathcal{R}(T_i)} \exp(\hat{h}_\beta(\mathbf{x}_j))}. \quad (21)$$

Here,  $T_i, E_i$ , and  $x_i$  are the respective event time, event indicator and covariates for the  $i$ -th observation. The risk set  $\mathbb{R}(t) = \{i : T_i \geq t\}$  is the set of patients still at risk of failure at time  $t$ .

**AB-MIL** Attention-based Deep Multiple Instance Learning proposes to use a weighted average of instances (low-dimensional embeddings) where weights are determined by a neural network. Additionally, the weights must sum to 1 to be invariant to the size of a bag. And MIL pooling as follows:

$$z = \sum_{k=1}^K \mathbf{a}_k \mathbf{h}_k \quad (22)$$

where:

$$\mathbf{a}_k = \frac{\exp(\mathbf{w}^T \tanh(\mathbf{V} \mathbf{h}_k^T))}{\sum_{j=1}^K \exp(\mathbf{w}^T \tanh(\mathbf{V} \mathbf{h}_j^T))}. \quad (23)$$

where  $\mathbf{w} \in \mathbb{R}^{L \times 1}$  and  $\mathbf{V} \in \mathbb{R}^{L \times M}$  are parameters. Moreover, we utilize the  $\tanh(\cdot)$  element-wise non-linearity to include both negative and positive values for proper gradient flow.

**PatchGCN** Patch-based Graph Convolutional Network presents a context-aware, spatially-resolved patch-based graph convolutional network that hierarchically aggregates instance-level histology features to model local and global-level topological structures in the tumor microenvironment. It makes  $F_{\text{GCN}}^{(l)}$  a residual mapping and stack multiple layers of  $F_{\text{GCN}}^{(l)}$  where the output of  $F_{\text{GCN}}^{(l)}$  additively combines with its input. The attention pooling of instance-level features is performed in local graph neighborhoods instead of across the entire bag.

$$\mathbf{G}^{(l+1)} = F_{\text{GCN}}^{(l)}(\mathbf{G}^{(l)}; \phi^{(l)}, \rho^{(l)}, \zeta^{(l)}) + \mathbf{G}^{(l)}. \quad (24)$$

in which  $\phi^{(l)}$  is the additively combined node and edge features followed by ReLU activation,  $\rho^{(l)}$  is a Softmax Aggregation scheme.  $\zeta^{(l)}(\mathbf{h}_v^{(l)}, \mathbf{m}_v^{(l)}) = MLP(\mathbf{h}_v^{(l)} + \mathbf{m}_v^{(l)}) \rightarrow \mathbf{h}_v^{(l+1)}$ .

## B.5 Evaluation Metric

**Accuracy** In the classification tasks, Accuracy measures how close a given set of observations are to their true value. Accuracy is the proportion of correctly classified instances (both positive and negative) out of the total number of instances.

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}, \quad (25)$$

where:

- **TP**: Number of correctly predicted positive instance.
- **TN**: Number of correctly predicted negative instance.
- **FP**: Number of incorrectly predicted positive instance.
- **FN**: Number of incorrectly predicted negative instance.

Accuracy measures the overall effectiveness of the model by determining the ratio of all correct predictions (both positive and negative) to the total number of cases.

**F1-Score** The F1-Score is the harmonic mean of Precision and Recall, providing a single metric that balances the trade-off between these two. It is beneficial when dealing with imbalanced datasets, as it considers both false positives and false negatives.

$$\text{F1} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (26)$$

where:

- **Precision** =  $\frac{\text{TP}}{\text{TP} + \text{FP}}$ : The proportion of correctly predicted positive instances out of all predicted positive instances.
- **Recall** =  $\frac{\text{TP}}{\text{TP} + \text{FN}}$ : The proportion of correctly predicted positive instances out of all actual positive instances.

**RMSE** The Root Mean Square Error is a commonly used metric for evaluating the reconstruction prediction error of a model. The formula for calculating RMSE is as follow:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}, \quad (27)$$

where  $n$  is the number of samples.  $y_i$  is the actual value of the  $i$ -th sample.  $\hat{y}_i$  is the predicted value of the  $i$ -th sample. The lower the RMSE, the better the model's predictive reconstruction performance.

**C-Index** The concordance index(C-index, CI) is the most frequently used survival model evaluation metric. It is a measure of rank correlation between predicted risk scores  $\hat{h}$  and observed time points  $T$  with event indicator  $\delta$  (1: death, recurrence, 0: not occurred). It is defined as the ratio of correctly ordered(concordant) pairs to comparable pairs.

$$\text{C-index} = \frac{\sum_{i,j} 1_{T_j < T_i} \cdot 1_{\hat{h}_j > \hat{h}_i} \cdot \delta_j}{\sum_{i,j} 1_{T_j < T_i} \cdot \delta_j}. \quad (28)$$

Two samples  $i$  and  $j$  are comparable if the sample with lower observed time  $T$  experienced an event, i.e., if  $T_j > T_i$  and  $\delta_i = 1$ . A comparable pair  $(i, j)$  is concordant if the estimated risk  $\hat{h}$  by a survival model is higher for subjects with lower survival time, i.e.,  $\hat{h}_i > \hat{h}_j \wedge y_j > y_i$ , otherwise the pair is discordant.

## B.6 Experiment Environment and Dependencies

All experiments are conducted on a workstation with 8 NVIDIA GeForce RTX 3090 (24 GB) GPUs equipped with 256 GB memory. Entity graph construction is implemented by Histocartography (Jaume et al. 2021a) equipped with SLIC (compactness is 10, blur kernel size is 1, threshold is 0.02). Our graph convolutional model is implemented by Pytorch Geometric (Fey and Lenssen 2019).

## C. Visualization and Interpretability

### C.1 More Kaplan-Meier analysis

We added the KM curve of the CRC-PM dataset and other backbones's result of *ESCA* and *KIRC* in Fig. 7.

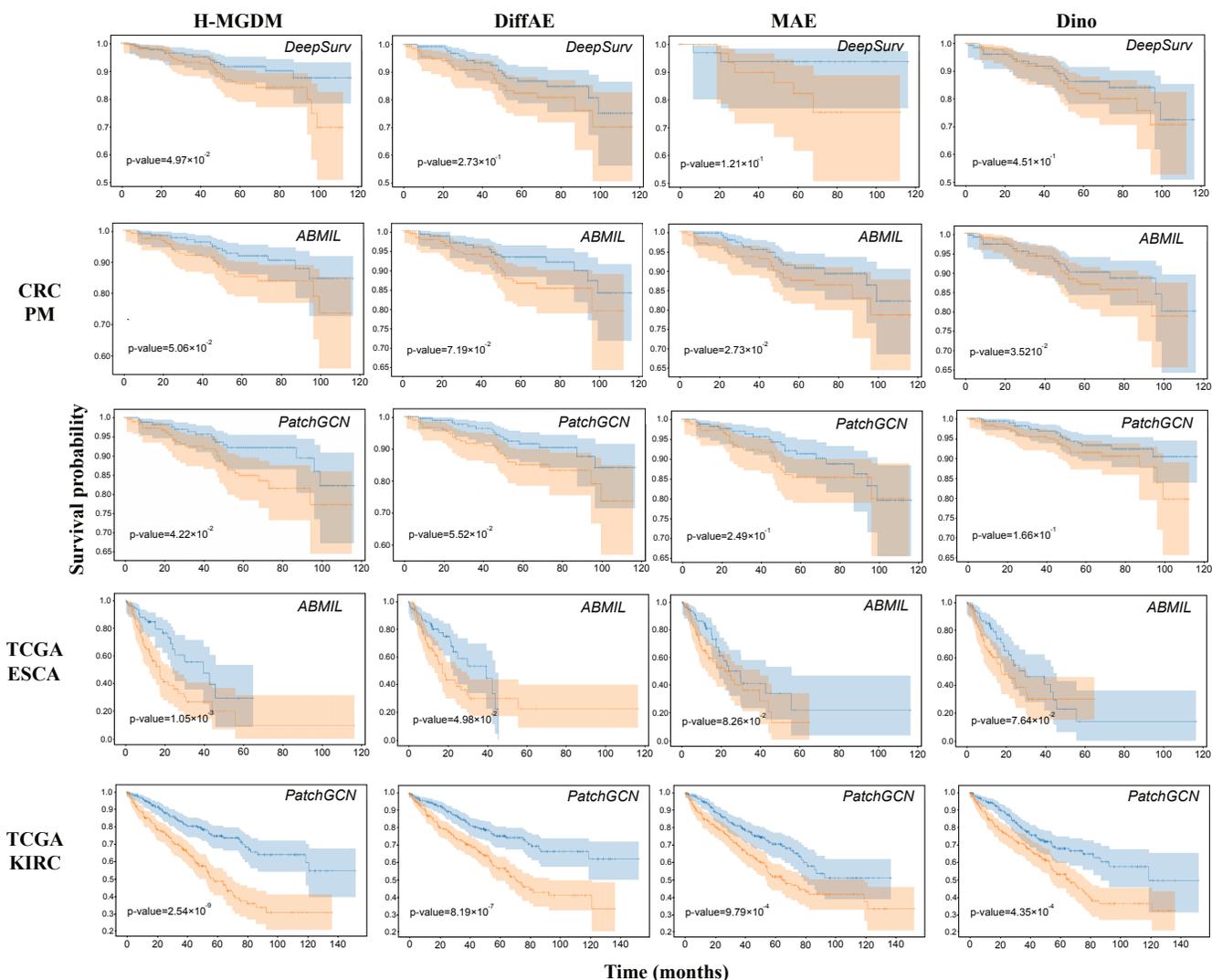


Figure 7: Kaplan-Meier Analysis of Comparison methods and our framework for three datasets on different backbones. Each cohort was split into a high-risk group (orange curve) and a low-risk group (blue curve) according to the median output of the prediction model.

## C.2 More Patch level Interpretability of representation

More patch-level entities interpretability are shown on Fig 8. Our H-MGDM can more effectively focus on pathological entities related to the category and automatically ignore non-effective areas such as the background.

## D. Limitations and Future Work

### D.1 Non-tile-level Entity Extraction

For each pixel  $s$ , a window of size  $a \times a$  centered at  $s$  is considered a vertex  $v$  in the pathological entity graph  $P$  in pixel space. Pixels within the window that do not belong to  $s$  are assigned the background color, which can cause information redundancy and interference during superpixel feature extraction due to the excess background information. Additionally, the window size is influenced by the size of the

superpixel: if the window is too small, it may not fully capture the superpixel's information, while if it is too large, it may include too much background, leaving less space for semantically meaningful parts. Therefore, a specialized feature extraction model for superpixels is urgently needed.

### D.2 Limited Amount of Data

The data available for pretraining our model is currently limited, which poses a challenge to achieving optimal performance. Expanding the dataset with more extensive collections, especially those that encompass a broader range of cancer types, would significantly enhance both the training process and the overall effectiveness of the model. A larger and more diverse dataset would provide the model with richer information, allowing it to learn more robust patterns and generalize better across different cancer types, ul-

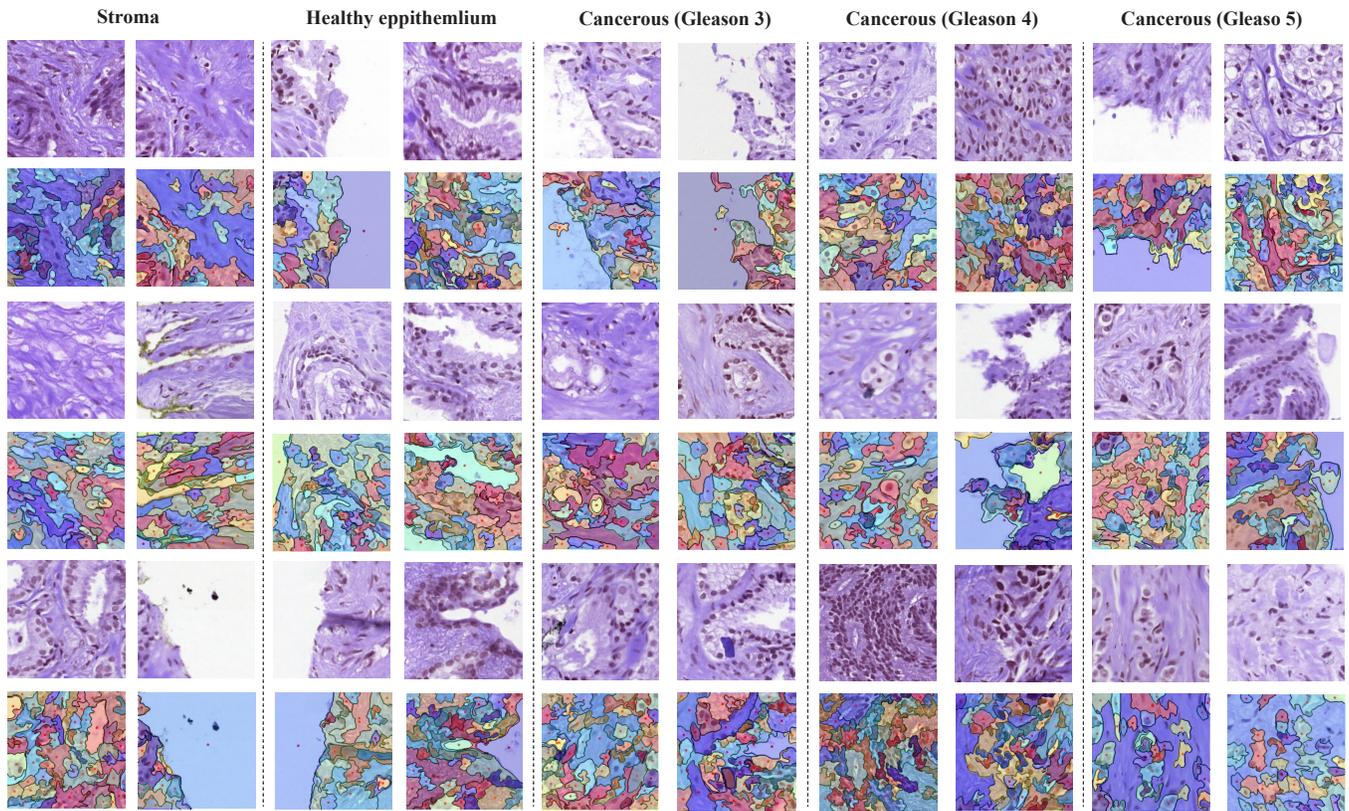


Figure 8: More Histocartographical Image Interpretability

timely leading to improved predictive accuracy and reliability in clinical applications.