# Weakly Supervised Joint Whole-Slide Segmentation and Classification in Prostate Cancer

Pushpak Pati[†1*], Guillaume Jaume[†2,3,4,5], Zeineb Ayadi[8], Kevin Thandiackal[1,6],
Behzad Bozorgtabar[7], Maria Gabrani[1], Orcun Goksel[6,9]

[1]*IBM Research Europe, Switzerland*
[2]*Department of Pathology, Brigham and Women's Hospital, Harvard Medical School, USA*
[3]*Department of Pathology, Massachusetts General Hospital, Harvard Medical School, USA*
[4]*Cancer Program, Broad Institute of Harvard and MIT, USA*
[5]*Data Science Program, Dana-Farber/Harvard Cancer Center, USA*
[6]*Computer-Assisted Applications in Medicine, ETH Zurich, Switzerland*
[7]*Signal Processing Laboratory 5, EPFL, Switzerland*
[8]*EPFL, Switzerland*
[9]*Department of Information Technology, Uppsala University, Sweden*

*Abstract*—The segmentation and automatic identification of histological regions of diagnostic interest offer a valuable aid to pathologists. However, segmentation methods are hampered by the difficulty of obtaining pixel-level annotations, which are tedious and expensive to obtain for Whole-Slide images (WSI). To remedy this, weakly supervised methods have been developed to exploit the annotations directly available at the image level. However, to our knowledge, none of these techniques is adapted to deal with WSIs. In this paper, we propose WHOLESIGHT, a weakly-supervised method, to simultaneously segment and classify WSIs of arbitrary shapes and sizes. Formally, WHOLESIGHT first constructs a tissue-graph representation of the WSI, where the nodes and edges depict tissue regions and their interactions, respectively. During training, a graph classification head classifies the WSI and produces node-level pseudo labels via post-hoc feature attribution. These pseudo labels are then used to train a node classification head for WSI segmentation. During testing, both heads simultaneously render class prediction and segmentation for an input WSI. We evaluated WHOLESIGHT on three public prostate cancer WSI datasets. Our method achieved state-of-the-art weakly-supervised segmentation performance on all datasets while resulting in better or comparable classification with respect to state-of-the-art weakly-supervised WSI classification methods. Additionally, we quantify the generalization capability of our method in terms of segmentation and classification performance, uncertainty estimation, and model calibration.

*Keywords*-Computational Pathology, Whole-Slide image segmentation, Weakly supervised learning, Weakly supervised classification, Weakly supervised segmentation

## I. INTRODUCTION

Prostate cancer is the second most frequently diagnosed cancer in men in the United States, with 250,000 new registered cases resulting in 35,000 deaths in 2022 [1]. Yet the number of pathologists, whose role is critical in the diagnosis and management of cancer patients, is gradually declining. In the United States, an 18% decrease was recorded between 2007 and 2017, resulting in a 42% increase in the average workload [2]. In addition, the practice of uro-pathology also has its share of challenges [3]. Although the diagnostic criteria for grading prostate cancer are established [4], the continuum of phenotypic features across the diagnostic spectrum leaves room for disparities, with significant intra- and interobserver variability [5], [6]. The manual inspection of slides is also tedious and time-consuming, and would benefit from automation and standardization. These elements justify the development of Computer-Aided Diagnosis (CAD) tools to automate the diagnostic workflow.

To this end, several Artificial Intelligence (AI) based CAD tools are proposed, including nucleus segmentation and classification [7]–[9], gland segmentation [9]–[11], and tumor detection [12], [13]. Albeit the remarkable performance, these tools often demand task- and tissue-specific annotations on large datasets, which are tedious, time-consuming and often infeasible to acquire. For reducing annotation requirements, different approaches are proposed, in particular, weakly-supervised methods based on Multiple Instance Learning (MIL) framework for the automatic *classification* of Whole-Slide Images (WSIs) [14], [15].

Although classification is useful, it remains limited in its role of supporting the pathologist's attention during diagnosis. In this context, semantic segmentation methods are preferable as they enable the generation of pixel-level delineation of the tissue constituents that can highlight diagnostically relevant regions. Such visualization allows for strengthening trust between pathologists and CAD tools. Additionally, the identified regions can be leveraged by a classifier to improve patient diagnosis. However, semantic

*Corresponding author: Pushpak Pati. Email: pus@zurich.ibm.com

1

segmentation generally requires pixel-level labels, which makes it more demanding in terms of annotations than classification tasks. For this reason, the development of weakly-supervised semantic segmentation (WSSS) methods appears as the most adequate response.

While WSSS has been successful on natural images, it encounters various challenges when applied to histopathology images [16], as they, (1) contain fine-grained objects with large intra-class variations [17]; (2) often include ambiguous boundaries among histology components [18]; (3) can be several giga-pixels with arbitrary tissue sizes. Nevertheless, some WSSS methods are proposed for various histology tasks. The methods by [19]–[25] performing WSSS at patch-level are limited by the need for patch-level annotations, and inability to perform global contextualized WSI segmentation. While [26], [27] scale to larger tiles, they pose high computational complexity and memory requirements for operating on WSIs. The methods by [25], [26] require *exact* tile annotations for model training, *i.e.*, a precise denomination of each lesion type in a tile, which requires pathologists to annotate images beyond standard clinical needs. On a different note, recent WSI classification methods use attention mechanisms or feature attributions to highlight salient regions [14]. Though these regions are informative for visual assessment, they are insufficient, incomplete, and blurry for accurately delineating relevant regions. Additionally, producing granular saliency requires densely overlapping patch predictions, which is computationally expensive while working with WSIs.

In view of the aforementioned limitations of the WSSS methods, we propose WHOLESIGHT, "Whole-slide SegmentatIon using Graphs for HisTopathology", that can simultaneously segment and classify arbitrarily large histopathology images by using WSI-level labels, and without any task-specific assumptions or post-processing. Formally, WHOLESIGHT transforms an image into a superpixel-based tissue-graph (TG), and considers the segmentation problem as a *node-classification* task. WHOLESIGHT incorporates both local and global tissue microenvironment to perform contextualized segmentation, principally in agreement with inter-pixel relation-based WSSS [28]. To summarize, our contributions are:

- WHOLESIGHT, a novel graph-based weakly-supervised method to jointly segment and classify WSIs using readily available WSI-level annotations.
- A comprehensive evaluation of WHOLESIGHT on 3 prostate cancer datasets for Gleason pattern segmentation and Gleason grading, and benchmarked against state-of-the-art WSI-level weakly-supervised methods.
- Thorough generalizability quantification of WHOLESIGHT on *in-* and *out-of-domain* cohorts in terms of segmentation and classification performance, uncertainty estimation, and calibration of neural network predictions.

A preliminary version of this work was presented as [29]. Our substantial extensions herein include, (1) an improved WHOLESIGHT method in terms of model architecture and automatic synthesis of node labels, (2) extensive evaluations on large cohorts of WSIs (approximately $100\times$), and (3) generalization assessment.

## II. RELATED WORK

### A. Weakly-supervised histopathology image classification

Weakly-supervised classification of WSIs has been mostly developed around MIL. In MIL, a WSI is first decomposed into a "bag" of patches and are encoded by a *neural encoder*, *e.g.*, a Convolutional Neural Network (CNN). Then, an *aggregator* pools the patch embeddings to produce a slide-level representation for mapping to a class label via a *neural predictor*. The *aggregator* can be based on an attention mechanism weighing the importance of each patch, as in [14], [30], or as recently proposed, it can take the form of a transformer [31], [32] or a Graph Neural Network (GNN) [33], enabling modeling inter-patch dependencies and global context. Differently, context can be modeled using multi-scale representations of WSIs, either via multi-magnification patch embeddings [23], [34] or by learning to automatically select important regions, as proposed in [35], [36]. Despite the success of these approaches, they cannot directly be extended for semantic segmentation.

### B. Weakly-supervised histopathology image segmentation

WSSS approaches in histopathology can be categorized by the type of supervision (or annotation), *e.g.*, point annotations, scribbles, or image-level labels, and the scale of operation, *e.g.*, patches, tiles, Tissue-Micro Array (TMAs), or WSIs. [37]–[39] utilized point annotations to segment cells and nuclei in histology patches. [23], [40] used scribble annotations to segment tissue and tumor regions, respectively, at patch-level. Both the approaches used U-Net [41], where [23] leveraged concentric patches across multiple magnifications for including relevant context information, and [40] modified the objective function to balance the contribution of the annotated pixels. The majority of the WSSS methods in histopathology utilized image-level supervision and are limited to operate with patch annotations. [19] proposed multiple clustered instance learning to process sliding patches for simultaneous grading and segmentation of colon TMAs. [21] trained a binary classifier for pixel-level predictions and afterward computed an image-level prediction from pixel labels via a softmax function. They optimized image prediction, such that pixel predictions were improved. [22] proposed CAMEL, a MIL-based label enrichment method. It split an image into latticed instances, generated instance labels, and assigned instance labels to corresponding pixels to enable supervised segmentation. [26] proposed HistoSegNet, which trained a CNN to predict tissue types in a tile and used feature attribution to derive

pixel-level predictions. It also employed a series of dedicated post-processing steps for prediction refinement. [24] used foreground proportion as the weak labels and combined a fully convolutional network and a graph convolutional network for tissue segmentation. [25] proposed a feature attribution-based model to generate pseudo labels, followed by a multi-layer pseudo-supervision network for segmenting tissue types. As a main limitation, these methods cannot perform WSSS on WSIs using only WSI labels. To perform WSSS beyond patch-level, [27] proposed WeGleNet, that scales to TMAs. It included a segmentation- and a global-aggregation layer to classify images during training, and up-sampled pixel-level softmax activations during inference for image segmentation. However, the method cannot precisely delineate lesions and highlight multiple lesion occurrences. It also requires processing densely overlapping patches for fine segmentation, and cannot scale to WSIs. In contrast, our WHOLESIGHT can perform WSSS by leveraging image-level supervision, while efficiently scaling to WSIs of arbitrary dimensions.

### C. Generalization quantification in histopathology

Generalizability of CAD tools in histopathology is affected by domain-level biases, which are introduced due to numerous reasons, such as different staining protocols, manufacturing devices, materials, and scanning devices with respective color response [42]. Though generalizable tools, that are robust to domain shifts, are desired, it is challenging to model and detect the domain shifts in Deep Learning (DL). Nevertheless, several approaches have been proposed to reduce such domain shifts via data- and model-level adaptation.

Data-level adaptation can be achieved via stain normalization [43]–[46], color augmentation [47], [48], or stain invariant feature learning [49], [50]. Model-level adaptation is typically done via domain adversarial training [51]–[54], which leverages target domain unlabeled data along with source domain data for modeling.

However, the aforementioned data- and model-level adaptation approaches do not exhaustively assess the generalization ability of their trained DL models beyond task performance. In this case, accurate *uncertainty estimation* and *model calibration* are crucial to know *when* to trust the model – a task known to be challenging for neural networks that often provide over-confident predictions [55]–[57]. To the best of our knowledge, computational pathology research in these directions is scarce and remains unexplored.

### III. METHODOLOGY

In this section, we present WHOLESIGHT for scalable WSSS of histopathology images. First, we transform a WSI into a TG representation, where nodes and edges of the graph denote tissue regions and their interactions, respectively (Section III-B). Next, a GNN contextualizes node

embeddings characterizing tissue regions (Section III-C), which are then processed by a *graph classification head* for Gleason grading (Section III-D). Finally, we generate node-level pseudo labels using feature attribution and a node selection strategy, which are used to train a *node classification head*. The *node-head* outputs the segmentation mask with pixel-level Gleason pattern assignment (Section III-E). An overview of the method is presented in Figure 1.

### A. Notation and preliminaries

We define a graph $G \in \mathcal{G}$ as $(V_G, E_G, H)$, where $V_G$ and $E_G$ denote the set of nodes and edges, respectively, $H \in \mathbb{R}^{|V| \times d}$ denotes $d$-dimensional node features (or denoted at node-level as $H_{v,.} := h(v) \in \mathbb{R}^d$), and $\mathcal{G}$ is the set of graphs. The neighborhood of a node $v \in V_G$ is denoted as $\mathcal{N}(v) := \{u \in V_G \mid (v, u) \in E_G \ \lor \ (u, v) \in E_G\}$. We represent the cardinality of a set as $|.|$, *e.g.*, $|\mathcal{N}(v)|$ indicates the number of neighbors of $v$.

GNNs [58] are a class of neural networks that learn from graph-structured data. Specifically, GNNs follow a two-step procedure to contextualize node features by including neighborhood node information. First, in an AGGREGATE step, for each node $v \in V_G$, the neighboring node features $\mathcal{N}(v)$ are aggregated by a differentiable and permutation-invariant function. Next, in an UPDATE step, the current features of $v$ and the aggregated features of $\mathcal{N}(v)$ are processed by a differentiable operator to update the features of $v$. This procedure is repeated $T$ times, where $T$ is the number of GNN layers.

In this work, we use the Graph Isomorphism Network (GIN) [59], where the AGGREGATE step is a *mean*-operator, and the UPDATE step includes a multi-layer perceptron (MLP). Formally, a GIN layer is given as,

$$h^{(t+1)}(v) = \text{MLP}\Big( h^{(t)}(v) + \frac{1}{|\mathcal{N}(v)|} \sum_{u \in \mathcal{N}(v)} h^{(t)}(u) \Big) \quad (1)$$

$T$ GIN layers, denoted as $\mathcal{F}_\theta$, are stacked to acquire context information up to $T$-hops for each $v$. For graph classification, a fix-sized graph-level embedding $h_G$ is derived by pooling the node embeddings $h^T(v), \ \forall v \in V_G$ by a READOUT step, *e.g.*, a *mean*-operation. Subsequently, $h_G$ is mapped to target classes by a classifier network, $\mathcal{F}_\phi$. Similarly, for node classification, $h^T(v), \ \forall v \in V_G$ can be classified by a classifier network $\mathcal{F}_\psi$.

Formally, classification aims to predict target label $y \in \mathcal{K}$ for an input $x \in \mathcal{X}$, where $\mathcal{K}$ and $\mathcal{X}$ denote the set of classes and inputs, respectively. Given a set of sample pairs $\{(x_i, y_i)\}_{i=1}^N$, where $N$ is the number of samples and $(x_i, y_i) \sim p(x, y)$, the data likelihood can be expressed as $p(Y|X, \theta, \phi) = \Pi_{i=1}^N p(y_i|x_i, \theta, \phi)$. The optimal parameters $(\hat{\theta}, \hat{\phi})$ are obtained by maximum likelihood estimation, or equivalently by minimizing the Negative Log-Likelihood (NLL) $-\sum_{i=1}^N \log p(y_i|x_i, \theta, \phi)$. For graph classification, a sample pair is denoted as $(y_G, G)$, $y_G \in \mathcal{K}_\mathcal{G}$, $G \in$
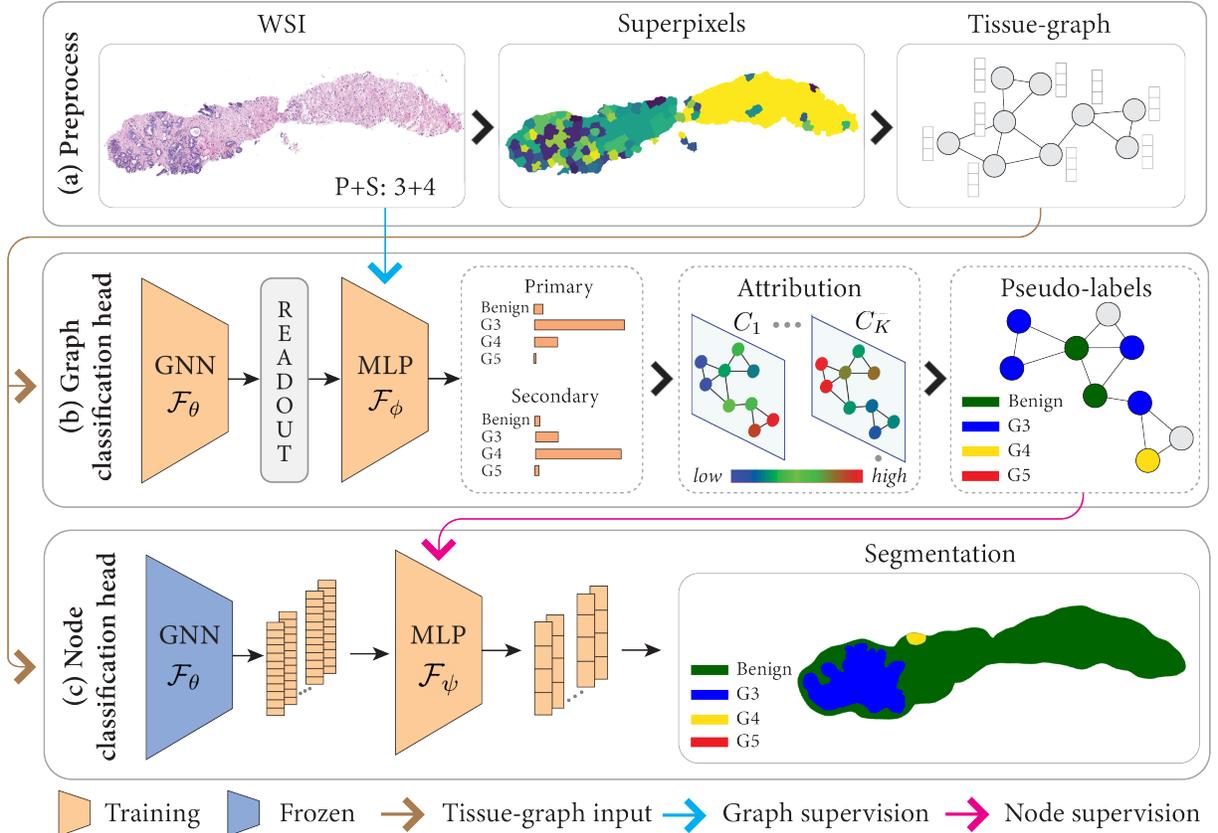
Figure 1: Overview of the proposed WHOLESIGHT method. (a) In the preprocessing step, a TG is constructed to represent a WSI, where the nodes and edges are defined by identifying superpixels and region adjacency connectivity, respectively. (b) The *graph classification head* classifies the TG into primary and secondary Gleason patterns. Subsequently, a feature attribution technique and a node selection strategy derive node-level pseudo-labels. (c) The *node classification head* learns on the pseudo-labels to classify the nodes, thereby resulting in the WSI segmentation.

$\mathcal{G}$. For node classification, a sample pair is denoted as $(y_V, v)$, $y_V \in \mathcal{K}_V$, $v \in \mathcal{V}$. For the task in this paper, the set of graph- and node-level classes are the same, *i.e.*, $\mathcal{K} := \mathcal{K}_\mathcal{G} = \mathcal{K}_\mathcal{V}$.

We further introduce the notion of model *calibration* [60]. Intuitively, the probability of outcomes, *i.e.*, confidence scores, of a calibrated model should match its performance. For example, the samples predicted with an average confidence of 60% by a model should have an average accuracy of 60%. Formally, for a given network, $f : \mathcal{X} \to \mathcal{K}$, and $p(X, Y)$ a joint distribution over the data and the labels, $f(x)$ is said to be calibrated with respect to $p$ if, $\mathbb{E}_p[Y | f(X) = \beta] = \beta$, $\forall \beta \in [0, 1]$. The *calibration* can be visualized with a *reliability diagram* [61]. Namely, all the samples in the dataset are assigned to bins according to their predicted confidence scores. Then, the model accuracy is computed for the samples in each bin. The network performance is plotted against the binned confidence scores, where deviations from the diagonal represent uncalibrated bins.

### B. Preprocessing and tissue-graph construction

First, we stain-normalize the input H&E stained images using the method by [44]. It reduces appearance variability across images caused during tissue preparation, *i.e.*, different specimen preparation techniques, staining protocols, fixation characteristics, and imaging device characteristics [62], [63]. Then, we transform the normalized images into TGs (Figure 1(a)), where the nodes and the edges of a TG denote tissue regions and inter-tissue interactions, respectively. Motivated by [64], [65], we consider superpixels as the visual primitives to encode tissue regions. Compared to rectangular patches, superpixels are more flexible to accommodate arbitrary shapes according to the local homogeneity of tissue. The homogeneity constraint also restricts the superpixels to span across multiple distinct structures and include different morphological regions.

TG construction follows [66], where the prominent steps are, (1) detection of superpixels to define nodes $V_G$, (2) characterization of superpixels to define node features $H$, and (3) building graph topology to define edges $E_G$. We

adopt a two-step process to identify superpixels in a WSI. First, we use Simple Linear Iterative Clustering (SLIC) [67] to produce over-segmented superpixels. Over-segmentation is conducted at a low magnification to capture homogeneous regions while offering a good compromise between granularity and smoothing-out noise. In the second step, the over-segmented superpixels are hierarchically merged according to their channel-wise color similarity at high magnification. Color similarity is quantified in terms of channel-wise 8-bin color histograms, mean, standard deviation, median, energy, and skewness. The resulting merged tissue regions form the nodes of the TG. The merged superpixels denote morphologically meaningful homogeneous regions. Additionally, merging reduces the node complexity of the TG, thus enables the scaling of TG to a large WSI and contextualization to distant tissue regions.

We characterize the TG nodes by morphology and spatial features. Considering the potentially arbitrary dimension of superpixels, we use a two-step process to derive morphology. First, we extract patches of size 144×144 pixels from a superpixel, resize them to 224×224 size, and encode them into 1280-dimensional features via MobileNetV2 network [68] pre-trained on ImageNet [69]. Superpixel-level features are computed as the mean of the patch-level features. Next, we compute spatial features for each node by normalizing the superpixel centroids by the image dimensions. Normalization ensures the invariability of the spatial features to the varying dimensions of input WSIs. Finally, we define the TG edges by constructing a region adjacency graph topology [70] using the spatial connectivity of superpixels. To this end, we assume that adjacent tissue regions biologically interact the most, and thus should be connected in a TG.

### C. Contextualization of node embeddings

Given a TG, we learn discriminative node embeddings (see Figure 1(b)) by using the node context information, *i.e.*, the tissue microenvironment and the inter-tissue interactions. Specifically, we use GIN [59] denoted as $\mathcal{F}_\theta$. Since GNNs can operate on graphs of arbitrary and varying sizes, they allow to encode histopathology images represented in form of TGs without needing tile-based processing. As the discriminative information of a node relies on its local subgraph structures and can lie at different abstraction levels in the GNN, we employ a Jumping Knowledge [71] strategy to utilize multi-level node representations. Namely, the final node-level embedding after $T$ GIN-layers is defined as,

$$h^{(T)}(v) = \text{CONCAT}(h^{(t)}(v),\ \forall t \in \{1,...,T\}) \quad (2)$$

where CONCAT denotes a concatenation operation.

### D. WSI classification

Following the contextualization of node features, a *graph-classification head* classifies the TG by using graph-level embeddings $h_G$ and graph/image-level supervision. To obtain a fix-sized $h_G$, we use a READOUT operation that averages the node embeddings $h^{(T)}(v),\ \forall v \in V_G$. Subsequently, $h_G$ is input to a multi-task classifier for primary and secondary Gleason grading. Specifically, the classifier includes two Multi-Layer Perceptrons (MLPs), denoted as $F_\phi = \{F_{\phi_1}, F_{\phi_2}\}$, to individually predict the primary, *i.e.*, the worst Gleason pattern, and secondary, *i.e.*, the second-worst Gleason pattern, in the WSI. Each MLP solves a multi-class problem with $|\mathcal{K}|$ Gleason patterns, *i.e.*, benign, Grade 3, Grade 4, and Grade 5. The final Gleason grade is derived as the sum of the predicted primary and secondary patterns. $\mathcal{F}_\theta$ and $\mathcal{F}_\phi$ are optimized jointly by minimizing the weighted cross-entropy loss,

$$\mathcal{L}_G = \lambda \mathcal{L}_{CE}(y_{G_P}, \hat{y}_{G_P}) + (1-\lambda)\mathcal{L}_{CE}(y_{G_S}, \hat{y}_{G_S}) \quad (3)$$

where, $P$ and $S$ denote primary and secondary labels of ground truth $y_G$ and prediction $\hat{y}_G$, and $\lambda \in [0,1]$ is a hyper-parameter balancing the two terms. Further, during training we introduce class-weights as $w := \{\log(\frac{\sum_i N_i}{N_i}),\ i = \{1,...,|\mathcal{K}|\}\}$, where $N_i$ is the count of class-wise Gleason patterns in the training WSIs. These weights take care of the class imbalance in Gleason grading by assigning a higher value to classes with lower frequency.

### E. Weakly supervised semantic segmentation

Nodes in a TG identify superpixels, *i.e.*, morphologically homogeneous tissue regions. Since each Gleason pattern is characterized by *distinct* morphological patterns, we assume that each tissue region, depicted by a node, includes a *unique* Gleason pattern. Thus, the WSI segmentation task is translated into classifying the nodes of the TG. In presence of only image supervision, the node classification is achieved in two steps. First, pseudo-node labels are generated by using the image labels, and then, the pseudo labels are used to train a node classifier.

*Pseudo node labels:* Following WSI classification, a post-hoc *feature attribution* technique is used to measure the importance of each node towards TG classification. Specifically, we use GRAPHGRAD-CAM [72], [73], an extension of GRAD-CAM [74] to operate with GNNs. Given a graph $G$, GRAPHGRAD-CAM produces class-wise node attribution maps, $A_k,\ \forall k \in \mathcal{K}$. These maps highlight the importance $\forall v \in V_G$ for classifying $G$ into $|\mathcal{K}|$, as shown in Figure 1. Provided the importance scores for $v$ towards $|\mathcal{K}|$, it can be assumed that the label of $v$ is $k \in \mathcal{K}$, if the highest importance score corresponds to class $k$. At this stage, an *argmax* operation across $A_k,\ \forall k \in \mathcal{K}$ can be considered to classify the nodes. However, such node labeling may be suboptimal, because,

- Some nodes marginally contribute and bear low importance scores $\forall k \in \mathcal{K}$ for classifying a graph. However, an *argmax* across the importance scores for a node

greedily selects the class with the highest score, even though the node label is not ascertained.

- A node highly contributing towards the prediction of a class is not necessarily part of this class. For example, a node can bear high importance if it provides useful complementary information for tie-breaking or ruling out another class possibility. Formally, if the set of nodes $V_k \subset V$ has high importance scores for class $k$, the labels of $V_k$ are not ensured to be $k$. Even, the labels of $v \in V_k$ are not guaranteed to be the same.
- A class attribution map does not necessarily highlight all the nodes belonging to the class. Depending on the task complexity, a classifier may utilize only a subset of the informative nodes from a class to predict the graph label. Formally, if the set of nodes $V_k \subset V$ have high importance scores for class $k$, then $V_k$ may not include all the nodes in $\mathcal{V}_k \subset V$ that have the actual label $k$, i.e., $V_k \subset \mathcal{V}_k$.
- In presence of several feature attribution techniques in literature, with different underlying mechanisms, can produce different attribution maps [73]. Thus, a single attribution technique may not be trusted for score-based node classification.

We, therefore, strategize to use the highlighted nodes by feature attribution as pseudo-labels to train a *node-classifier*. For a graph $G$ with Gleason score $P+S$, $P, S \in \mathcal{K}$, we compute node importance scores $I_P$ and $I_S$, $\forall v \in V_G$ using GRAPHGRAD-CAM. As the scores by GRAPHGRAD-CAM are unbounded, we normalize the scores using min-max. Then, we select the top $n\%$ nodes above a threshold $t$, denoted as $V_P$ and $V_S$, where $n$ and $t$ are hyperparameters tuned during training. It selects the most informative nodes for downstream node classification. For a node $v \in V_P$ and $v \in V_S$, we use $\arg\max(I_P(v), I_S(v))$ to ensure $V_P \cap V_S = \emptyset$. Finally, classes with the highest scores are assigned as pseudo labels $y_{\tilde{V}}$ to the nodes. Pursuing the process for all the TGs in the dataset renders pseudo labels $Y_{\tilde{V}}$.

*Node classification:* $Y_{\tilde{V}}$ is used to train the *node-classification head*, as shown in Figure 1. Specifically for a graph $G$, we get the node embeddings $h^{(T)}(v)$, $\forall v \in V_G$ using $\mathcal{F}_{\hat{\theta}}$, where $\hat{\theta}$ are the parameters of the GNN. $\mathcal{F}_{\hat{\theta}}$ is frozen during node classification such that the same GNN backbone is used for both segmentation and classification, thereby reducing the number of trainable parameters. $h^{(T)}(v)$, $\forall v$ are processed by an MLP $\mathcal{F}_{\psi}$ to predict $Y_{\tilde{V}}$. $\mathcal{F}_{\psi}$ is trained by optimizing a weighted multi-class cross-entropy objective. Similar to the graph classification, class-weights are defined as $w := \{\log(\frac{\sum_i N_i}{N_i}), \ i = \{1, ..., |\mathcal{K}|\}\}$, where $N_i$ is the number of annotated nodes of class $i$. The node-wise predicted class labels are used to obtain the final segmentation prediction. Noticeably, WHOLESIGHT does not include any customized post-processing, unlike [26],

thus being applicable to various tissues and segmentation tasks.

Notably, the *graph-* and the *node-classification heads* address complementary tasks for a graph $G$, i.e., at graph-level and at node-level, respectively. Therefore, following the training of $\mathcal{F}_{\psi}$, we unfreeze $\mathcal{F}_{\theta}$, and jointly fine-tune $\hat{\theta}$ and $\hat{\psi}$ with a small learning rate. The complementarity of the tasks provides an additional informative signal to further improve the segmentation and classification performance of WHOLESIGHT.

## IV. EXPERIMENTS

### A. Datasets

We evaluate our method on three datasets containing whole-slide prostate cancer needle biopsies for Gleason pattern segmentation and Gleason grading. Gleason patterns include grade 3 (G3)- moderately differentiated nuclei and poorly-formed cribriform glands, grade 4 (G4)- poorly differentiated nuclei and irregular masses, and grade 5 (G5)- less differentiated nuclei and lack or only occasional glands. Normal glands and non-epithelial tissues are labeled as benign (B). Gleason grade depicts the worst (*primary*, P) and the second-worst (*secondary*, S) Gleason patterns in a WSI. Dataset details are as follows:

*Sicap dataset:* The dataset [75] contains 18,783 patches of size 512×512 with *complete* pixel-level annotations and slide-level Gleason grades for 155 WSIs from 95 patients. The original slides and masks were reconstructed by stitching the patches. The WSIs were scanned at 40× magnification by Ventana iS-can Coreo scanner and downsampled to 10× magnification. The slides were annotated by expert urogenital pathologists at the Hospital Clínico of Valencia, Spain.

*Radboud dataset:* [76] includes 5,759 needle biopsies from 1,243 patients at the Radboud University Medical Center, Netherlands. The slides were scanned with a 3D Histech Panoramic Flash II 250 scanner at 20× magnification (resolution 0.24$\mu$m/pixel) and were downsampled to 10×. Annotations include WSI Gleason grades and noisy pixel-level Gleason pattern masks, released as part of the Prostate cANcer graDe Assessment (PANDA) challenge [77]. The masks were cleaned for segmentation using standard image manipulation techniques, i.e., contextualized noise removal, hole filling, and edge smoothing. In absence of large public datasets with pixel-level annotated prostate cancer WSIs, we used this dataset for developing and evaluating our method.

*Karolinska dataset:* The dataset [78] comprises of 5,662 core needle biopsies from 1,222 patients at various hospitals in Stockholm, Sweden. The slides were scanned with a Hamamatsu C9600-12 and an Aperio Scan Scope AT2 scanner at 20× magnification with a pixel resolution of 0.45202$\mu$m and 0.5032$\mu$m, respectively. The biopsies were annotated by an expert uro-pathologist for Gleason grading.
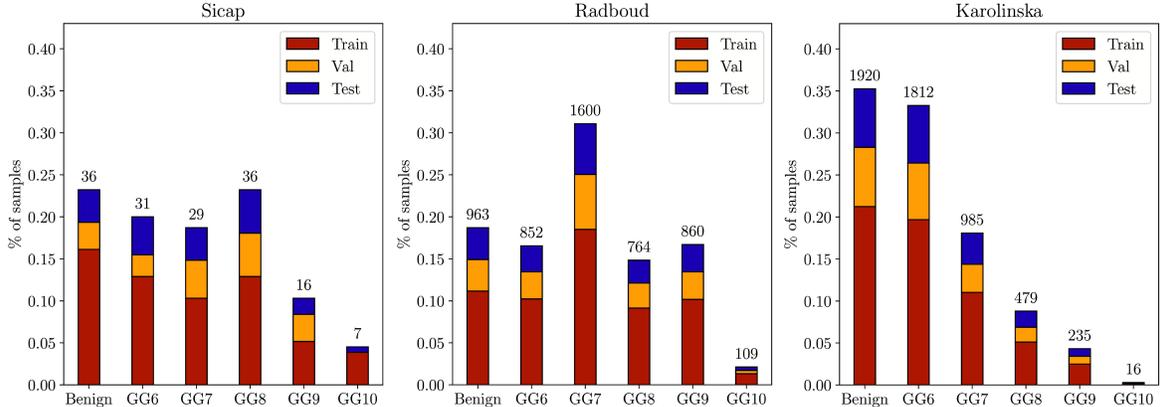
Figure 2: Gleason grade-wise data distribution across train, validation, and test in Karolinska, Radboud and Sicap datasets.

Each dataset is split into train, validation, and test in a ratio of 60%, 20%, and 20% at Gleason grade level, using a random stratification that preserves the percentage of classes in each split. The dataset distributions and splits are displayed in Figure 2, which highlights the class-level imbalances.

### B. Implementation and evaluation

We implemented WHOLESIGHT using PyTorch [79], DGL [80], and Histocartography [81], and conducted experiments on NVIDIA Tesla P100 GPU and POWER9 CPU.

To develop the WHOLESIGHT network, $\mathcal{F}_\theta$, $\mathcal{F}_\phi$, and $\mathcal{F}_\psi$ were designed by optimizing their respective hyperparameters. First, $\mathcal{F}_\theta$ and $\mathcal{F}_\phi$ were trained by using image-level labels, and then pseudo-node labels were created to train $\mathcal{F}_\psi$ to produce segmentation output. The number of GIN layers in $\mathcal{F}_\theta$ were optimized for the values $\{3, 4, 5\}$, where the UPDATE function was defined as a 2-layer MLP with 64 hidden units and ReLU activations. $\mathcal{F}_\phi$ contains two heads for classifying *primary* and *secondary* Gleason grades, where each head consists of a 2-layer MLP with 128 hidden units and ReLU activations. $\mathcal{F}_\psi$ contains a 2-layer MLP with 128 hidden units and ReLU activations.

Considering the small size of the Sicap dataset, node-level augmentations were employed to augment the training dataset. Specifically, random node rotations $\{90, 180, 270\}$ degrees, and horizontal and vertical mirroring were used for augmenting the nodes. Batch size and learning rate were optimized from $\{4, 8, 16\}$ and $\{10^{-4}, 5 \times 10^{-4}, 10^{-3}\}$, respectively. Dropout layers with rates 0.2, 0.5 and 0.5 were included in the MLPs belonging to $\mathcal{F}_\theta$, $\mathcal{F}_\phi$, and $\mathcal{F}_\psi$, respectively. The pseudo-node labels were extracted for selection percentages in $\{5, 10, 15, 20\}$ and thresholds $\{0.5, 0.6, 0.7\}$. Following the hyperparameter tuning, ten WHOLESIGHT models were trained with different network initializations. Validation weighted-F1 was used for model selection. The reported results correspond to the mean and standard deviation over these ten models.

*Classification metrics:* Classification performance was measured by the weighted-F1 score of Gleason grade and the quadratic kappa score ($\kappa^2$) of ISUP grade [82], [83]. ISUP is an alternate grading system whose correspondence with Gleason grading is defined as, Benign → ISUP-0, GG-(3+3) → ISUP-1, GG-(3+4) → ISUP-2, GG-(4+3) → ISUP-3, GG-8 → ISUP-4, and GG≥9 → ISUP-5. $\kappa^2$ captures the degree of disagreement between the prediction and ground truth labels. For example, a grade 6 sample predicted as grade 10 is penalized more than predicting grade 7.

*Segmentation metrics:* Segmentation performance was measured by Dice score. Given the imbalance of the Gleason patterns in the datasets, we also reported the per-pattern Dice score.

*Uncertainty metrics:* Following the work of [84], we evaluated the classification uncertainties in terms of Brier score $s_B$ (lower is better) and the NLL $s_{NLL}$ (lower is better) over a set of $N$ test samples, defined as,

$$s_B = \frac{1}{N} \sum_{n=1}^{N} \sum_{i=1}^{|\mathcal{K}|} (y_i - \hat{y}_i)^2, \quad s_{NLL} = -\frac{1}{N} \sum_{n=1}^{N} \sum_{i=1}^{|\mathcal{K}|} p(y_i) \log \hat{p}(y_i)$$
(4)

*Calibration metrics:* Reliability diagrams provide an intuitive understanding of model calibration. To quantify these observations, we used the Expected Calibration Error (ECE) metric [85], which computes the weighted average deviation of the confidence scores over all the bins, *i.e.*,

$$c_{ECE} = \sum_{b=1}^{B} \frac{N_b}{N} |\text{acc}(b) - \text{conf}(b)|$$
(5)

where $n_b$ is the number of samples in bin $b$, $\text{acc}(b)$ and $\text{conf}(b)$ are the accuracy and the average confidence of samples in $b$.

### C. Baselines

We compared WHOLESIGHT with state-of-the-art WSI classification methods and two variants of WHOLESIGHT.

*FSConv:* We implemented the two-step method proposed by [75] for WSI classification. First, we extracted patches of size 256×256 from WSIs and classified them using FSConv+global-max pooling. The patches were labeled using the Gleason pattern masks, and patches with >90% homogeneous pattern were selected for classifier training. During inference, dense patch predictions produced the output segmentation masks. An MLP was trained on the Gleason grade percentages over the WSI patches for Gleason grading.

WHOLESIGHT*(Graph,* GRAPHGRAD-CAM*):* In comparison to WHOLESIGHT, this baseline contained only $\mathcal{F}_\theta$ and $\mathcal{F}_\phi$. It did not create or utilize pseudo labels, and the segmentation output was obtained by taking the *argmax* over the class-wise GRAPHGRAD-CAM attribution maps.

WHOLESIGHT*(Multiplex, NC):* This variant used both image- and pixel-level supervision during training and acts as the upper bound for WHOLESIGHT. As pixel-level annotations were available, $\mathcal{F}_\psi$ was trained using ground-truth node-level labels, instead of generating pseudo-node labels. The model consisted of the same $\mathcal{F}_\theta$, $\mathcal{F}_\phi$, and $\mathcal{F}_\psi$ as the WHOLESIGHT architecture. In this baseline, $\mathcal{F}_\theta$, $\mathcal{F}_\phi$, and $\mathcal{F}_\psi$ were trained jointly by optimizing a multi-task objective, *i.e.*, WSI-level primary and secondary Gleason score prediction along with node-level Gleason pattern prediction. This variant of WHOLESIGHT was proposed in our preliminary work, as described in [29].

*Multiple Instance Learning (MIL):* MIL methods are state-of-the-art for WSI classification. In particular, we compared to ABMIL [30], which used an attention mechanism to aggregate patch embeddings into a fix-sized WSI embedding that was fed to a classifier for Gleason grading. We also included CLAM [14], a method built on ABMIL by including an additional constrain to cluster similar patch embeddings. Our experiments followed the public implementations * with adjustments to enable multi-task classification.

For all the baselines, hyper-parameters are thoroughly tuned to use the best learning rate and batch size, if applicable. Subsequently, ten models were re-trained from scratch with the optimal parameters. We report the mean and standard deviation over these runs for each experiment.

### D. *WSSS performance analysis*

We studies the classification and segmentation performance of WHOLESIGHT and the competing methods by independently training and testing them on Sicap, Radboud, and Karolinska.

*Performance analysis:* Table I presents the results on Sicap. The analyses are grouped into two supervision settings, *i.e.*, *complete* ($\mathcal{C}$) and *weak* ($\mathcal{W}$). Setting-$\mathcal{C}$ utilizes both image- and pixel-level annotations, whereas, Setting-$\mathcal{W}$ only uses image-level labels. WHOLESIGHT reached

37.6% average Dice score, which significantly outperforms WHOLESIGHT (Graph, GRAPHGRAD-CAM) by +6.6% in absolute. WHOLESIGHT (Multiplex, NC), that acts as the upper bound, produced a significant gain in segmentation compared to WHOLESIGHT. The per-class Dice scores indicate that the benign patterns that constitute most tissue areas have a high detection rate compared to less occurring Gleason patterns. For the classification task, WHOLESIGHT outperformed ABMIL and CLAM, both in terms of Gleason grade weighted-F1 and ISUP $\kappa^2$. Notably,

Table II presents the results on Radboud. WHOLESIGHT rendered an absolute gain of +4.8% in average Dice score over WHOLESIGHT (Graph, GRAPHGRAD-CAM). This confirms the utility of pseudo-node labels for superior segmentation. Similar to the observations on Sicap, benign patterns had a high detection rate, followed by G3, G4, and G5 patterns. As Radboud dataset includes more G5 patterns than Sicap, we observed a significant gain in detecting high-grade Gleason patterns. For the classification task, the observations were also consistent with the observations on Sicap.

Table III presents the results on Karolinska. In absence of ground truth pixel-level annotations, the segmentation performances could not be computed. WHOLESIGHT (Graph) outperformed the baselines in terms of classification performance.

The observations across Table I, III and III conclude that, jointly optimizing classification and segmentation objectives provide complementary information to improve the overall classification performance, *i.e.*, WHOLESIGHT (Multiplex) > WHOLESIGHT > WHOLESIGHT (Graph).

### E. *Generalization: performance, uncertainty, and calibration*

We studied the generalization ability of WHOLESIGHT following a modified training setting. Specifically, we used Radboud and Karolinska training WSIs for model training. Thus, the training set encompassed better sample variability and diagnostically more challenging cases than the standalone training counterparts on individual datasets. Testing was performed individually on Radboud and Karolinska test WSIs, herein studying the *in-domain* performance. Further, we tested on the entire Sicap dataset, which consisted of *out-of-domain* WSIs.

*Performance analysis:* Table IV compared the classification performance of WHOLESIGHT and the competing baselines. In terms of the weighted-F1 score on the *in-domain* test set, WHOLESIGHT outperformed ABMIL and performed better or comparable to CLAM. Similar patterns were also observed for ISUP $\kappa^2$. However, the variances of the classification for WHOLESIGHT is consistently lower than ABMIL and CLAM. When tested on the *out-of-domain* Sicap dataset, WHOLESIGHT achieved significantly better classification than ABMIL and CLAM.

---

*CLAM publicly available code: https://github.com/mahmoodlab/CLAM

Table I: Classification and segmentation results on **Sicap** dataset. The best performances for using image-level supervision are highlighted in **bold**.

| Annot. | Method | per-class Dice | | | | avg. Dice | GG wF1 | ISUP $\kappa^2$ |
|---|---|---|---|---|---|---|---|---|
| | | Benign | Grade3 | Grade4 | Grade5 | | | |
| $\mathcal{C}$ | FSConv [75] | 65.7±0.5 | 24.4±1.4 | 29.0±1.4 | 8.4±0.6 | 31.9±0.5 | 48.7±3.4 | 50.9±3.4 |
| | WHOLESIGHT (Multiplex, NC) | 92.5±0.3 | 35.4±2.3 | 51.6±2.0 | 11.1±2.2 | 47.6±2.1 | 59.6±4.1 | 84.6±3.2 |
| $\mathcal{W}$ | ABMIL [30] | - | - | - | - | - | 50.2±6.3 | 67.8±5.2 |
| | CLAM [14] | - | - | - | - | - | 51.4±5.5 | 75.2±4.7 |
| | WHOLESIGHT (Graph, GRAPHGRAD-CAM) | 64.4±6.1 | 23.1±2.0 | 32.8±6.9 | 3.7±1.0 | 31.0±2.7 | 53.3±5.3 | 81.9±6.7 |
| | WHOLESIGHT (Graph + Pseudo nodes, Node class.) | **67.6±0.5** | **28.6±0.3** | **49.1±0.4** | **5.0±0.4** | **37.6±0.3** | **58.6±6.2** | **89.0±1.2** |

Table II: Classification and segmentation results on **Radboud** dataset. The best performances for using image-level supervision are highlighted in **bold**.

| Annot. | Method | per-class Dice | | | | avg. Dice | GG wF1 | ISUP $\kappa^2$ |
|---|---|---|---|---|---|---|---|---|
| | | Benign | Grade3 | Grade4 | Grade5 | | | |
| $\mathcal{C}$ | FSConv [75] | 84.3±0.1 | 53.3±0.5 | 62.5±0.3 | 36.9±0.5 | 59.2±0.1 | 45.9±1.3 | 69.4±0.6 |
| | WHOLESIGHT (Multiplex, NC) | 91.5±0.1 | 63.9±0.4 | 66.2±0.3 | 36.8±1.2 | 64.6±0.4 | 68.9±0.9 | 83.8±0.9 |
| $\mathcal{W}$ | ABMIL [30] | - | - | - | - | - | 59.3±1.7 | 79.7±1.2 |
| | CLAM [14] | - | - | - | - | - | 60.3±1.6 | 80.2±1.5 |
| | WHOLESIGHT (Graph, GRAPHGRAD-CAM) | 71.3±2.2 | 26.9±1.0 | 24.9±1.2 | 15.1±0.5 | 34.6±0.6 | 66.0±1.0 | 82.2±0.5 |
| | WHOLESIGHT (Graph + Pseudo nodes, Node class.) | **75.9±0.3** | **32.9±1.0** | **29.1±1.5** | **19.6±0.7** | **39.4±0.3** | **67.9±0.3** | **83.0±0.2** |

Table III: Classification results on **Karolinska** dataset. The best performances for using image-level supervision are highlighted in **bold**.

| Annot. | Method | GG wF1 | ISUP $\kappa^2$ |
|---|---|---|---|
| $\mathcal{W}$ | ABMIL [30] | 65.0±2.0 | 79.1±1.2 |
| | CLAM [14] | 63.6±2.5 | 77.6±2.0 |
| | WHOLESIGHT (Graph) | **70.5±0.6** | **80.2±0.7** |

Table IV: Classification and segmentation results on Radboud, Karolinska, and Sicap datasets for models trained using both Radboud and Karolinska datasets.

| Annot. | Method | Radboud | | | Karolinska | | Sicap | | |
|---|---|---|---|---|---|---|---|---|---|
| | | avg. Dice | GG wF1 | ISUP $\kappa^2$ | GG wF1 | ISUP $\kappa^2$ | avg. Dice | GG wF1 | ISUP $\kappa^2$ |
| $\mathcal{C}$ | FSConv [75] | 59.2±0.1 | 45.9±1.3 | 69.4±0.6 | 34.5±1.1 | 40.1±1.3 | 49.5±0.4 | 52.1±2.1 | 53.8±1.7 |
| | WHOLESIGHT (Multiplex, NC) | 64.5±0.3 | 69.0±1.0 | 83.6±0.9 | 71.2±0.7 | 82.5±1.3 | 60.0±0.5 | 65.5±2.5 | 85.6±2.8 |
| $\mathcal{W}$ | ABMIL [30] | - | 57.6±2.3 | 73.8±2.3 | 65.5±1.3 | 77.3±2.8 | - | 56.4±2.7 | 75.0±7.5 |
| | CLAM [14] | - | 61.7±2.1 | 78.6±1.3 | 69.3±1.3 | **82.8±1.0** | - | 53.1±3.8 | 74.6±4.2 |
| | WHOLESIGHT (Graph, GRAD-CAM) | 34.6±0.6 | 66.0±1.0 | 82.2±0.5 | 69.2±0.9 | 80.3±0.9 | 30.4±1.0 | 65.1±2.3 | 86.1±2.5 |
| | WHOLESIGHT (Graph + Pseudo nodes, Node class.) | **44.6±0.2** | **66.2±0.1** | **82.9±0.1** | **70.2±0.1** | 81.3±0.1 | **42.0±0.3** | **65.2±0.1** | **86.6±0.1** |

Confusion matrices for the best Gleason grading, ISUP grading, primary- and secondary classification with WHOLESIGHT are presented in Figure 3 on the three datasets. It can be observed that most misclassifications lie close to the diagonal. Majority of the confusion occurred between GG6 and GG7, *i.e.*, GG(3 + 3) versus GG(3 + 4) and GG(4 + 3). Such ambiguity is prevalent among pathologists, as shown in [86], [87]. High-grade Gleason grading better on Radboud than Karolinska due to more number of high-grade samples in Radboud. Primary- and secondary classification weighted-F1 for Radboud, Karolinska and Sicap were 79.3%, 81.7%, 81.6% and 62.5%, 69.7%, 64.5%, respectively. This indicated that identifying secondary Gleason pattern is more challenging. Table IV
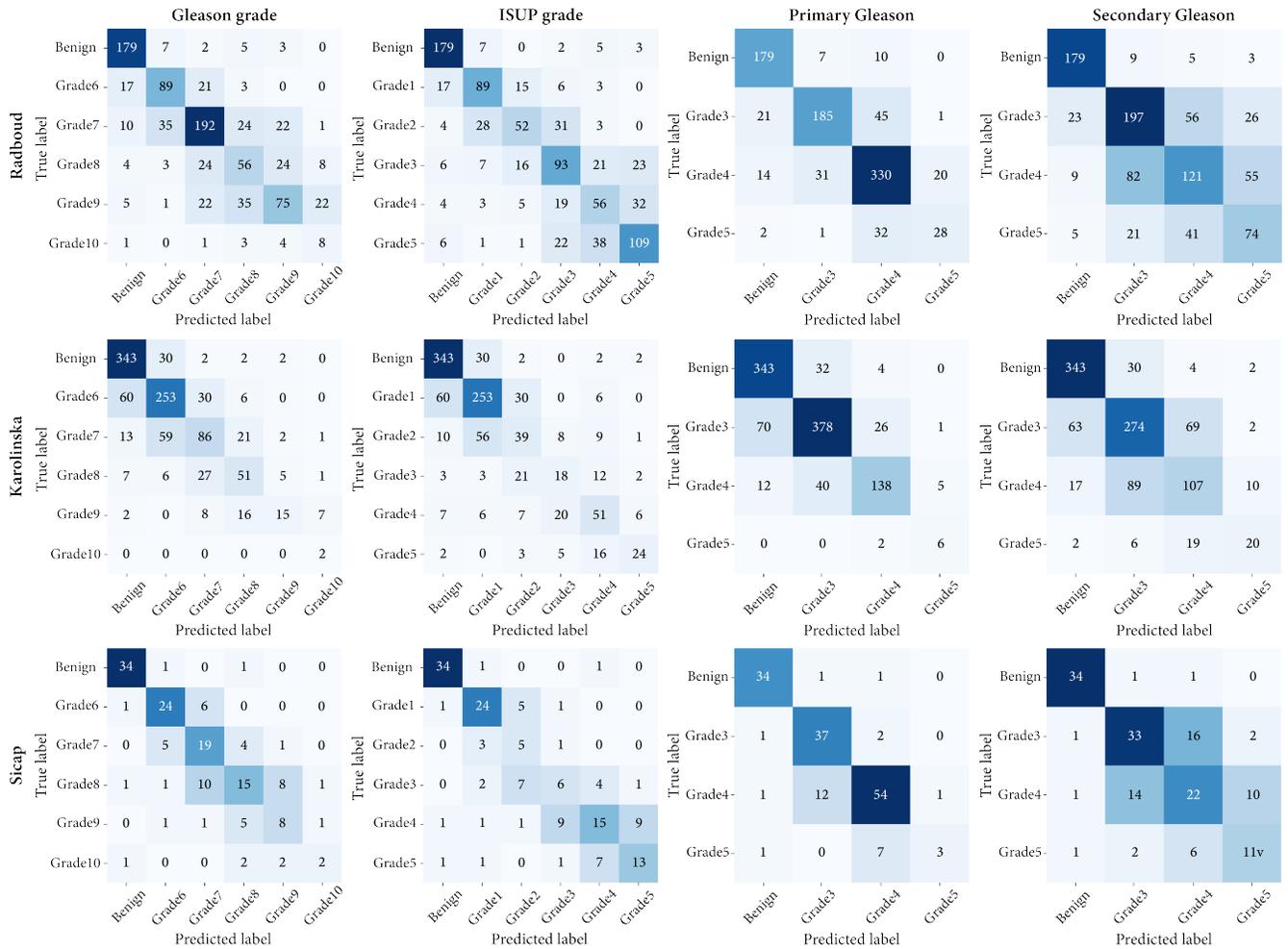
Figure 3: Confusion matrices for Gleason grading, ISUP grading, primary- and secondary Gleason classification on Radboud, Karolinska, and Sicap datasets for with the best WHOLESIGHT model trained using Radboud and Karolinska training datasets.

also presents the generalizability assessment of segmentation. WHOLESIGHT consistently performed better than WHOLESIGHT (Graph, GRAPHGRAD-CAM). Noticeably, the Dice scores on Radboud and Sicap datasets improved over the segmentation results in Table II and I by 5.2% and 4.4% for WHOLESIGHT. It can be reasoned to the usage of more training WSIs, which indicate that WSSS can be improved by utilizing more weak supervision.

*Uncertainty analysis:* Figure 4 presents the classification uncertainty analysis of WHOLESIGHT, WHOLESIGHT (Graph, GRAD-CAM), and WHOLESIGHT (Multiplex, NC), in terms of NLL and Brier score, on Radboud and Karolinska datasets. WHOLESIGHT (Multiplex, NC) rendered a significantly lower NLL than WHOLESIGHT across all datasets for primary, secondary, and Gleason grade (P+S) classification. Noticeably, the NLL and Brier scores were consistently higher for predicting the secondary Gleason

patterns than the primary patterns. This resonates with the fact that identifying secondary patterns is more challenging with higher ambiguity.

*Model calibration analysis:* A model with good uncertainty estimate should be well-calibrated, *i.e.*, the model confidence should be close to the model performance. Figure 4 presents the reliability diagrams of the primary classification head on Karolinska and Radboud datasets. WHOLESIGHT showed consistently better calibration than WHOLESIGHT (Graph, GRAD-CAM) and similar calibration with respect to WHOLESIGHT (Multiplex, NC). ECE also metric quantitatively supported this observation. However, we observed that still all models remains over-confident as the model accuracies over the confidence bins remained lower than the expected calibration (in blue).
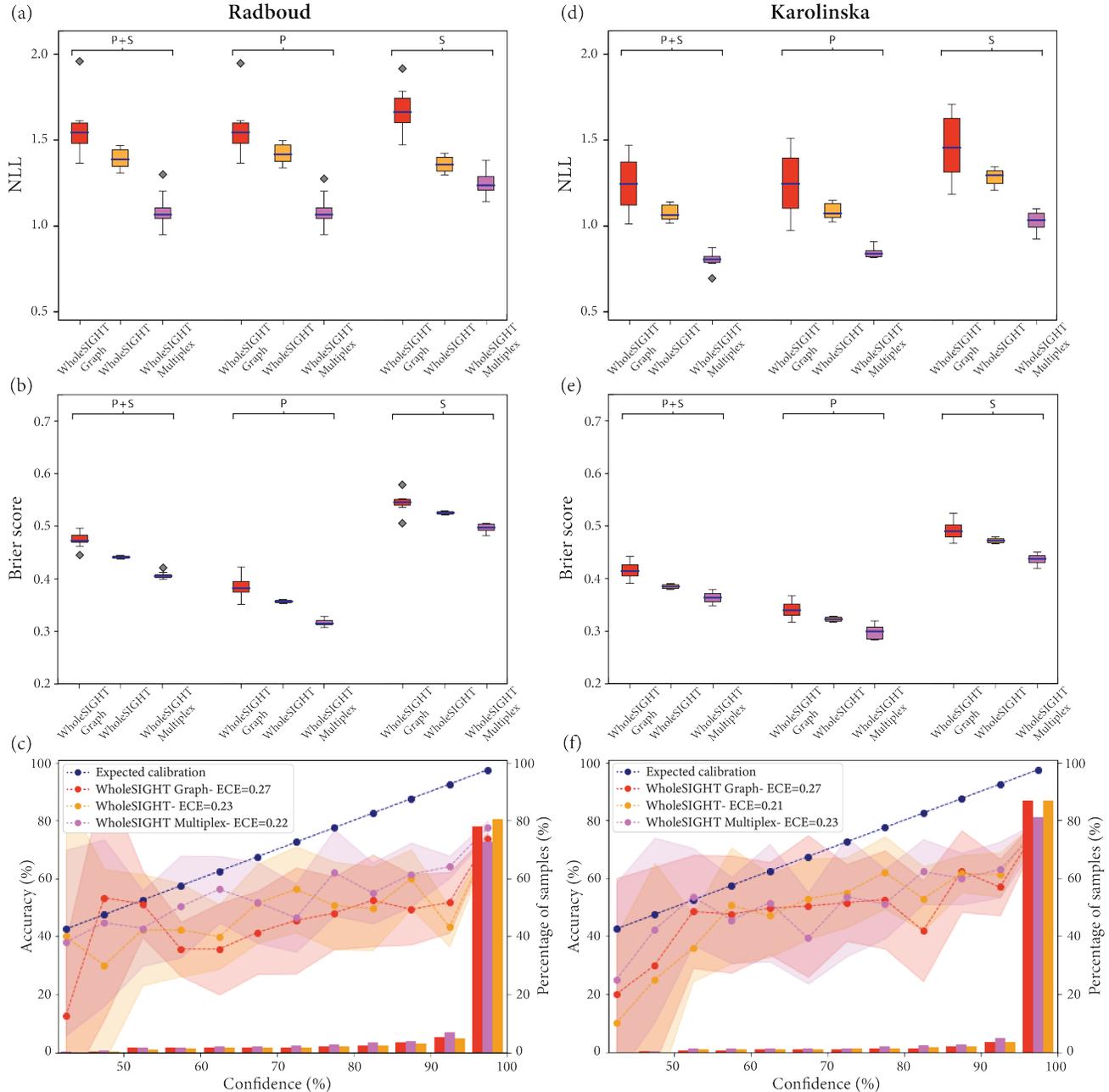
Figure 4: Uncertainty and model calibration analysis of WHOLeSIGHT, WHOLeSIGHT (Graph, GRAD-CAM), and WHOLeSIGHT (Multiplex, NC) models for Radboud (a, b, c) and Karolinska (d, e, f) datasets. (a, d) and (b, e) present NLL (lower is better) and Brier scores (lower is better), respectively. (c, f) present reliability diagrams of the primary Gleason classification head. Expected calibration (blue) highlights a perfectly calibrated model. Calibrations of WHOLeSIGHT, WHOLeSIGHT (Graph, GRAD-CAM), and WHOLeSIGHT (Multiplex, NC) are in orange, red, and purple, respectively, along with the number of samples (in %) in each bin.

*F. Qualitative analysis*

We qualitatively analyze the results of WHOLeSIGHT by (1) visualizing overlaid segmentation masks on WSIs, (2) analyzing the t-distributed stochastic neighbor (t-SNE) [88]

node embeddings, and (3) correlating the segmentation outputs with pathological reasonings.

*Segmentation mask visualization:* Figure 5 demonstrates segmentation predictions obtained with WHOLeSIGHT and WHOLeSIGHT(Multiplex, NC) on Sicap dataset. We can
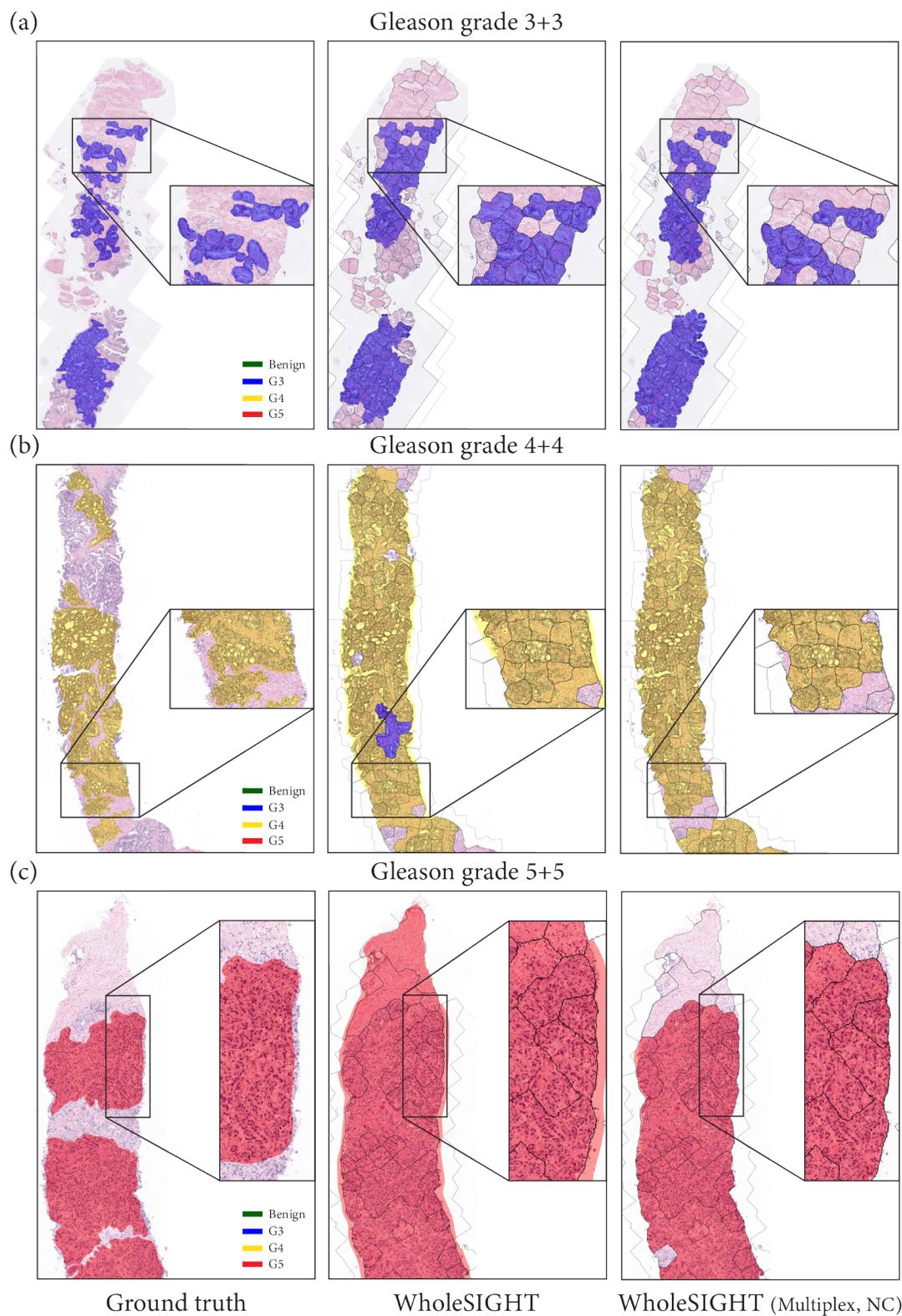
Figure 5: Sample segmentation maps from Sicap dataset. Ground truth is shown on the left, WHOLESIGHT predictions in the middle, and WHOLESIGHT(Multiplex, NC) on the right. Tissue regions, *i.e.*, TG nodes, are represented by black overlay. (a, b, c) display GG(3+3), GG(4+4), and GG(5+5) samples, respectively. For better visualization, benign areas are not highlighted in the segmentation maps.
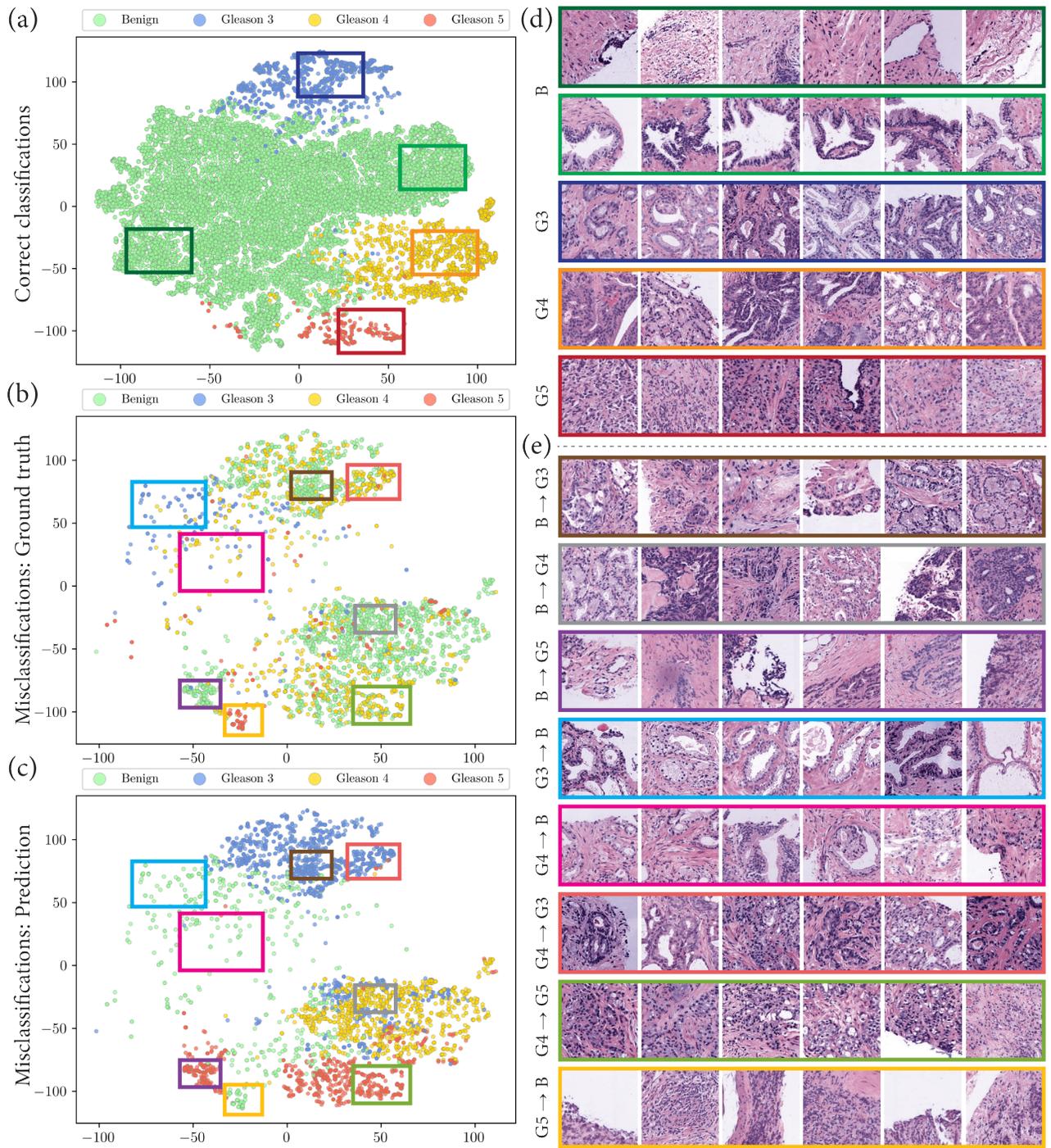
Figure 6: t-SNE visualization of tissue-graph node embeddings and example patches from several regions on the two-dimensional t-SNE feature space for Sicap dataset. (a) t-SNE visualization of the correctly classified nodes. (b) and (c) display the t-SNE visualization of misclassified nodes, where (b) and (c) highlight the ground truth and predicted node labels, respectively. (d) and (e) demonstrate square patches of size $224 \times 224$ at $10\times$ magnification cropped around the node centroids selected from different regions on the t-SNE embedding space. (d) and (e) highlight the correctly and incorrectly classified node patches, respectively. The labels of the patches in (e) are formatted as $Y \rightarrow \hat{Y}$, where $Y$ and $\hat{Y}$ denote the ground truth and the predicted class label. The colored rectangles around the patches in (d) and (e) correspond to respective colored rectangles in (a), (b), and (c).

13

observe that WHOLESIGHT correctly delineates the cancerous regions in the WSIs. Zooming into different regions conclude that the tissue regions of TG, *i.e.*, the nodes of TG, (outlined in black in Figure 5) encode meaningful units of *homogeneous* tissue. It substantiates the relevance of using TG representations for segmenting tissue regions into Gleason patterns. We further notice that WHOLESIGHT, in a few cases, predicts benign regions adjacent to cancerous patterns as cancerous. For example, the benign region, primarily consisting of stroma, in Figure 5(c) is predicted as G5. We argue that these false positive detections do not inhibit the applicability of the method, as neighboring cancerous regions are correctly detected. In a few other cases, WHOLESIGHT correctly detects missed cancerous regions in the ground truth annotations. For instance, in Figure 5(b), the missing G4 region in the upper part of the WSI is correctly identified.

Comparing WHOLESIGHT with WHOLESIGHT (Multiplex, NC), we observe that several false positives are removed, *e.g.*, in Figure 5(a), thus offering more accurate segmentation. However, the improvements by WHOLESIGHT (Multiplex, NC) are achieved at the cost of training with pixel-level annotations that are hardly available in real-world practice. Thus, WHOLESIGHT appears to be an appealing compromise between segmentation performance and annotation requirement.

*Visualizing t-SNE feature space:* A t-SNE visualization of the learned tissue-level embeddings is demonstrated in Figure 6 for Sicap. t-SNE projects the GNN node embeddings onto a two-dimensional feature space, allowing to analyze the connection between node embeddings and the Gleason pattern distribution.

Figure 6(a) displays the t-SNE feature space for the *correctly* classified nodes, which highlights demarcated clusters for each Gleason pattern. The large cluster of benign nodes indicates the variability of the benign tissue. Several patches from each Gleason pattern cluster are presented in Figure 6(d). We can observe the reduced nuclei differentiation across the patches from benign to Gleason grade 5. Further, Figure 6(b) and (c) display the t-SNE feature space for the misclassified nodes. Specifically, Figure 6(b) presents the ground truth node labels, and Figure 6(c) the predicted node labels. Different embedding locations are further selected and highlighted by different colored rectangles and put in relation with corresponding patches to indicate the inter-class ambiguities, as demonstrated in Figure 6(e). For example, the first row in Figure 6(e) showcases patches that are benign but are predicted as G3. We can visually compare these patches with the G3 patches in the third row of Figure 6(d). Similar ambiguities between other pairs of Gleason patterns are also included in Figure 6(e).

*Interpreting model outcomes via predicted segmentations:* Predicted segmentations provide human-understandable *interpretability* maps. For researchers, the segmentations allow to, (1) identify morphological patterns responsible for WSI classification, (2) analyze failure cases by inspecting pixel-level predictions, and ultimately (3) better understand the model behavior towards biomarker discovery. For pathologists, they assist to, (1) put in relation the predicted WSI-level Gleason scores and the highlighted pixel-level Gleason patterns, (2) confirm that the morphology of the identified cancerous regions aligns with pre-established diagnostic criteria.

Additionally, in the perspective of developing AI-assisted human-in-the-loop tools, a Gleason grading system that can simultaneously *classify* and *segment* WSIs is closer to the latest pathological standards. Indeed, recent revisions of the Gleason grading system [83] emphasized the importance of reporting the percentage of each grade for better patient stratification and treatment selection [89]–[92]. These percentages can be trivially derived from the predicted segmentation maps by counting the number of pixels belonging to each pattern. Naturally, such information is not available in mere WSI classification systems. Reporting per-grade percentage is particularly important in ambiguous and borderline cases. For instance, consider two patients with Gleason score 3+4. When a small percentage of pattern-4 is present, *e.g.*, 10%, the case can be considered as an intermediate risk cancer where active patient surveillance is enough [93]. However, a larger secondary pattern may require specific treatments. Reporting percentages of each grade allows us to discriminate between these two scenarios easily.

Similarly, consider a Gleason score 4+3 with a small secondary Gleason pattern, *e.g.*, 90% and 10% area for primary and secondary patterns, respectively. This case will be scored as 4+3, even though it is close to a score of 4+4, which would lead to a different treatment protocol. By explicitly reporting the Gleason pattern percentages, such corner cases can be avoided.

## V. CONCLUSION

Accurate delineation of patterns in whole-slide histopathology images typically demands pixel-level annotations, which are hard to acquire in a real-world scenario. Nonetheless, the semantic segmentation of diagnostically relevant patterns is crucial for disease diagnosis and treatment selection. To this end, we proposed a novel weakly-supervised semantic segmentation method, WHOLESIGHT, that can segment the relevant patterns of interest in histopathology images by leveraging only image-level supervision. To our knowledge, WHOLESIGHT is the first weakly-supervised semantic segmentation method that can operate in an end-to-end manner on histopathology images of arbitrary shape and size. We evaluated our proposed method on three publicly available prostate needle biopsy datasets for Gleason grade classification

and Gleason pattern segmentation. On comparing state-of-the-art methods for histopathology applications, we demonstrated the classification and segmentation superiority of WHOLESIGHT. Additionally, WHOLESIGHT is a modular approach that can utilize both image-level and pixel-level supervision to simultaneously perform image classification and segmentation tasks. Though we have evaluated our method for H&E stained prostate cancer needle biopsies, the technology is easily extendable to other tissue types, imaging techniques, and image modalities.

## REFERENCES

[1] R. Siegel, K. Miller, H. Fuchs, and A. Jemal, "Cancer statistics, 2022," *CA: A Cancer Journal for Clinicians*, vol. 72, pp. 7–33, 2022.

[2] M. Wilson, K. Fleming, M. Kuti, L. Looi, N. Lago, and K. Ru, "Access to pathology and laboratory medicine services: A crucial gap," *Lancet*, vol. 391, no. 10133, pp. 1927–1938, 2018.

[3] M. Amin, S. Smith, V. Reuter, J. Epstein, D. Grignon, D. Hansel, O. Lin, J. McKenney, R. Montironi, G. Paner, H. Al-Ahmadie, F. Algaba, S. Ali, I. Alvarado-Cabrero, L. Bubendorf, L. Cheng, J. Cheville, G. Kristiansen, R. Cote, B. Delahunt, J. Eble, E. Genega, C. Gulmann, A. Hartmann, C. Langner, A. Lopez-Beltran, C. Magi-Galluzzi, J. Merce, G. Netto, E. Oliva, P. Rao, J. Ro, J. Srigley, S. Tickoo, T. Tsuzuki, S. Umar, T. van der Kwast, R. Young, and M. Soloway, "Update for the practicing pathologist: The international consultation on urologic disease-european association of urology consultation on bladder cancer," *Modern Pathology*, vol. 28, no. 5, pp. 612–630, 2015.

[4] P. Tan, I. Ellis, K. Allison, E. Brogi, S. Fox, S. Lakhani, A. Lazar, E. Morris, A. Sahin, R. Salgado, A. Sapino, H. Sasano, S. Schnitt, C. Sotiriou, P. Diest, V. White, M. Lokuhetty, and I. Cree, "The 2019 world health organization classification of tumours of the breast," *Histopathology*, vol. 77, no. 2.

[5] D. Gomes, S. Porto, D. Balabram, and H. Gobbi, "Inter-observer variability between general pathologists and a specialist in breast pathology in the diagnosis of lobular neoplasia, columnar cell lesions, atypical ductal hyperplasia and ductal carcinoma in situ of the breast," *Diagnostic Pathology*, vol. 9, no. 1, pp. 1–9, 2014.

[6] J. Elmore, G. Longton, P. Carney, B. Geller, T. Onega, A. Tosteson, H. Nelson, M. Pepe, K. Allison, S. Schnitt, F. O'Malley, and D. Weaver, "Diagnostic concordance among pathologists interpreting breast biopsy specimens," *JAMA*, vol. 313, no. 11, pp. 1122–1132, 2015.

[7] S. Graham, Q. Vu, S. Raza, A. Azam, Y. Tsang, J. Kwak, and N. Rajpoot, "Hover-net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images," *Medical Image Analysis*, vol. 58, p. 101563, 2019.

[8] P. Pati, A. Foncubierta-Rodríguez, O. Goksel, and M. Gabrani, "Reducing annotation effort in digital pathology: A co-representation learning framework for classification tasks," *Medical Image Analysis*, vol. 67, p. 101859, 2021.

[9] S. Graham, Q. D. Vu, M. Jahanifar, S. E. A. Raza, F. Minhas, D. Snead, and N. Rajpoot, "One model is all you need: multi-task learning enables simultaneous histology image segmentation and classification," *Medical Image Analysis*, vol. 83, p. 102685, 2023.

[10] K. Sirinukunwattana, J. Pluim, H. Chen, X. Qi, P.-A. Heng, Y. Guo, L. Y. Wang, B. Matuszewski, E. Bruni, U. Sanchez, A. Böhm, O. Ronneberger, B. Ben Cheikh, D. Racoceanu, P. Kainz, M. Pfeiffer, M. Urschler, D. Snead, and N. Rajpoot, "Gland segmentation in colon histology images: The glas challenge contest," *Medical Image Analysis*, vol. 35, pp. 489–502, 2017.

[11] T. Binder, E. Tantaoui, P. Pati, R. Catena, A. Set-Aghayan, and M. Gabrani, "Multi-organ gland segmentation using deep learning," *Frontiers in Medicine*, vol. 6, p. 173, 2019.

[12] B. Ehteshami Bejnordi, M. Veta, P. Johannes van Diest, B. van Ginneken, N. Karssemeijer, G. Litjens, J. A. W. M. van der Laak, and the CAMELYON16 Consortium, "Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer," *JAMA*, vol. 318, no. 22, p. 2199–2210, 2019.

[13] G. Aresta, T. Araújoab, S. Kwok, S. Chennamsetty, M. Safwan, A. Varghese, B. Marami, M. Prastawa, M. Chan, M. Donovan, G. Fernandez, J. Zeineh, M. Kohl, C. Walz, F. Ludwig, S. Braunewell, M. Baust, Q. Vu, M. To, E. Kim, J. Kwak, S. Galal, V. Sanchez-Freire, N. Brancati, M. Frucci, D. Riccio, Y. Wang, L. Sun, K. Ma, J. Fang, I. Kone, L. Boulmane, A. Campilho, C. Eloy, A. Polónia, and P. Aguiar, "Bach: Grand challenge on breast cancer histology images," *Medical Image Analysis*, vol. 56, pp. 122–139, 2019.

[14] M. Lu, D. Williamson, T. Chen, R. Chen, M. Barbieri, and F. Mahmood, "Data efficient and weakly supervised computational pathology on whole slide images," *Nature Biomedical Engineering*, vol. 5, no. 6, pp. 555–570, 2021.

[15] G. Campanella, M. Hanna, L. Geneslaw, A. Miraflor, V. Silva, K. Busam, E. Brogi, V. Reuter, D. Klimstra, and T. Fuchs, "Clinical-grade computational pathology using weakly supervised deep learning on whole slide images," *Nature Medicine*, vol. 25, no. 8, p. 1301–1309, 2019.

[16] L. Chan, M. Hosseini, and K. Plataniotis, "A comprehensive analysis of weakly-supervised semantic segmentation in different image domains," *International Journal of Computer Vision*, vol. 129, no. 2, pp. 361–384, 2021.

[17] J. Xie, R. Liu, J. Luttrell, and C. Zhang, "Deep learning based analysis of histopathological images of breast cancer," vol. 10, p. 80, 2019.

[18] Y. Xu, Z. Jia, L.-B. Wang, Y. Ai, F. Zhang, M. Lai, and E. Chang, "Large scale tissue histopathology image classification, segmentation, and visualization via deep convolutional activation features," *Bioinformatics*, vol. 18, no. 281, 2017.

[19] Y. Xu, J.-Y. Zhu, E. Chang, M. Lai, and Z. Tu, "Weakly supervised histopathology cancer image segmentation and classification," *Medical Image Analysis*, vol. 18, no. 3, pp. 591–604, 2014.

[20] L. Hou, D. Samaras, T. Kurc, Y. Gao, J. Davis, and J. Saltz, "Patch-based convolutional neural network for whole slide tissue image classification," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition CVPR*, 2016, p. 2424–2433.

[21] Z. Jia, H. Xingyi, E. Chang, and Y. Xu, "Constrained deep weak supervision for histopathology image segmentation," *IEEE Transactions on Medical Imaging*, vol. 36, no. 11, pp. 2376–2388, 2017.

[22] G. Xu, Z. Song, Z. Sun, C. Ku, Z. Yang, C. Liu, S. Wang, J. Ma, and W. Xu, "Camel: A weakly supervised learning framework for histopathology image segmentation," in *International Conference on Computer Vision (ICCV)*, 2019, pp. 10 681–10 690.

[23] D. Ho, D. Yarlagadda, T. D'Alfonso, M. Hanna, A. Grabenstetter, P. Ntiamoah, E. Brogi, L. Tan, and T. Fuchs, "Deep multi-magnification networks for multi-class breast cancer image segmentation," *Computerized Medical Imaging and Graphics*, vol. 88, p. 101866, 2021.

[24] J. Zhang, Z. Hua, K. Yan, K. Tian, J. Yao, E. Liu, M. Liu, and X. Han, "Joint fully convolutional and graph convolutional networks for weakly-supervised segmentation of pathology images," *Medical Image Analysis*, vol. 73, p. 102183, 2021.

[25] C. Han, J. Lin, J. Mai, Y. Wang, Q. Zhang, B. Zhao, X. Chen, X. Pan, Z. Shi, Z. Xu, S. Yao, L. Yan, H. Lin, X. Huang, C. Liang, G. Han, and Z. Liu, "Multi-layer pseudo-supervision for histopathology tissue semantic segmentation using patch-level classification labels," *Medical Image Analysis*, p. 102487, 2022.

[26] L. Chan, M. Hosseini, C. Rowsell, K. Plataniotis, and S. Damaskinos, "Histosegnet: Semantic segmentation of histological tissue type in whole slide images," in *International Conference on Computer Vision (ICCV)*, 2019, pp. 10 661–10 670.

[27] J. Silva-Rodríguez, A. Colomer, and V. Naranjo, "Weglenet: A weakly-supervised convolutional neural network for the semantic segmentation of gleason grades in prostate histology images," *Computerized Medical Imaging and Graphics*, vol. 88, p. 101846, 2021.

[28] J. Ahn, S. Cho, and S. Kwak, "Weakly supervised learning of instance segmentation with inter-pixel relations," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 2204–2213.

[29] V. Anklin, P. Pati, G. Jaume, B. Bozorgtabar, A. Foncubierta-Rodriguez, J. Thiran, M. Sibony, M. Gabrani, and O. Goksel, "Learning whole-slide segmentation from inexact and incomplete labels using tissue graphs," in *Medical Image Computing and Computer Assisted Intervention (MICCAI)*, 2021, pp. 636–646.

[30] M. Ilse, J. Tomczak, and M. Welling, "Attention-based deep multiple instance learning," in *International Conference on Machine Learning (ICML)*, 2018, pp. 2127–2136.

[31] A. Myronenko, Z. Xu, D. Yang, H. Roth, and D. Xu, "Accounting for dependencies in deep learning based multiple instance learning for whole slide imaging," in *Medical Image Computing and Computer Assisted Intervention (MICCAI)*, 2021, pp. 329–338.

[32] Z. Shao, H. Bian, Y. Chen, Y. Wang, J. Zhang, X. Ji, and Y. Zhang, "Transmil: Transformer based correlated multiple instance learning for whole slide image classification," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.

[33] Y. Lee, J. Park, S. Oh, K. Shin, J. Sun, M. Jung, C. Lee, H. Kim, J. Chung, K. Moon, and S. Kwon, "Derivation of prognostic contextual histopathological features from whole-slide images of tumours via graph deep learning," *Nature Biomedical Engineering*, pp. 1–15, 2022.

[34] J. Lipkova, T. Chen, M. Lu, R. Chen, M. Shady, M. Williams, J. Wang, Z. Noor, R. Mitchell, M. Turan, G. Coskun, F. Yilmaz, D. Demir, D. Nart, K. Başak, N. Turhan, S. Ozkara, Y. Banz, K. Odening, and F. Mahmood, "Deep learning-enabled assessment of cardiac allograft rejection from endomyocardial biopsies," *Nature Medicine*, vol. 28, pp. 575–582, 2022.

[35] K. Thandiackal, B. Chen, P. Pati, G. Jaume, D. Williamson, M. Gabrani, and O. Goksel, "Differentiable zooming for multiple instance learning on whole-slide images," in *European Conference on Computer Vision (ECCV)*, 2022.

[36] F. Kong and R. Henao, "Efficient classification of very large images with tiny objects," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 2384–2394.

[37] H. Lee and W. Jeong, "Scribble2label: Scribble-supervised cell segmentation via self-generating pseudo-labels with consistency," in *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2020, pp. 14–23.

[38] H. Qu, P. Wu, Q. Huang, J. Yi, Z. Yan, K. Li, G. M. Riedlinger, S. De, S. Zhang, and D. N. Metaxas, "Weakly supervised deep nuclei segmentation using partial points annotation in histopathology images," *IEEE Transactions on Medical Imaging*, vol. 39, no. 11, pp. 3655–3666, 2020.

[39] T.-A. Yen, H.-C. Hsu, P. Pati, M. Gabrani, A. Foncubierta-Rodríguez, and P.-C. Chung, "Ninepins: Nuclei instance segmentation with point annotations," *arXiv preprint arXiv:2006.13556*, 2020.

[40] J.-M. Bokhorst, H. Pinckaers, P. van Zwam, I. Nagtegaal, J. van der Laak, and F. Ciompi, "Learning from sparsely annotated data for semantic segmentation in histopathology images," in *Medical Imaging with Deep Learning (MIDL)*, 2018.

[41] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, N. Navab, J. Hornegger, W. M. Wells, and A. F. Frangi, Eds. Cham: Springer International Publishing, 2015, pp. 234–241.

[42] M. Aubreville, "Quantifying the scanner-induced domain gap in mitosis detection," in *Medical Imaging with Deep Learning (MIDL)*, 2021.

[43] M. Macenko, M. Niethammer, J. Marron, D. Borland, J. Woosley, X. Guan, C. Schmitt, and N. Thomas, "A method for normalizing histology slides for quantitative analysis," in *IEEE International Symposium on Biomedical Imaging (ISBI)*, 2009, pp. 1107–1110.

[44] A. Vahadane, T. Peng, S. Albarqouni, M. Baust, K. Steiger, A. Schlitter, A. Sethi, I. Esposito, and N. Navab, "Structure-preserving color normalization and sparse stain separation for histological images," *IEEE Transactions on Medical Imaging*, vol. 35, no. 8, pp. 1962–1971, 2016.

[45] M. Stanisavljevic, A. Anghel, N. Papandreou, S. Andani, P. Pati, H. Ruschoff, P. Wild, M. Gabrani, and H. Pozidis, "A fast and scalable pipeline for stain normalization of whole-slide images in histopathology," in *European Conference on Computer Vision (ECCV) Workshops*, 2018.

[46] J. Ren, I. Hacihaliloglu, E. Singer, D. Foran, and X. Qi, "Unsupervised domain adaptation for classification of histopathology whole-slide images," *Frontiers in Bioengineering and Biotechnology*, vol. 7, p. 102, 2019.

[47] D. Tellez, M. Balkenhol, N. Karssemeijer, G. Litjens, J. van der Laak, and F. Ciompi, "H and e stain augmentation improves generalization of convolutional networks for histopathological mitosis detection," in *Medical Imaging*, vol. 10581, 2018, pp. 264–270.

[48] K. Faryna, J. van der Laak, and G. Litjens, "Tailoring automated data augmentation to h&e-stained histopathology," in *Medical Imaging with Deep Learning (MIDL)*, 2021.

[49] S. Otálora, M. Atzori, V. Andrearczyk, A. Khan, and H. Müller, "Staining invariant features for improving generalization of deep convolutional neural networks in computational pathology," *Frontiers in Bioengineering and Biotechnology*, 2019.

[50] R. Yamashita, J. Long, S. Banda, J. Shen, and D. L. Rubin, "Learning domain-agnostic visual representation for computational pathology using medically-irrelevant style transfer augmentation," *IEEE Transactions on Medical Imaging*, vol. 40, no. 12, pp. 3945–3954, 2021.

[51] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 2096–2030, 2016.

[52] M. Aubreville, C. Bertram, S. Jabari, C. Marzahl, R. Klopfleisch, and A. Maier, "Inter-species, inter-tissue domain adaptation for mitotic figure assessment: Learning new tricks from old dogs," *Bildverarbeitung für die Medizin 2020*, 2020.

[53] N. Hashimoto, D. Fukushima, R. Koga, Y. Takagi, K. Ko, K. Kohno, M. Nakaguro, S. Nakamura, H. Hontani, and I. Takeuchi, "Multi-scale domain-adversarial multiple-instance cnn for cancer subtype classification with unannotated histopathological images," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 3852–3861.

[54] M. E. Tschuchnig, G. J. Oostingh, and M. Gadermayr, "Generative adversarial networks in digital pathology: a survey on trends and future potential," *Patterns*, vol. 1, no. 6, 2020.

[55] C. Guo, G. Pleiss, Y. Sun, and K. Weinberger, "On calibration of modern neural networks," in *International Conference on Machine Learning (ICML)*, 2019, pp. 1321–1330.

[56] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

[57] S. Fort, H. Hu, and B. Lakshminarayanan, "Deep ensembles: A loss landscape perspective," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

[58] W. L. Hamilton, "Graph representation learning," *Synthesis Lectures on Artifical Intelligence and Machine Learning*, vol. 14, no. 3, pp. 1–159, 2020.

[59] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" in *International Conference on Learning Representations (ICLR)*, 2019.

[60] J. Nixon, M. Dusenberry, L. Zhang, G. Jerfel, and D. Tran, "Measuring calibration in deep learning." in *CVPR Workshops*, vol. 2, no. 7, 2019.

[61] M. DeGroot and S. Fienberg, "The comparison and evaluation of forecasters," *Journal of the Royal Statistical Society: Series D (The Statistician)*, vol. 32, no. 1-2, pp. 12–22, 1983.

[62] M. Veta, J. Pluim, P. Diest, and M. Viergever, "Breast cancer histopathology image analysis: A review," *IEEE Transactions on Biomedical Engineering*, no. 5, pp. 1400–1411, 2014.

[63] D. Tellez, G. Litjens, P. Bándi, W. Bulten, J.-M. Bokhorst, F. Ciompi, and J. van der Laak, "Quantifying the effects of data augmentation and stain color normalization in convolutional neural networks for computational pathology," *Medical Image Analysis*, vol. 58, p. 101544, 2019.

[64] B. Bejnordi, G. Litjens, M. Hermsen, N. Karssemeijer, and J. van der Laak, "A multi-scale superpixel classification approach to the detection of regions of interest in whole slide histopathology images," in *Medical Imaging 2015: Digital Pathology*, vol. 9420, 2015, pp. 99–104.

[65] P. Pati, G. Jaume, A. F. Rodriguez, and M. Gabrani, "Interpretation of whole-slide images in digital pathology," 2022, uS Patent App. 16/953,377.

[66] P. Pati, G. Jaume, A. Foncubierta-Rodríguez, F. Feroce, A. Anniciello, G. Scognamiglio, N. Brancati, M. Fiche, E. Dubruc, D. Riccio, M. Bonito, G. Pietro, G. Botti, J. Thiran, M. Frucci, O. Goksel, and M. Gabrani, "Hierarchical graph representations in digital pathology," *Medical Image Analysis*, vol. 75, p. 102264, 2022.

[67] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Süsstrunk, "Slic superpixels compared to state-of-the-art superpixel methods," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 11, pp. 2274–2282, 2012.

[68] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 4510–4520.

[69] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009, pp. 248–255.

[70] F. Potjer, "Region adjacency graphs and connected morphological operators," in *Mathematical Morphology and its Applications to Image and Signal Processing. Computational Imaging and Vision*, vol. 5, 1996, p. 111–118.

[71] K. Xu, C. Li, Y. Tian, T. Sonobe, K.-i. Kawarabayashi, and S. Jegelka, "Representation learning on graphs with jumping knowledge networks," in *International Conference on Machine Learning (ICML)*, 2018.

[72] P. Pope, S. Kolouri, M. Rostami, C. Martin, and H. Hoffmann, "Explainability methods for graph convolutional neural networks," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 10 764–10 773.

[73] G. Jaume, P. Pati, B. Bozorgtabar, A. Foncubierta-Rodríguez, F. Feroce, A. Anniciello, T. Rau, J. Thiran, M. Gabrani, and O. Goksel, "Quantifying explainers of graph neural networks in computational pathology," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 8106–8116.

[74] R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, and D. Batra, "Grad-CAM : Visual Explanations from Deep Networks," in *International Conference on Computer Vision (ICCV)*, 2017, pp. 618–626.

[75] J. Silva-Rodríguez, A. Colomer, M. Sales, R. Molina, and V. Naranjo, "Going deeper through the gleason scoring scale: An automatic end-to-end system for histology prostate grading and cribriform pattern detection," *Computer Methods and Programs in Biomedicine*, vol. 195, p. 105637, 2020.

[76] W. Bulten, H. Pinckaers, H. Boven, R. Vink, T. Bel, B. Ginneken, J. van der Laak, C. Hulsbergen-van de Kaa, and G. Litjens, "Automated deep-learning system for gleason grading of prostate cancer using biopsies: a diagnostic study," *Lancet Oncology*, vol. 21, no. 2, pp. 233–241, 2020.

[77] W. Bulten, K. Kartasalo, P. Chen, P. Ström, H. Pinckaers, K. Nagpal, Y. Cai, D. Steiner, H. van Boven, R. Vink, C. Hulsbergen-van de Kaa, J. van der Laak, M. Amin, A. Evans, T. van der Kwast, R. Allan, P. Humphrey, H. Grönberg, H. Samaratunga, B. Delahunt, T. Tsuzuki, T. Häkkinen, L. Egevad, M. Demkin, S. Dane, F. Tan, M. Valkonen, G. Corrado, L. Peng, C. Mermel, P. Ruusuvuori, G. Litjens, M. Eklund, and the PANDA challenge consortium, "Artificial intelligence for diagnosis and gleason grading of prostate cancer: the panda challenge," *Nature Medicine*, vol. 28, p. 154–163, 2022.

[78] P. Ström, K. Kartasalo, H. Olsson, L. Solorzano, B. Delahunt, D. Berney, D. Bostwick, A. Evans, D. Grignon, P. Humphrey, K. Iczkowski, J. Kench, G. Kristiansen, T. Van der Kwast, K. Leite, J. McKenney, J. Oxley, C.-C. Pan, H. Samaratunga, and M. Eklund, "Pathologist-level grading of prostate biopsies with artificial intelligence," *Bioinformatics*, 2019.

[79] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *Neural Information Processing Systems (NeurIPS)*, 2019, pp. 8024–8035.

[80] M. Wang, L. Yu, D. Zheng, Q. Gan, Y. Gai, Y. Zihao, M. Li, J. Zhou, Q. Huang, C. Ma, Z. Huang, Q. Guo, H. Zhang, H. Lin, J. Zhao, J. Li, A. Smola, and Z. Zhang, "Deep graph library: Towards efficient and scalable deep learning on graphs," in *International Conference on Learning Representations (ICLR) workshop*, 2019.

[81] G. Jaume, P. Pati, V. Anklin, A. Foncubierta, and M. Gabrani, "Histocartography: A toolkit for graph analytics in digital pathology," in *MICCAI Workshop on Computational Pathology (MICCAI-W)*, 2021, pp. 117–128.

[82] J. Epstein, W. Allsbrook, M. Amin, L. Egevad, and the ISUP Grading Committee, "The 2005 international society of urological pathology (isup) consensus conference on gleason grading of prostatic carcinoma," *The American Journal of Surgical Pathology*, vol. 29, no. 9, pp. 1228–1242, 2005.

[83] J. Epstein, L. Egevad, M. Amin, B. Delahunt, J. Srigley, P. Humphrey, and the Grading Committee, "The 2014 international society of urological pathology (isup) consensus conference on gleason grading of prostatic carcinoma: Definition of grading patterns and proposal for a new grading system," *The American Journal of Surgical Pathology*, vol. 40, no. 2, pp. 244–252, 2016.

[84] A. Gomariz, T. Portenier, C. Nombela-Arrieta, and O. Goksel, "Probabilistic spatial analysis in quantitative microscopy with uncertainty-aware cell detection using deep bayesian regression of density maps," *Science Advances*, vol. 8, no. 5, 2021.

[85] G. Liang, Y. Zhang, X. Wang, and N. Jacobs, "Trainable calibration measures for neural networks from kernel mean embeddings," in *International Conference on Machine Learning (ICML)*, 2018.

[86] T. A. Ozkan, A. T. Eruyar, O. O. Cebeci, O. Memik, L. Ozcan, and I. Kuskonmaz, "Interobserver variability in gleason histological grading of prostate cancer," *Scandinavian journal of urology*, vol. 50, no. 6, pp. 420–424, 2016.

[87] E. N. Salmo, "An audit of inter-observer variability in gleason grading of prostate cancer biopsies: The experience of central pathology review in the north west of england," *Integr Cancer Sci Ther*, vol. 2, no. 2, pp. 104–106, 2015.

[88] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne," *Journal of machine learning research*, vol. 9, no. 11, 2008.

[89] L. Cheng, D. Davidson, H. Lin, and M. Koch, "Percentage of gleason pattern 4 and 5 predicts survival after radical prostatectomy," *Cancer*, vol. 110, no. 9, pp. 1967–1972, 2007.

[90] C. Huang, M. Kong, M. Zhou, A. Rosenkrantz, S. Taneja, J. Melamed, and F. Deng, "Gleason score 3+4=7 prostate cancer with minimal quantity of gleason pattern 4 on needle biopsy is associated with low-risk tumor in radical prostatectomy specimen," *The American Journal of Surgical Pathology*, vol. 38, no. 8, pp. 1096–1101, 2014.

[91] B. Choy, S. Pearce, B. Anderson, A. Shalhav, G. Zagaja, S. Eggener, and G. Paner, "Prognostic significance of percentage and architectural types of contemporary gleason pattern 4 prostate cancer in radical prostatectomy," *The American Journal of Surgical Pathology*, vol. 40, no. 10, pp. 1400–1406, 2016.

[92] M. Sharma and H. Miyamoto, "Percent gleason pattern 4 in stratifying the prognosis of patients with intermediate-risk prostate cancer," *Translational Andrology and Urology*, vol. 7, no. 4, 2018.

[93] M. Amin, D. Lin, J. Gore, J. Srigley, H. Samaratunga, L. Egevad, M. Rubin, J. Nacey, H. Carter, L. Klotz, H. Sandler, A. Zietman, S. Holden, R. Montironi, P. Humphrey, A. Evans, J. Epstein, B. Delahunt, J. McKenney, D. Berney, T. Wheeler, A. Chinnaiyan, L. True, B. Knudsen, and M. Hammond, "The critical role of the pathologist in determining eligibility for active surveillance as a management option in patients with prostate cancer: consensus statement with recommendations supported by the college of american pathologists, international society of urological pathology, association of directors of anatomic and surgical pathology, the new zealand society of pathologists, and the prostate cancer foundation," *Archives of Pathology & Laboratory Medicine*, vol. 138, no. 10, pp. 1387–1405, 2014.