

# ResGAT: A Residual Graph Attention Network for Cancer Subtype Classification in Whole Slide Images

Zhenhan Lin<sup>1</sup>

Hao Tong<sup>3</sup>

Yunfei Hu<sup>1</sup>

Xianyong Gui<sup>4</sup>

Jeanne Shen<sup>5</sup>

Byrne Lee<sup>6</sup>

Lu Zhang<sup>7</sup>

Daniel Moyer<sup>1</sup>

Mu Zhou<sup>8</sup>

Xin Maizie Zhou<sup>1,2,\*</sup>

Konstantinos Votanopoulos<sup>3,\*</sup>

ZHENHAN.LIN@VANDERBILT.EDU

TONGH@ALUMNI.WFU.EDU

YUNFEI.HU@VANDERBILT.EDU

XIANYONG.GUI@WFUSM.EDU

JEANNES@STANFORD.EDU

BYRNELEE@STANFORD.EDU

ERICLUZHANG@HKBU.EDU.HK

DANIEL.MOYER@VANDERBILT.EDU

MUZHOU1@GMAIL.COM

MAIZIE.ZHOU@VANDERBILT.EDU

KVOTANOP@WAKEHEALTH.EDU

<sup>1</sup> Department of Computer Science, Vanderbilt University, Nashville, TN, United States

<sup>2</sup> Department of Biomedical Engineering, Vanderbilt University, Nashville, TN, United States

<sup>3</sup> Department of General Surgery, Wake Forest University, Winston-Salem, NC, United States

<sup>4</sup> Department of Pathology, Wake Forest University, Winston-Salem, NC, United States

<sup>5</sup> Department of Pathology, Stanford University School of Medicine, Palo Alto, CA, United States

<sup>6</sup> Department of Surgery, Stanford University, Palo Alto, CA, United States

<sup>7</sup> Department of Computer Science, Hong Kong Baptist University, Hong Kong

<sup>8</sup> Department of Computer Science, Rutgers University, New Brunswick, NJ, United States

**Editors:** Under Review for MIDL 2026

## Abstract

Multiple instance learning (MIL) provides a weakly supervised framework for whole slide image (WSI) classification, enabling slide-level prediction from gigapixel images with only slide-level labels. However, WSI subtype classification in realistic settings is still challenging. In this work, we propose ResGAT, a residual graph attention framework that operates on hybrid  $k$ -NN patch graphs and models WSI representations with stacked residual graph attention blocks. ResGAT is evaluated on the subtype classification task across a rare, class-imbalanced appendiceal cancer cohort, BRACS and two TCGA datasets. It outperforms SOTA MIL baselines on the appendiceal cancer and BRACS cohorts, and remains competitive on the TCGA datasets. On the appendiceal cancer cohort, we further assess cross-site generalization via few-shot adaptation under source shift, showing that ResGAT adapts effectively to new domains with limited labels. An ablation study is provided to validate the effectiveness of key architectural components of our method.

**Keywords:** whole slide image classification, multiple instance learning, residual graph attention framework, cross-site generalization

## 1. Introduction

As histopathology digitization becomes routine, incorporating computational models into diagnostic workflows is increasingly feasible (Hanna et al., 2019; Kumar et al., 2020;

Zhang et al., 2025). These computational models provide slide-level classification results together with interpretable justifications, promoting consistent decisions and transparent verification (Tizhoosh and Pantanowitz, 2018; Yilmaz et al., 2024). This is particularly valuable for rare diseases, where expert diagnosticians are scarce. However, a fundamental challenge lies in the gigapixel scale of whole-slide images (WSIs), which prevents them from being processed as a single image. In practice, the standard approach involves tiling tissue regions into thousands of patches, formulating the task as a Multiple Instance Learning (MIL) problem.

The evolution of MIL for WSI classification has shifted from simple feature pooling to sophisticated context modeling. Initial frameworks adopted static aggregation strategies, such as max-pooling (Campanella et al., 2019) and mean-pooling. While computationally efficient, these methods often lose critical contextual information by focusing only on the extreme feature or diluting signals through averaging. The introduction of Attention-based MIL (ABMIL) (Ilse et al., 2018) marked a pivotal advancement by using trainable weights to rank instances. Subsequent research has sought to address overfitting and attention concentration through advanced strategies: pseudo-bag augmentation and feature distillation methods like DTFD-MIL (Zhang et al., 2022); and attention-challenging frameworks such as ACMIL (Zhang et al., 2024) and MHIM (Tang et al., 2023) that mitigate attention concentration by suppressing high-confidence instances to encourage the discovery of comprehensive diagnostic patterns. Despite these improvements, the attention mechanisms often treat instances as independent and identically distributed (i.i.d.). To explicitly capture inter-instance correlations, recent sequence-based works like TransMIL (Shao et al., 2021) and the Mamba-based architecture (Yang et al., 2024) leverage self-attention and selective scan mechanisms to explicitly model long-range dependencies, marking a paradigm shift towards correlated feature learning.

Running parallel to sequence-based advancements, Graph Neural Networks (GNNs) have emerged as a distinct paradigm focused on explicitly encoding the structural topology of the tissue (Brussee et al., 2025). By representing patches as nodes and their interactions as edges, these methods avoid flattening the spatial structure into a sequence. Early implementations employed  $k$ -nearest neighbor ( $k$ NN) algorithms to construct spatial graphs, demonstrating that explicitly modeling local neighborhoods enhances diagnostic accuracy (Chen et al., 2021; Zheng et al., 2022). Subsequent research has explored more intricate graph constructions, including hierarchical formulations for multi-resolution reasoning (Hou et al., 2022) and heterogeneous graphs that distinguish between different tissue components (Chan et al., 2023). However, the "over-smoothing" phenomenon (Chen et al., 2020) is challenging for graph-based MIL approaches. Stacking multiple message passing layers induces node representations to become homogenized, losing the discriminative power essential for classification. This degradation poses an obstacle in realistic clinical settings, which are characterized by extreme heterogeneity in tissue scale. In such diverse scenarios, the fact that applying standard readout functions to homogenized features yields inconsistent diagnostic profiles across varying graph sizes, harming the reliability required for clinical deployment.

Motivated by these challenges, we propose ResGAT, a weakly supervised MIL framework for whole slide image subtype classification. The whole slide image is represented as a hybrid  $k$ -NN patch graph with nodes initialized by extracted patch features and connected

via spatial and feature proximity. ResGAT processes the patch graphs with stacked residual graph attention blocks, where each block features a dual-branch design combining multi-head graph attention with a parallel linear projection. This design preserves patch-specific information while adaptively aggregating contextual information, yielding representations that support effective slide-level prediction. In comparative evaluations against representative MIL baselines, our model achieves superior classification performance on both a rare, class-imbalanced appendiceal cancer cohort and the multi-class BRACS dataset, while it remains competitive on two public TCGA datasets. On the appendiceal cancer cohort, we also introduce a benchmarking protocol to assess cross-site generalization and few-shot adaptation, demonstrating that ResGAT maintains strong performance when labeled data are limited in new domains. An ablation study is provided to examine the effectiveness of the core components of ResGAT. Furthermore, the framework supports qualitative interpretation through heatmaps that highlight prediction-relevant regions.

## 2. Method

### 2.1. Problem Formulation

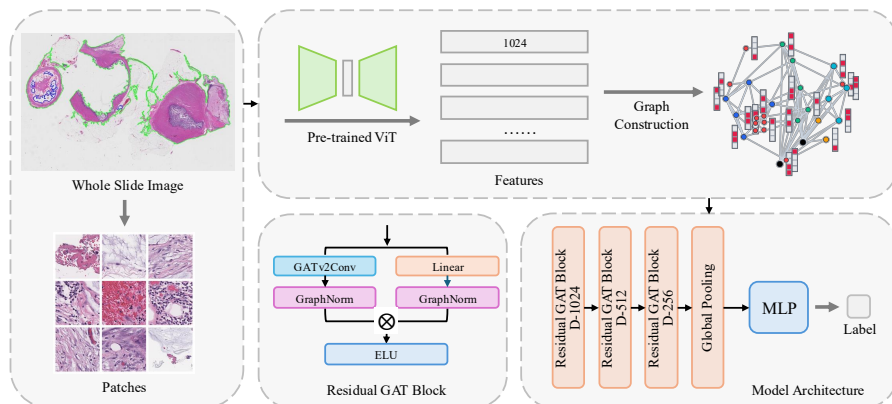


Figure 1: Overview of the ResGAT pipeline for WSI classification. The framework consists of three main components: (1) tissue segmentation and patch extraction, (2) patch-level feature encoding and graph construction, and (3) slide-level representation learning and prediction.

The whole slide image is treated as a bag of patch embeddings in the multiple instance learning (MIL) setting. Given a slide  $s$ , the foreground tissue is segmented and tiled into patches at a fixed magnification. Each patch is then encoded into a feature vector  $\mathbf{x}_i \in \mathbb{R}^D$  using a large-scale pre-trained pathology encoder. This yields a set  $\mathcal{B}_s = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  with a slide-level label  $y_s$  indicating the cancer subtype. Our goal is to learn a permutation-invariant function  $f_\theta : \mathcal{B}_s \mapsto y_s$  for subtype classification.

Following previous graph-based MIL methods, we represent each slide as a patch graph  $\mathcal{G}_s = (\mathcal{V}_s, \mathcal{E}_s)$ . Each node  $v_i \in \mathcal{V}_s$  corresponds to a patch embedding  $\mathbf{x}_i$ , and the edges in  $\mathcal{E}_s$

are constructed based on both spatial proximity and feature similarity between the patches. ResGAT takes the graph as input, updates node features with stacked residual graph attention blocks, and aggregates them into a slide-level representation for classification. Fig. 1 shows the overall architecture of ResGAT.

## 2.2. Graph Construction

To establish the graph topology  $\mathcal{E}_s$ , we introduce a hybrid  $k$ -NN edge construction procedure. Each node  $v_i$  is associated with a spatial coordinate  $\mathbf{p}_i \in \mathbb{R}^2$  derived from the patch location on the WSI. Initially, we identify the  $d_{spa}$  nearest spatial neighbors of  $v_i$  measured by Euclidean distance between coordinates, denoted as the set  $\mathcal{N}_{spa}(v_i)$ , and its  $d_{feat}$  nearest feature neighbors measured by cosine distance between node features, denoted as the set  $\mathcal{N}_{feat}(v_i)$ . We define the candidate pool as the intersection

$$\mathcal{C}(v_i) = \mathcal{N}_{spa}(v_i) \cap \mathcal{N}_{feat}(v_i),$$

which is subsequently ranked by node feature similarity. The top  $k$  candidates are selected as the final connected neighbors of  $v_i$ . In cases of a sparse or empty intersection ( $|\mathcal{C}(v_i)| < k$ ), the adjacency list is padded with up to three auxiliary nearest feature neighbors to maintain robust connectivity. The resulting patch graph  $\mathcal{G}_s$  is treated as undirected. The hyperparameters  $d_{spa}, d_{feat}, k$  jointly determine the graph density and the node degree variance. We adopt a general configuration with  $k = 6, d_{feat} = 50, d_{spa} \in \{15, 24\}$  in our main evaluations. A comprehensive sensitivity analysis of these parameters is provided in Section 3.4.1.

## 2.3. ResGAT Architecture and Training Objective

**Node Updates.** Given the graph  $\mathcal{G}_s$ , we initialize the node features as  $\mathbf{h}_i^{(0)} = \mathbf{x}_i$  for  $i = 1, \dots, N$ . Let  $\mathbf{h}_i^{(\ell)}$  denote the representation of node  $v_i$  at layer  $\ell$ . ResGAT applies a stack of  $L = 3$  residual blocks to obtain the updated node representations  $\mathbf{h}_i^{(L)}$ . Each residual block updates node features through a linear projection in parallel with a multi-head graph attention convolution (GATv2Conv (Brody et al., 2021)). Let  $\mathcal{N}(i) = \{j \mid (i, j) \in \mathcal{E}_s\}$  denote the neighbors of node  $v_i$ . For each layer, the following combined update is applied to all nodes:

$$\begin{aligned} e_{ij}^{(k)} &= \mathbf{a}^{(k)\top} \text{LeakyReLU}(\mathbf{W}_s^{(k)} \mathbf{h}_i^{(\ell)} + \mathbf{W}_t^{(k)} \mathbf{h}_j^{(\ell)}), \quad j \in \mathcal{N}(i), \\ \alpha_{ij}^{(k)} &= \frac{\exp(e_{ij}^{(k)})}{\sum_{u \in \mathcal{N}(i)} \exp(e_{iu}^{(k)})}, \\ \mathbf{m}_i^{(\ell)} &= \left\| \sum_{k=1}^K \alpha_{ij}^{(k)} \mathbf{W}^{(k)} \mathbf{h}_j^{(\ell)} \right\|, \\ \mathbf{h}_i^{(\ell+1)} &= \phi\left(\text{GN}(\mathbf{m}_i^{(\ell)}) + \text{GN}(\mathbf{W}_{\text{res}}^{(\ell)} \mathbf{h}_i^{(\ell)})\right), \end{aligned} \tag{1}$$

where  $d_h = D_{\ell+1}/K$  is the output dimension of each attention head,  $\mathbf{W}_s^{(k)}, \mathbf{W}_t^{(k)}, \mathbf{W}^{(k)} \in \mathbb{R}^{d_h \times D_\ell}$  are learnable projections for head  $k$ ,  $\mathbf{a}^{(k)} \in \mathbb{R}^{d_h}$  is the corresponding attention vector,

and  $\mathbf{W}_{\text{res}}^{(\ell)} \in \mathbb{R}^{D_{\ell+1} \times D_{\ell}}$  is the learnable linear projection on the residual path. The operator  $\parallel$  denotes concatenation over  $K$  heads.  $\text{GN}(\cdot)$  denotes GraphNorm and is applied separately to the two branches, and  $\phi$  is the ELU non-linearity. This formulation accommodates progressively decreasing dimensions (e.g.,  $1024 \rightarrow 512 \rightarrow 256$ ).

**Graph Normalization.** Each residual block employs GraphNorm (Cai et al., 2021) to stabilize training against the severe variations in graph size and topological structure across different slides. Given the intermediate node representations at the layer  $\ell$ , GraphNorm defines the operation as

$$\mathbf{u}_i^{(\ell)} = \gamma \odot \frac{\mathbf{f}_i^{(\ell)} - \boldsymbol{\alpha} \odot \boldsymbol{\mu}^{(\ell)}}{\sqrt{(\boldsymbol{\sigma}^{(\ell)})^2 + \epsilon}} + \boldsymbol{\beta}, \quad (2)$$

where  $\boldsymbol{\mu}^{(\ell)}$  and  $(\boldsymbol{\sigma}^{(\ell)})^2$  are the mean and variance of  $\{\mathbf{f}_i^{(\ell)}\}_{i=1}^N$  over nodes in the graph, and  $\gamma, \boldsymbol{\beta}, \boldsymbol{\alpha}$  are learnable parameters shared across nodes. The operator  $\odot$  denotes element-wise multiplication. Intuitively,  $\gamma$  and  $\boldsymbol{\beta}$  provide a channel-wise affine re-parametrization of the normalized features, while  $\boldsymbol{\alpha}$  modulates the strength of graph-level centering on each feature dimension.

**Pooling and Loss.** Following residual blocks, we apply the global mean pooling over the updated node representations  $\{\mathbf{h}_i^{(L)}\}_{i=1}^N$  to obtain the slide-level representation  $\mathbf{z}_s \in \mathbb{R}^{D_L}$ . This vector is fed into an MLP classifier to produce the logit vector  $\hat{\mathbf{y}}_s \in \mathbb{R}^C$ , where  $C$  is the number of cancer subtypes. The predicted probabilities are obtained via a Softmax function. Given the ground-truth label encoded as a one-hot vector  $\mathbf{y}_s \in \{0, 1\}^C$ , we train the model using the standard cross-entropy loss:

$$\mathcal{L} = -\frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} \sum_{c=1}^C y_{s,c} \log \left( \frac{\exp(\hat{y}_{s,c})}{\sum_{c'=1}^C \exp(\hat{y}_{s,c'})} \right). \quad (3)$$

An ablation study of the two-branch residual block design is provided in Section 3.4.2.

## 2.4. Graph Class Activation Mapping

We adapt Grad-CAM++ (Chattopadhyay et al., 2018) to our graph-based pipeline to generate heatmaps that highlight prediction-relevant regions. Given a target class  $c$ , let  $h_{i,d}^{(L)}$  denote the  $d$ -th feature channel of the final node representation  $\mathbf{h}_i^{(L)} \in \mathbb{R}^{D_L}$ . We compute channel-wise importance weights  $w_d^c$  from the gradients of the class logit  $\hat{y}_c$ :

$$w_d^c = \alpha_d^c \cdot \frac{1}{N} \sum_{i=1}^N \text{ReLU} \left( \frac{\partial \hat{y}_c}{\partial h_{i,d}^{(L)}} \right), \quad \alpha_d^c = \frac{\sum_i \left( \frac{\partial \hat{y}_c}{\partial h_{i,d}^{(L)}} \right)^2}{2 \sum_i \left( \frac{\partial \hat{y}_c}{\partial h_{i,d}^{(L)}} \right)^2 + N \sum_i \left( \frac{\partial \hat{y}_c}{\partial h_{i,d}^{(L)}} \right)^3 + \epsilon}. \quad (4)$$

The saliency score  $M_i^c$  for each node is computed via a weighted combination:

$$M_i^c = \text{ReLU} \left( \sum_{d=1}^{D_L} w_d^c h_{i,d}^{(L)} \right). \quad (5)$$

These scores are min-max normalized and mapped back to the corresponding patch locations on the WSI.

### 3. Experiments

#### 3.1. Dataset and Experimental Setup

##### 3.1.1. DATASET

**Appendiceal cancer cohort.** This cohort consists of 141 diagnostic WSIs of 92 patients with low-grade appendiceal mucinous neoplasm (LAMN) and mucinous adenocarcinoma (MAC). It is significantly imbalanced (LAMN:MAC = 32:15). Sourcing from both Wake Forest and Stanford introduces additional domain shift challenges. Clinically, MAC is regarded as the more aggressive subtype with worse prognosis than LAMN, so it is treated as the positive class when computing AUC and the reported F1-score corresponds to the positive label.

**TCGA datasets.** Two public datasets, NSCLC and ESCA, were curated from The Cancer Genome Atlas (TCGA) program (Tomczak et al., 2015). Both datasets involve distinguishing adenocarcinoma from squamous cell carcinoma: LUAD vs LUSC for NSCLC, and EAC vs ESCC for the esophagus ESCA. For evaluation, the clinically more aggressive squamous cell carcinoma was treated as the positive class when computing the AUC.

**BRACS dataset.** The BRACS dataset (Brancati et al., 2022) is a public dataset with 526 WSIs of various breast lesions. It contains seven diagnostic categories: normal (N), benign lesions (PB), usual ductal hyperplasia (UDH), atypical ductal hyperplasia (ADH), flat epithelial atypia (FEA), ductal carcinoma in situ (DCIS), and invasive carcinoma (IC). Due to severe class imbalance, both the AUC and F1-score were computed using macro-averaging across all categories.

Slide counts per diagnostic category are summarized in Table 6 in the Appendix. All tissue segmentation and patch extraction were performed at 20× magnification.

##### 3.1.2. EVALUATION PROTOCOLS

All experiments were conducted on an NVIDIA RTX A6000 GPU with 48GB memory. For feature extraction, we adopted the CLAM (Lu et al., 2021) preprocessing pipeline with HSV-based tissue segmentation and contour-based spatial sampling to identify tissue regions. Patch-level 1024-dimensional feature vectors were extracted using UNI (Chen et al., 2024)(ViT-L/16 via DINOv2) with standard ImageNet normalization (Deng et al., 2009).

For each cohort, we performed slide-level 5-fold cross-validation with patient-wise splits. In each split, three folds were used for training, one for validation and one for testing. We report the mean and standard deviation of metrics over the five test folds. Specifically, balanced accuracy is reported for the highly imbalanced Appendiceal cancer and BRACS cohorts, while overall accuracy is used for the TCGA cohorts.

For the domain adaptation analysis on the appendiceal cancer cohort, WF slides formed the source domain and SF slides the target domain. The WF data were partitioned into training, validation and test subsets in a 70/15/15 ratio for pre-training each model. For the target domain, we defined a fixed SF test set of 12 slides (10 LAMN and 2 MAC); this SF test set was used for all zero-shot and few-shot evaluations. Zero-shot performance was obtained

Table 1: Subtype classification performance (mean<sub>std</sub>, %) across four datasets: appendiceal cancer, TCGA-NSCLC, TCGA-ESCA, and BRACS. Results are reported as balanced accuracy (BAcc), accuracy, AUC, F1 and F1-macro.

Method	Appendiceal Cancer			BRACS		
	BAcc	AUC	F1	BAcc	AUC	F1-macro
CLAM-SB	90.09 <sub>6.47</sub>	94.96 <sub>8.79</sub>	86.25 <sub>8.15</sub>	31.25 <sub>4.19</sub>	77.29 <sub>3.21</sub>	27.11 <sub>4.45</sub>
CLAM-MB	88.62 <sub>10.68</sub>	<b>96.82<sub>4.13</sub></b>	85.36 <sub>14.95</sub>	30.48 <sub>3.30</sub>	74.72 <sub>3.12</sub>	27.56 <sub>3.89</sub>
DSMIL	78.92 <sub>13.86</sub>	90.58 <sub>9.89</sub>	68.44 <sub>24.77</sub>	32.15 <sub>6.74</sub>	68.44 <sub>3.66</sub>	29.31 <sub>6.84</sub>
TransMIL	84.07 <sub>10.71</sub>	92.47 <sub>7.64</sub>	76.87 <sub>14.33</sub>	30.18 <sub>4.07</sub>	76.09 <sub>2.57</sub>	26.80 <sub>3.89</sub>
WiKG	84.31 <sub>7.39</sub>	94.37 <sub>6.51</sub>	79.16 <sub>11.19</sub>	28.08 <sub>1.86</sub>	70.75 <sub>2.93</sub>	22.58 <sub>2.42</sub>
PatchGCN	87.41 <sub>9.48</sub>	95.35 <sub>7.42</sub>	83.84 <sub>12.29</sub>	28.93 <sub>3.85</sub>	71.57 <sub>4.68</sub>	21.87 <sub>5.11</sub>
DTFD-MIL	86.22 <sub>9.56</sub>	93.27 <sub>11.35</sub>	80.08 <sub>13.02</sub>	26.82 <sub>3.63</sub>	73.37 <sub>2.10</sub>	21.36 <sub>3.48</sub>
MHIM-DSMIL	86.42 <sub>12.74</sub>	97.03 <sub>2.72</sub>	81.15 <sub>19.45</sub>	33.29 <sub>7.18</sub>	77.47 <sub>3.89</sub>	<b>31.64<sub>5.95</sub></b>
MHIM-TransMIL	87.49 <sub>9.45</sub>	91.59 <sub>11.69</sub>	84.94 <sub>13.47</sub>	27.16 <sub>2.17</sub>	68.89 <sub>3.55</sub>	25.32 <sub>1.53</sub>
<b>ResGAT (ours)</b>	<b>92.56<sub>6.36</sub></b>	96.41 <sub>1.94</sub>	<b>90.98<sub>7.98</sub></b>	<b>33.76<sub>6.07</sub></b>	<b>77.61<sub>0.95</sub></b>	28.74 <sub>6.08</sub>

Method	TCGA-NSCLC		TCGA-ESCA	
	Accuracy	AUC	Accuracy	AUC
CLAM-SB	<b>93.72<sub>1.72</sub></b>	<b>97.55<sub>1.44</sub></b>	<b>98.04<sub>1.60</sub></b>	99.83 <sub>0.34</sub>
CLAM-MB	92.70 <sub>1.53</sub>	97.39 <sub>1.57</sub>	96.11 <sub>3.16</sub>	<b>100.00<sub>0.00</sub></b>
DSMIL	92.29 <sub>1.40</sub>	97.08 <sub>1.53</sub>	95.42 <sub>4.45</sub>	97.72 <sub>2.65</sub>
TransMIL	92.29 <sub>2.13</sub>	97.15 <sub>0.82</sub>	93.51 <sub>4.08</sub>	99.39 <sub>0.52</sub>
WiKG	92.09 <sub>1.94</sub>	96.35 <sub>1.45</sub>	93.48 <sub>3.59</sub>	99.63 <sub>0.74</sub>
PatchGCN	93.00 <sub>2.08</sub>	97.13 <sub>1.53</sub>	92.84 <sub>2.37</sub>	99.17 <sub>1.45</sub>
DTFD-MIL	93.61 <sub>1.75</sub>	97.41 <sub>1.38</sub>	96.11 <sub>3.16</sub>	99.39 <sub>0.65</sub>
MHIM-DSMIL	92.70 <sub>1.23</sub>	97.48 <sub>1.23</sub>	94.82 <sub>4.37</sub>	98.88 <sub>1.80</sub>
MHIM-TransMIL	92.40 <sub>1.31</sub>	97.30 <sub>1.51</sub>	94.82 <sub>4.37</sub>	99.73 <sub>0.22</sub>
<b>ResGAT (ours)</b>	93.51 <sub>0.75</sub>	97.15 <sub>1.47</sub>	98.02 <sub>1.62</sub>	99.91 <sub>0.17</sub>

by applying the WF-pretrained model directly to the SF test set. For few-shot adaptation, we fine-tuned the pretrained model on small labeled SF subsets with 3, 6 and 9 training slides and separate validation sets of 3, 3 and 5 slides, respectively. After adaptation, we report overall accuracy on the SF test set and quantify adaptation efficacy using backward transfer (BWT) and forward transfer (FWT). We use overall accuracy in this setting rather than balanced accuracy to provide a clearer view of adaptation trends. BWT is defined as the change in WF test accuracy before and after fine-tuning, where large negative BWT values indicate catastrophic forgetting. FWT is computed as the improvement of SF test accuracy over the zero-shot baseline, where positive values indicate successful adaptation.

See Appendix A for more implementation details.

### 3.2. Comparison with state-of-the-art methods

We compared our method with nine strong MIL baselines that cover diverse design paradigms: attention-based pooling MIL (CLAM-SB and CLAM-MB (Lu et al., 2021)), transformer-based MIL (TransMIL (Shao et al., 2021)), dual-stream MIL (DSMIL (Li et al., 2021)), distillation-based MIL (DTFD-MIL (Zhang et al., 2022)), graph-based MIL

Table 2: Domain adaptation performance comparison. Source refers to pre-trained test accuracy from WF dataset. Zero-shot refers to SF test performance on two classes data separately without adaptation. FWT measures forward transfer (target improvement), BWT measures backward transfer (source performance retention). Class 0 and Class 1 represent LAMN and MAC respectively.

Method	Source(WF)	Zero-shot (SF)		3-shot (SF)			6-shot (SF)			9-shot (SF)		
	Accuracy	class 0	class 1	Acc	FWT $\uparrow$	BWT $\uparrow$	Acc	FWT $\uparrow$	BWT $\uparrow$	Acc	FWT $\uparrow$	BWT $\uparrow$
WiKG	89.47	100	0	83.33	0	10.53	83.33	0	5.26	83.33	0	5.26
TransMIL	84.21	100	0	83.33	0	0	83.33	0	0	83.33	0	0
DSMIL	73.68	70	50	75.0	8.33	0	75.0	8.33	0	75.0	8.33	0
MHIM-DSMIL	84.21	90	0	75.0	0	5.26	75.0	0	5.26	83.33	8.33	0
MHIM-TransMIL	89.47	100	0	83.33	0	0	91.67	8.33	5.26	100	16.67	5.26
CLAM-MB	89.47	100	0	83.33	0	0	83.33	0	0	83.33	0	5.26
CLAM-SB	<b>94.74</b>	90	0	75.0	0	0	75.0	0	0	75.0	0	0
DTFD-MIL	89.47	90	100	91.67	0	0	100	8.33	0	100	8.33	5.26
PatchGCN	91.67	<b>100</b>	<b>100</b>	100	0	0	100	0	0	100	0	2.64
<b>ResGAT</b>	<b>92.86</b>	100	50	<b>91.67</b>	<b>8.33</b>	<b>0</b>	<b>100</b>	<b>8.33</b>	<b>0</b>	<b>100</b>	<b>8.33</b>	<b>0</b>

(WiKG (Li et al., 2024) and PatchGCN (Chen et al., 2021)), and hard-instance-mining MIL (MHIM-DSMIL and MHIM-TransMIL (Tang et al., 2023, 2026)). Table 1 reports mean and standard deviation over five folds for all metrics on the four datasets.

On the appendiceal cancer cohort, ResGAT achieves the highest balanced accuracy at  $92.56 \pm 6.36\%$ , outperforming the best baseline CLAM-SB by roughly 2.5% and yielding the lowest standard deviation across folds. It also attains the highest F1-score and a high AUC, indicating reliable detection of the clinically more aggressive MAC subtype. On TCGA-NSCLC and TCGA-ESCA, CLAM-SB attains the highest mean accuracy, while ResGAT remains competitive: its accuracy is only 0.21% and 0.02% below CLAM-SB on TCGA-NSCLC and TCGA-ESCA, respectively. Notably, ResGAT’s low standard deviations on TCGA cohorts shows stable performance across folds. On BRACS, a challenging seven-class fine-grained classification task, ResGAT achieves the highest balanced accuracy and AUC among all methods, while MHIM-DSMIL obtains the highest F1-macro. The overall low balanced accuracy across all methods reflects the inherent difficulty of fine-grained breast lesion subtyping. Overall, these results indicate that ResGAT performs well on the class-imbalanced and label-noisy appendiceal cancer cohort and the BRACS dataset, while remaining comparable to competitive MIL baselines on other datasets.

The results also highlight the complementary strengths of other MIL approaches. On the two TCGA cohorts, DTFD-MIL obtains the second highest accuracies and AUCs, with CLAM-MB generally close behind. The MHIM variants (MHIM-DSMIL and MHIM-TransMIL) consistently improve over their backbones, and show the effectiveness of the hard-instance mining strategy.

### 3.3. Domain Adaptation Analysis

In this experiment, we evaluated cross-site robustness on the appendiceal cancer cohort, where WF and SF correspond to different acquisition sites (see Section 3.1.1 for details). Such cross-site settings often introduce substantial distribution shift due to differences in

scanners, staining protocols and local practice, and models trained on a single site can experience a marked performance drop when deployed elsewhere (Liu et al., 2025; Pocevičiute et al., 2024). We therefore used this scenario to assess generalization ability of methods, which is a critical consideration for realistic clinical deployment. We first evaluated zero-shot performance, where a model trained on the source site is directly applied to the target site. Subsequently, we evaluate few-shot adaptation, where only a small number of labeled SF slides are available for finetuning the source-trained model (see Section 3.1.2 for details).

### 3.3.1. CROSS-DOMAIN GENERALIZATION

Table 2 compares our method with the same nine MIL baselines. While most MIL baselines achieve reasonably high accuracy on the WF source test set, their zero-shot performance on the SF target set is highly variable and often subtype-imbalanced. Several baselines, including WiKG, TransMIL and the CLAM variants, fail to correctly predict MAC samples during cross-site transfer, indicating a strong bias toward the majority subtype when crossing sites. In comparison, PatchGCN and DTFD-MIL achieve strong zero-shot performance on the SF test set, with per-class accuracies exceeding 90%, suggesting robust initial cross-site generalization. ResGAT achieves the second-highest source-domain accuracy on the WF test set and provides competitive zero-shot accuracy on the SF test set, establishing a solid foundation for further adaptation.

### 3.3.2. FEW-SHOT ADAPTATION

In this setting, we analyzed how pre-trained models adapt to target data when finetuned on a small number of labeled SF slides. ResGAT demonstrates superior adaptation efficiency, reaching 100% overall accuracy on the SF test set at the 3-shot setting and maintaining this performance across the 6-shot and 9-shot settings. Its already high source test performance remains robust across all settings ( $BWT = 0$ ), showing that adaptation does not induce forgetting on the source domain. This result suggests that ResGAT can be effectively adapted to a new site with only a small number of labeled slides, which is especially valuable in rare-disease scenarios where annotation is costly and limited.

PatchGCN maintains its perfect accuracy throughout all few-shot settings, suggesting that its graph-based representation captures site-invariant tissue structures. DTFD-MIL and MHIM-TransMIL show steady improvements on SF test accuracy as more target slides become available, alongside positive forward transfer (FWT) and BWT, indicating stable learning and knowledge retention under additional target supervision. By contrast, CLAM-SB and CLAM-MB, despite their strong performance on general classification tasks, show little change in SF test accuracy across all few-shot settings, suggesting that their architectures are less responsive to limited target supervision during cross-site adaptation.

## 3.4. Ablation Study

### 3.4.1. EFFECTIVENESS OF PROPOSED EDGE CONSTRUCTION

To assess the contribution of the proposed graph construction, we compared several topology variants: Feature kNN (edges based on feature similarity), Spatial kNN (edges

Table 3: Ablation on graph construction. Values are reported as mean<sub>std</sub> over 5-fold cross-validation (%).

Graph Variant	Appendiceal Cancer		TCGA-NSCLC		TCGA-ESCA		BRACS	
	BAcc	AUC	Accuracy	AUC	Accuracy	AUC	BAcc	AUC
Feature kNN	90.97 <sub>4.82</sub>	96.45 <sub>3.86</sub>	92.70 <sub>1.71</sub>	97.21 <sub>1.17</sub>	98.02 <sub>1.62</sub>	99.91 <sub>0.17</sub>	32.30 <sub>2.55</sub>	77.13 <sub>1.27</sub>
Spatial kNN	91.79 <sub>6.22</sub>	96.06 <sub>4.22</sub>	93.10 <sub>1.24</sub>	97.35 <sub>1.37</sub>	97.35 <sub>2.49</sub>	99.45 <sub>0.68</sub>	33.41 <sub>3.49</sub>	77.49 <sub>1.70</sub>
Hybrid ( $d_{spa}=24$ )	<b>92.56<sub>6.36</sub></b>	<b>96.41<sub>1.94</sub></b>	92.60 <sub>1.56</sub>	<b>97.65<sub>0.90</sub></b>	<b>98.02<sub>1.62</sub></b>	<b>99.91<sub>0.17</sub></b>	<b>33.76<sub>6.07</sub></b>	<b>77.61<sub>0.95</sub></b>
Hybrid ( $d_{spa}=15$ )	90.78 <sub>6.05</sub>	94.23 <sub>3.69</sub>	<b>93.51<sub>0.75</sub></b>	97.15 <sub>1.47</sub>	98.02 <sub>1.62</sub>	99.54 <sub>0.93</sub>	31.54 <sub>2.58</sub>	77.00 <sub>2.09</sub>
Node-permuted	92.46 <sub>6.14</sub>	96.08 <sub>3.53</sub>	92.80 <sub>0.99</sub>	97.70 <sub>1.00</sub>	98.02 <sub>1.62</sub>	99.83 <sub>0.21</sub>	33.53 <sub>3.98</sub>	76.95 <sub>3.98</sub>

based on spatial proximity), Hybrid (our method with two settings of  $d_{spa}$ ) and Node-permuted (hybrid adjacency with features randomly reassigned to nodes). For all graph variants, we use  $k = 6$ ; for the hybrid case, we vary the  $d_{spa}$  hyperparameter while keeping all other settings fixed (see Section 2.2 for details). As shown in Table 3, the hybrid graph consistently provides the strongest overall performance across datasets, indicating that combining spatial proximity and feature similarity yields a more effective graph topology than using either criterion alone. Notably, the node-permuted variant remains competitive. This suggests that the adjacency structure itself provides structural regularization that stabilizes representation learning and mitigates overfitting.

We further investigated the robustness of the hybrid topology through a sensitivity analysis of its hyperparameters. Specifically, we performed a grid search over the number of spatial neighbors ( $d_{spa} \in \{15, 24, 36, 48, 60\}$ ), feature neighbors ( $d_{feat} \in \{35, 50, 65, 90, 105\}$ ), and the maximum neighborhood size ( $k \in \{6, 8\}$ ). The resulting heatmap in Appendix Fig. 2 visualizes the evaluation metrics across all four datasets. ResGAT exhibits stability over a broad spectrum of parameter combinations. Although the best parameter choices vary by dataset, the general configuration consistently provides strong performance across all datasets. Overall, these results demonstrate that the proposed edge construction is effective and robust across diverse datasets.

### 3.4.2. EFFECTIVENESS OF RESIDUAL BLOCK DESIGN

We conducted a set of experiments to evaluate key architectural design choices, including dual-branch architecture, normalization strategy, layer depth, and graph convolution type. First, we compared the performance of the full ResGAT model against two variants: one ablating the linear branch and another removing all inter-node edges, which degenerates the model into a node-wise MLP. As shown in Table 4, ablating the linear branch leads to a consistent drop in performance across datasets, indicating that direct patch-level feature propagation meaningfully complements graph aggregation. Disconnecting the inter-node edges also resulted in a decline in accuracy and balanced accuracy. Together, these findings confirm that the dual-branch architecture is essential. Preserving patch-specific features and aggregating topological context are both important for forming effective slide-level representations. Table 5 shows that GraphNorm provides the most favorable performance within ResGAT compared to LayerNorm and InstanceNorm. Specifically, it outperforms both alternatives on the appendiceal cancer and BRACS datasets, where its graph-level

Table 4: Ablation on the linear branch and inter-node edges. Values are reported as mean<sub>std</sub> over 5-fold cross-validation (%).

Setting	Appendiceal Cancer		TCGA-NSCLC		TCGA-ESCA		BRACS	
	BAcc	AUC	Acc	AUC	Acc	AUC	BAcc	AUC
w/o Inter-node Edges	88.95 <sub>6.71</sub>	<b>96.52</b> <sub>3.08</sub>	93.11 <sub>0.59</sub>	<b>97.94</b> <sub>0.74</sub>	97.38 <sub>2.43</sub>	99.82 <sub>2.20</sub>	31.29 <sub>2.89</sub>	75.34 <sub>1.63</sub>
w/o Linear Branch	87.08 <sub>6.77</sub>	93.79 <sub>3.56</sub>	92.70 <sub>1.60</sub>	97.08 <sub>1.66</sub>	93.51 <sub>3.49</sub>	99.29 <sub>0.60</sub>	29.41 <sub>4.11</sub>	74.66 <sub>3.89</sub>
<b>ResGAT</b>	<b>92.56</b> <sub>6.36</sub>	96.41 <sub>1.94</sub>	<b>93.51</b> <sub>0.75</sub>	97.15 <sub>1.47</sub>	<b>98.02</b> <sub>1.62</sub>	<b>99.91</b> <sub>0.17</sub>	<b>33.76</b> <sub>6.07</sub>	<b>77.61</b> <sub>0.95</sub>

Table 5: Ablation on normalization layers in ResGAT. Values are mean<sub>std</sub> over 5-fold cross-validation (%).

Normalization	Appendiceal Cancer		TCGA-NSCLC		TCGA-ESCA		BRACS	
	BAcc	AUC	Accuracy	AUC	Accuracy	AUC	BAcc	AUC
InstanceNorm	89.23 <sub>7.70</sub>	95.84 <sub>2.19</sub>	<b>93.51</b> <sub>0.66</sub>	<b>97.18</b> <sub>1.41</sub>	98.02 <sub>1.62</sub>	99.91 <sub>0.17</sub>	33.45 <sub>2.33</sub>	76.03 <sub>1.91</sub>
LayerNorm	81.32 <sub>11.15</sub>	91.31 <sub>7.19</sub>	91.48 <sub>1.98</sub>	96.63 <sub>1.77</sub>	93.46 <sub>4.08</sub>	99.30 <sub>0.57</sub>	25.22 <sub>3.96</sub>	69.76 <sub>4.77</sub>
GraphNorm	<b>92.56</b> <sub>6.36</sub>	<b>96.41</b> <sub>1.94</sub>	93.51 <sub>0.75</sub>	97.15 <sub>1.47</sub>	<b>98.02</b> <sub>1.62</sub>	<b>99.91</b> <sub>0.17</sub>	<b>33.76</b> <sub>6.07</sub>	<b>77.61</b> <sub>0.95</sub>

normalization statistics and learnable shift parameter offer a more expressive normalization strategy than the per-node or per-feature counterparts.

Additionally, Appendix Table 7 summarizes the impact of different layer depths and graph convolution types. For the layer depth study, the 2-layer variant removes the intermediate residual block with appropriate dimension alignment, while the 4-layer variant adds an additional block that preserves the output dimension of the third layer. The results show that a 3-layer configuration provides the best trade-off between performance and computational cost: 2-layer models generally underperform due to limited receptive fields, whereas increasing the depth to 4 layers incurs higher computational cost without yielding consistent improvements. Regarding the graph convolution type, we compared GAT against GCN, GIN, and GraphSAGE. While GIN performs comparably in specific instances, we adopt GAT as the default because it achieves the highest performance across the majority of metrics and remains the most stable choice across datasets.

We further evaluated computational efficiency by comparing the throughput of ResGAT against two other graph-based methods, WiKG and Patch-GCN, under the same training protocol. As detailed in Appendix Table 8, ResGAT achieves a throughput comparable to Patch-GCN and WiKG, demonstrating that the multi-layer residual block design enhances performance without incurring significant computational overhead.

### 3.5. Qualitative Results

We applied graph-adapted Grad-CAM++ (Section 2.4) to visualize WSI heatmaps. Appendix Fig. 3 illustrates three representative MAC cases, where the primary heatmaps in the first row were computed as a confidence-weighted average of the top cross-validation models. The regions with high saliency scores are outlined in yellow. While heatmaps from individual folds exhibit spatial variation, the top-performing models show consensus in the regions

they highlight, indicating that the network captures stable diagnostic patterns. Clinical review by our pathologist confirmed that the high-scoring patches predominantly correspond to tumor and stromal tissue, where tumor morphology and its spatial relationship with the stroma inform subtyping. This indicates that ResGAT’s predictions are informed by histologically meaningful features.

#### 4. Conclusion

In this work, we propose ResGAT, a graph-based MIL framework for WSI subtype classification. The architecture features a dual-branch residual graph attention design that preserves patch-specific features while adaptively aggregating graph-based context, helping mitigate the feature homogenization commonly associated with standard message passing. Our ablation study further shows that this dual-branch structure is effective, with direct patch-level feature propagation meaningfully complementing graph aggregation. Additionally, this study reveals that the proposed hybrid kNN graph topology, together with GraphNorm and a 3-layer GAT configuration, contributes to the overall performance of ResGAT. Our main results demonstrate that ResGAT outperforms SOTA MIL baselines on the class-imbalanced, label-noisy appendiceal cancer cohort and the challenging multi-class BRACS dataset, while remaining competitive on TCGA-NSCLC and TCGA-ESCA datasets. To assess model robustness under realistic deployment conditions, we introduce a cross-site evaluation protocol on the appendiceal cancer cohort that measures zero-shot generalization and few-shot adaptation across acquisition sites. In this setting, ResGAT reaches full target-site accuracy with only a few labeled slides and without forgetting the source domain. Notably, several MIL methods that perform well in general classification task fail to adapt under the cross-site setting. These observations suggest that, while general benchmarking provides a valuable baseline, extending evaluations to realistic diagnostic settings gives a more complete picture of a model’s clinical efficacy.

## Acknowledgments

This work was supported by the NIGMS Maximizing Investigators' Research Award (MIRA) R35 GM146960.

## References

- Nadia Brancati, Anna Maria Anniciello, Pushpak Pati, Daniel Riccio, Giosuè Scognamiglio, Guillaume Jaume, Giuseppe De Pietro, Maurizio Di Bonito, Antonio Foncubierta, Gerardo Botti, et al. Bracs: A dataset for breast carcinoma subtyping in h&e histology images. *Database*, 2022:baac093, 2022.
- Shaked Brody, Uri Alon, and Eran Yahav. How attentive are graph attention networks? *arXiv preprint arXiv:2105.14491*, 2021.
- Siemen Brussee, Giorgio Buzzanca, Anne MR Schrader, and Jesper Kers. Graph neural networks in histopathology: Emerging trends and future directions. *Medical Image Analysis*, page 103444, 2025.
- Tianle Cai, Shengjie Luo, Keyulu Xu, Di He, Tie-yan Liu, and Liwei Wang. Graphnorm: A principled approach to accelerating graph neural network training. In *International Conference on Machine Learning*, pages 1204–1215. PMLR, 2021.
- Gabriele Campanella, Matthew G Hanna, Luke Geneslaw, Allen Mirafior, Vitor Werneck Krauss Silva, Klaus J Busam, Edi Brogi, Victor E Reuter, David S Klimstra, and Thomas J Fuchs. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nature medicine*, 25(8):1301–1309, 2019.
- Tsai Hor Chan, Fernando Julio Cendra, Lan Ma, Guosheng Yin, and Lequan Yu. Histopathology whole slide image analysis with heterogeneous graph representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15661–15670, 2023.
- Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE winter conference on applications of computer vision (WACV)*, pages 839–847. IEEE, 2018.
- Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph convolutional networks. In *International conference on machine learning*, pages 1725–1735. PMLR, 2020.
- Richard J Chen, Ming Y Lu, Muhammad Shaban, Chengkuan Chen, Tiffany Y Chen, Drew FK Williamson, and Faisal Mahmood. Whole slide images are 2d point clouds: Context-aware survival prediction using patch-based graph convolutional networks. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 339–349. Springer, 2021.

- Richard J Chen, Tong Ding, Ming Y Lu, Drew FK Williamson, Guillaume Jaume, Andrew H Song, Bowen Chen, Andrew Zhang, Daniel Shao, Muhammad Shaban, et al. Towards a general-purpose foundation model for computational pathology. *Nature medicine*, 30(3): 850–862, 2024.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- Matthew G Hanna, Victor E Reuter, Jennifer Samboy, Christine England, Lorraine Corsale, Samson W Fine, Narasimhan P Agaram, Evangelos Stamelos, Yukako Yagi, Meera Hameed, et al. Implementation of digital pathology offers clinical and operational increase in efficiency and cost savings. *Archives of pathology & laboratory medicine*, 143(12):1545–1555, 2019.
- Wentai Hou, Lequan Yu, Chengxuan Lin, Helong Huang, Rongshan Yu, Jing Qin, and Liansheng Wang. H<sup>2</sup>-mil: exploring hierarchical representation with heterogeneous multiple instance learning for whole slide image analysis. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 933–941, 2022.
- Maximilian Ilse, Jakub Tomczak, and Max Welling. Attention-based deep multiple instance learning. In *International conference on machine learning*, pages 2127–2136. PMLR, 2018.
- Neeta Kumar, Ruchika Gupta, and Sanjay Gupta. Whole slide imaging (wsi) in pathology: current perspectives and future directions. *Journal of digital imaging*, 33(4):1034–1040, 2020.
- Bin Li, Yin Li, and Kevin W Eliceiri. Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14318–14328, 2021.
- Jiawen Li, Yuxuan Chen, Hongbo Chu, Qiehe Sun, Tian Guan, Anjia Han, and Yonghong He. Dynamic graph representation with knowledge-aware attention for histopathology whole slide image analysis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11323–11332, 2024.
- Jingsong Liu, Han Li, Chen Yang, Michael Deutges, Ario Sadafi, Xin You, Katharina Breininger, Nassir Navab, and Peter J Schüffler. Hasd: hierarchical adaption for pathology slide-level domain-shift. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 332–342. Springer, 2025.
- Ming Y Lu, Drew FK Williamson, Tiffany Y Chen, Richard J Chen, Matteo Barbieri, and Faisal Mahmood. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature biomedical engineering*, 5(6):555–570, 2021.
- M Pocevičiute, G Eilertsen, S Garvin, and C Lundstrom. Detecting domain shift in multiple instance learning for digital pathology using fréchet domain distance. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 157–167, 2024.

- Zhuchen Shao, Hao Bian, Yang Chen, Yifeng Wang, Jian Zhang, Xiangyang Ji, et al. Transmil: Transformer based correlated multiple instance learning for whole slide image classification. *Advances in neural information processing systems*, 34:2136–2147, 2021.
- Wenhao Tang, Sheng Huang, Xiaoxian Zhang, Fengtao Zhou, Yi Zhang, and Bo Liu. Multiple instance learning framework with masked hard instance mining for whole slide image classification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4078–4087, 2023.
- Wenhao Tang, Sheng Huang, Heng Fang, Fengtao Zhou, Bo Liu, and Qingshan Liu. Multiple instance learning framework with masked hard instance mining for gigapixel histopathology image analysis. *International Journal of Computer Vision*, 134(1):41, 2026.
- Hamid Reza Tizhoosh and Liron Pantanowitz. Artificial intelligence and digital pathology: challenges and opportunities. *Journal of pathology informatics*, 9(1):38, 2018.
- Katarzyna Tomczak, Patrycja Czerwińska, and Maciej Wiznerowicz. Review the cancer genome atlas (tcga): an immeasurable source of knowledge. *Contemporary Oncology/Współczesna Onkologia*, (1):68–77, 2015.
- Shu Yang, Yihui Wang, and Hao Chen. Mambamil: Enhancing long sequence modeling with sequence reordering in computational pathology. In *International conference on medical image computing and computer-assisted intervention*, pages 296–306. Springer, 2024.
- Fazilet Yilmaz, Arlen Brickman, Fedaa Najdawi, Evgeny Yakirevich, Robert Egger, and Murray B Resnick. Advancing artificial intelligence integration into the pathology workflow: Exploring opportunities in gastrointestinal tract biopsies. *Laboratory Investigation*, 104(5):102043, 2024.
- Hongrun Zhang, Yanda Meng, Yitian Zhao, Yihong Qiao, Xiaoyun Yang, Sarah E Coupland, and Yalin Zheng. Dtf-d-mil: Double-tier feature distillation multiple instance learning for histopathology whole slide image classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18802–18812, 2022.
- Yuchen Zhang, Zeyu Gao, Kai He, Chen Li, and Rui Mao. From patches to wsis: A systematic review of deep multiple instance learning in computational pathology. *Information Fusion*, 119:103027, 2025.
- Yunlong Zhang, Honglin Li, Yunxuan Sun, Sunyi Zheng, Chenglu Zhu, and Lin Yang. Attention-challenging multiple instance learning for whole slide image classification. In *European conference on computer vision*, pages 125–143. Springer, 2024.
- Yi Zheng, Rushin H Gindra, Emily J Green, Eric J Burks, Margrit Betke, Jennifer E Beane, and Vijaya B Kolachalama. A graph-transformer for whole slide image classification. *IEEE transactions on medical imaging*, 41(11):3003–3015, 2022.

## Appendix A. Implementation Details

For ResGAT model, we trained for 30 epochs using Adam optimizer with learning rate  $3 \times 10^{-4}$  and weight decay  $1 \times 10^{-4}$ . To account for randomness, each experiment was repeated with two random seeds 3 and 3407; the best-performing run is reported. Following standard MIL practice, we applied batch size of 1. For baseline methods, we used their recommended hyperparameters from official implementations to ensure fair comparison.

## Appendix B. Dataset Table

Table 6: Dataset statistics for Appendiceal Cancer, TCGA-NSCLC, TCGA-ESCA and BRACS. WSIs passing quality control are included. For the TCGA cohorts, slides with missing or ambiguous histologic labels were excluded.

Dataset	Label	Diagnosis	Number of WSIs	
			Site 1	Site 2
Appendiceal Cancer	0	LAMN	74	22
	1	MAC	40	5
TCGA-NSCLC	0	LUAD	496	
	1	LUSC	490	
TCGA-ESCA	0	EAC	63	
	1	ESCC	90	
BRACS	0	N	34	
	1	PB	142	
	2	UDH	70	
	3	ADH	46	
	4	FEA	41	
	5	DCIS	61	
	6	IC	132	

## Appendix C. Supplementary Results

Table 7: Ablation study on layer depth (left) and graph convolution type (right) across different datasets. All results use the same evaluation protocol as the main results.

Metric	Number of Layers			Graph Convolution Type			
	2	3	4	GCN	GIN	SAGE	GAT
<b>Appendiceal Cancer</b>							
<b>BAcc</b>	90.44 <sub>5.08</sub>	<b>92.56</b> <sub>6.36</sub>	88.09 <sub>7.97</sub>	88.05 <sub>8.07</sub>	91.02 <sub>5.03</sub>	89.53 <sub>4.47</sub>	<b>92.56</b> <sub>6.36</sub>
<b>AUC</b>	94.56 <sub>3.73</sub>	<b>96.41</b> <sub>1.94</sub>	95.83 <sub>2.12</sub>	93.39 <sub>3.39</sub>	<b>97.45</b> <sub>1.63</sub>	96.16 <sub>2.31</sub>	96.41 <sub>1.94</sub>
<b>F1</b>	87.51 <sub>5.81</sub>	<b>90.98</b> <sub>7.98</sub>	86.31 <sub>10.53</sub>	83.76 <sub>11.67</sub>	87.79 <sub>5.84</sub>	85.57 <sub>7.00</sub>	<b>90.98</b> <sub>7.98</sub>
<b>TCGA-ESCA</b>							
<b>Accuracy</b>	<b>98.04</b> <sub>2.59</sub>	98.02 <sub>1.62</sub>	98.02 <sub>1.62</sub>	96.73 <sub>2.04</sub>	96.09 <sub>2.40</sub>	98.02 <sub>1.62</sub>	<b>98.02</b> <sub>1.62</sub>
<b>AUC</b>	99.64 <sub>0.44</sub>	<b>99.91</b> <sub>0.17</sub>	99.82 <sub>0.22</sub>	99.72 <sub>0.56</sub>	99.74 <sub>0.51</sub>	<b>100.00</b> <sub>0.0</sub>	99.91 <sub>0.17</sub>
<b>TCGA-NSCLC</b>							
<b>Accuracy</b>	92.39 <sub>1.35</sub>	93.51 <sub>0.75</sub>	<b>94.01</b> <sub>1.63</sub>	91.58 <sub>1.60</sub>	93.51 <sub>2.33</sub>	91.99 <sub>2.27</sub>	<b>93.51</b> <sub>0.75</sub>
<b>AUC</b>	97.31 <sub>1.05</sub>	97.15 <sub>1.47</sub>	<b>97.69</b> <sub>1.15</sub>	97.51 <sub>0.99</sub>	<b>97.97</b> <sub>0.85</sub>	97.10 <sub>1.37</sub>	97.15 <sub>1.47</sub>
<b>BRACS</b>							
<b>BAcc</b>	31.77 <sub>6.18</sub>	<b>33.76</b> <sub>6.07</sub>	31.95 <sub>1.66</sub>	33.26 <sub>3.03</sub>	<b>34.01</b> <sub>2.67</sub>	33.76 <sub>5.33</sub>	33.76 <sub>6.07</sub>
<b>AUC</b>	76.16 <sub>2.36</sub>	77.61 <sub>0.95</sub>	<b>77.67</b> <sub>0.73</sub>	76.30 <sub>0.88</sub>	76.26 <sub>0.92</sub>	75.95 <sub>3.12</sub>	<b>77.61</b> <sub>0.95</sub>
<b>F1-macro</b>	28.06 <sub>6.00</sub>	28.74 <sub>6.08</sub>	<b>29.14</b> <sub>2.45</sub>	26.77 <sub>2.57</sub>	<b>32.10</b> <sub>2.25</sub>	29.93 <sub>4.91</sub>	28.74 <sub>6.08</sub>

Table 8: Computational efficiency comparison of three graph-based WSI methods on the same data splitting. Throughput is measured as the average number of graph batches processed per second (approximately WSIs per second) over a 30-epoch run.

Method	Average Throughput
Patch-GCN	7.03
WiKG	5.40
ResGAT	6.91

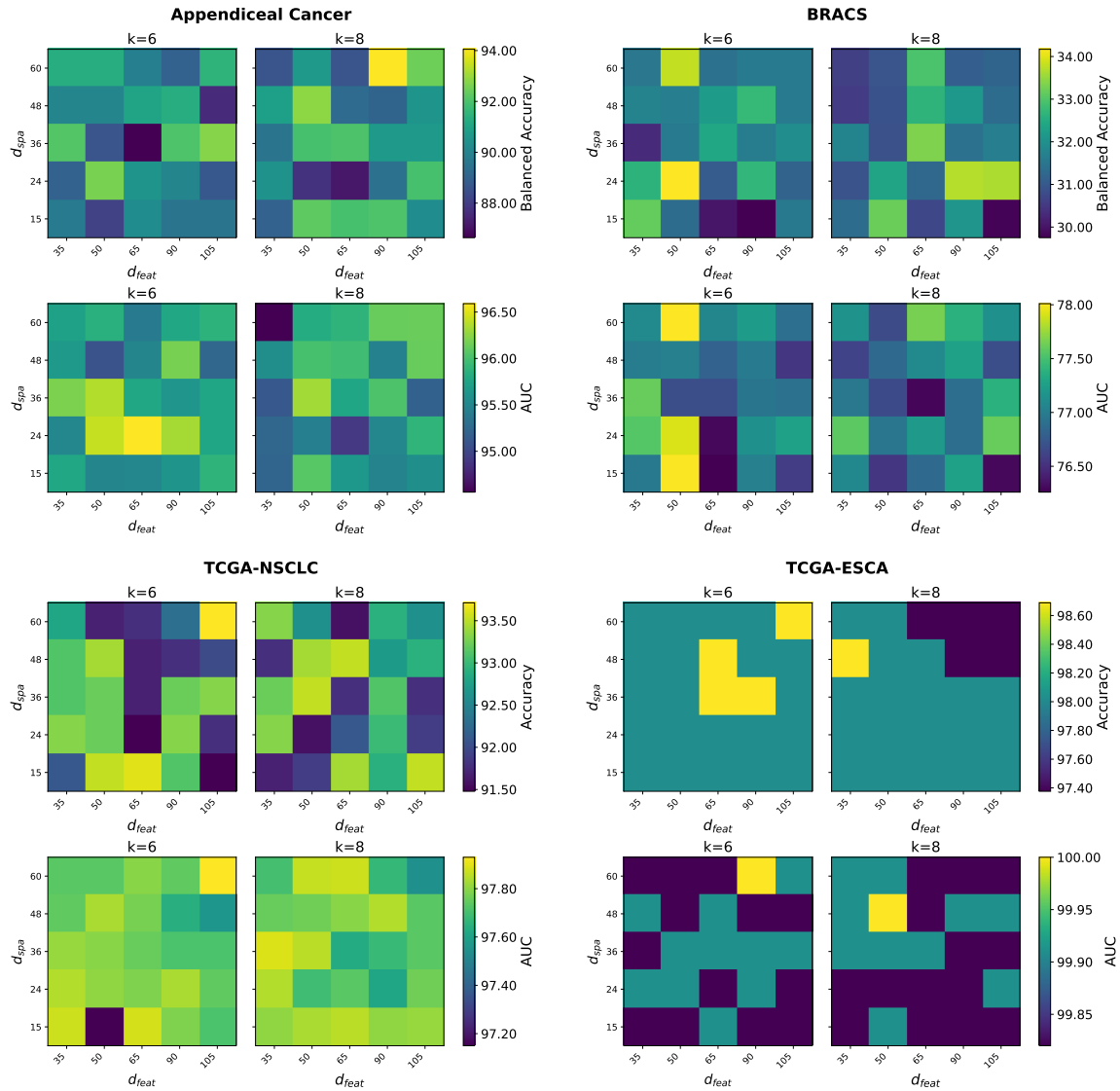


Figure 2: Hyperparameter sensitivity heatmaps across four datasets (Appendiceal Cancer, BRACS, TCGA-NSCLC, TCGA-ESCA). For each dataset, we visualize the performance over the hyperparameters ( $d_{spa}$ ,  $d_{feat}$ ) grid at two graph sparsity settings ( $k = 6, 8$ ). The top row reports the primary metric (Balanced Accuracy for Appendiceal Cancer and BRACS datasets; Accuracy for others), and the bottom row reports AUC score. Within each dataset, the two heatmaps in the same row share the colorbar to enable direct comparison between  $k$  values; brighter colors indicate better performance.

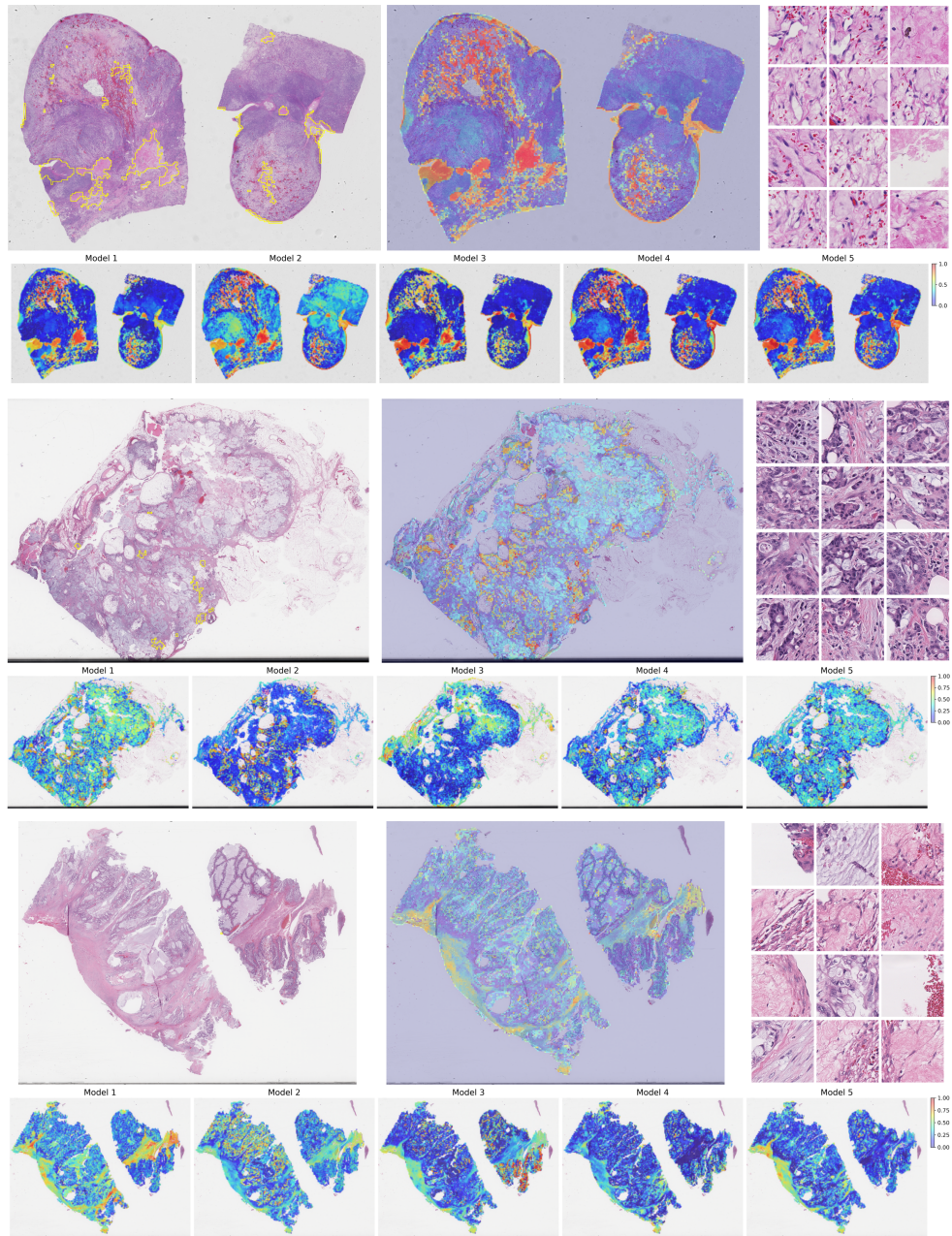


Figure 3: Heatmap visualizations for representative MAC cases (WF53, S22, S36). The first row shows the aggregated heatmap and the corresponding high-contribution regions outlined in yellow, computed as a confidence-weighted average of Models 1, 2, and 5. The second row displays heatmaps from the five cross-validation models. Selected high-contribution patches are shown for localized inspection.

