# Quantifying Explainers of Graph Neural Networks in Computational Pathology

Guillaume Jaume[1,2]*, Pushpak Pati[1,3,*], Behzad Bozorgtabar[2],
Antonio Foncubierta[1], Anna Maria Anniciello[4], Florinda Feroce[4], Tilman Rau[5],
Jean-Philippe Thiran[2], Maria Gabrani[1], Orcun Goksel[3,6]
[1]IBM Research Zurich, [2]EPFL Lausanne, [3]ETH Zurich,
[4]Fondazione Pascale, [5]University of Bern, [6]Uppsala University

{gja,pus}@zurich.ibm.com

## Abstract

*Explainability of deep learning methods is imperative to facilitate their clinical adoption in digital pathology. However, popular deep learning methods and explainability techniques (explainers) based on pixel-wise processing disregard biological entities' notion, thus complicating comprehension by pathologists. In this work, we address this by adopting biological entity-based graph processing and graph explainers enabling explanations accessible to pathologists. In this context, a major challenge becomes to discern meaningful explainers, particularly in a standardized and quantifiable fashion. To this end, we propose herein a set of novel quantitative metrics based on statistics of class separability using pathologically measurable concepts to characterize graph explainers. We employ the proposed metrics to evaluate three types of graph explainers, namely the layer-wise relevance propagation, gradient-based saliency, and graph pruning approaches, to explain Cell-Graph representations for Breast Cancer Subtyping. The proposed metrics are also applicable in other domains by using domain-specific intuitive concepts. We validate the qualitative and quantitative findings on the BRACS dataset, a large cohort of breast cancer RoIs, by expert pathologists. The code, data, and models can be accessed here[1].*

## 1. Introduction

Histopathological image understanding has been revolutionized by recent machine learning advancements, especially deep learning (DL) [8, 55]. DL has catered to increasing diagnostic throughput as well as a need for high predictive performance, reproducibility and objectivity. However, such advantages come at the cost of a reduced transparency in decision-making processes [30, 63, 23]. Considering the

---

*denotes equal contribution

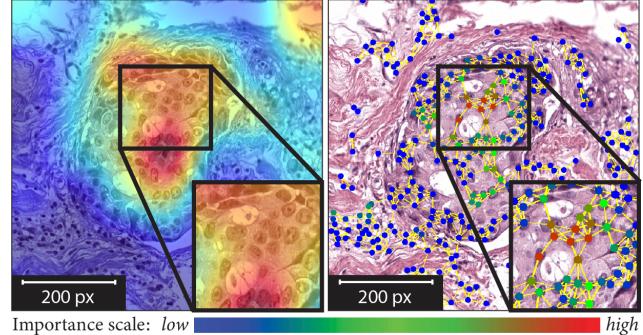[1]https://github.com/histocartography/patho-quant-explainer



Figure 1. Sample explanations produced by pixel- and entity-based explainability techniques for a ductal carcinoma *in situ* RoI.

need for reasoning any clinical decision, it is imperative to enable the explainability of DL decisions to pathologists.

Inspired by the explainability techniques (explainers) for DL model decisions on natural images [59, 70, 69, 6, 42, 54, 35, 75, 11, 34], several explainers have been implemented in digital pathology, such as feature attribution [10, 9, 23], concept attribution [21], and attention-based learning [40]. However, pixel-level explanations, exemplified in Figure 1, pose several notable issues, including: (1) a pixel-wise analysis disregards the notion of biological tissue entities, their topological distribution, and inter-entity interactions; (2) a typical patch-based DL processing and explainer fail to accommodate complete tumor macro-environment information; and (3) pixel-wise visual explanations tend to be blurry. Explainability in entity space is thus a natural choice to address the above issues. To that end, an entity graph representation is built for a histology image, where nodes and edges denote biological entities and inter-entity interactions followed by a Graph Neural Network (GNN) [37, 66]. The choice of entities, such as cells [22, 74, 46], tissues [46] or others, can be task-dependent. Subsequently, explainers for graph-structured data [7, 48, 67] applied to the entity graphs highlight responsible entities for the concluded diagnosis, thereby generating intuitive explanations for pathologists.

In the presence of various graph explainers producing distinct explanations for an input, it is crucial to discern the explainer that best fits the explainability definition [4]. In the context of computational pathology, explainability is defined as making the DL decisions understandable to pathologists [30]. To this end, the qualitative evaluation of explainers' explanations by pathologists is the candid measure. However, it requires evaluations by task-specific expert pathologists, which is subjective, time-consuming, cumbersome, and expensive. Additionally, though the explanations are intuitive, they do not relate to pathologist-understandable terminologies, *e.g.* "How big are the important nuclei?", "How irregular are their shape?" etc., which toughens the comprehensive analysis. These bottlenecks undermine not only any qualitative assessment but also quantitative metrics requiring user interactions [41]. Furthermore, expressing the quantitative metrics in user-understandable terminologies [4] is fundamental to achieve interpretability [16, 43]. To this end, the most popular quantitative metric, explainer *fidelity* [50, 15, 51, 29, 41, 48], is not satisfactory. Moreover, explainers intrinsically maintain high-*fidelity*, *e.g.* GNNEXPLAINER [67] produces an explanation to match the GNN's prediction on the original graph.

Thus, we propose a set of novel user-independent quantitative metrics expressing pathologically-understandable *concepts*. The proposed metrics are based on class separability statistics using such *concepts*, and they are applicable in other domains by incorporating domain-specific *concepts*. We use the proposed metrics to evaluate three types of graph-explainers, (1) graph pruning: GNNEXPLAINER [67, 31], (2) gradient-based saliency: GRAPHGRAD-CAM [54, 48], GRAPHGRAD-CAM++ [11], (3) layer-wise relevance propagation: GRAPHLRP [6, 42, 53], for explaining Cell-Graphs [22] in Breast Cancer Subtyping as shown in Figure 1. Our specific contributions in this work are:

- A set of novel quantitative metrics based on the statistics of class separability using domain-specific *concepts* to characterize graph explainability techniques. To the best of our knowledge, our metrics are the first of their kind to quantify explainability based on domain-understandable terminologies;

- Explainability in computational pathology using pathologically intuitive entity graphs;

- Extensive qualitative and quantitative assessment of various graph explainability techniques in computational pathology, with a validation of the findings by expert pathologists.

## 2. Related work

**Graphs in Digital Pathology:** Graph-based tissue image analysis effectively describes a tissue environment by incorporating morphology, topology, and tissue components interactions. To this end, cell-graph (CG) is the most popular graph representation, where nodes and edges depict cells and cellular interactions [22]. Cell morphology is embedded in the nodes via hand-crafted features [22, 74, 46] or DL features [12, 47]. The graph topology is heuristically defined using k-Nearest Neighbors, probabilistic modeling, Waxman model etc. [58] Subsequently, the CGs are processed by classical machine learning [58, 57, 56] or GNN [74, 12, 18, 46] to map the tissue structure to function relationship. Recently, improved graph-representations using patches [5], tissue components [46], and hierarchical cell-to-tissue relations [46] are proposed to enhance the structure-function mapping. Other graph-based applications in computational pathology include cellular community detection [32], whole-slide image classification [72, 1] etc. Intuitively, a graph representation utilizes pathologically relevant entities to represent a tissue specimen, which allows pathologists to readily relate with the input, also enabling them to include any task-specific prior knowledge.

**Explainability in Digital Pathology:** Explainability is an integral part of pathological diagnosis. Though DL solutions have achieved remarkable diagnostic performance, their lack of explainability is unacceptable in the medical community [63]. Recent studies have proposed visual explanations [23] and salient regions [10, 23] using feature-attribution techniques [54, 11]. Differently, concept-attribution technique [21] evaluates the sensitivity of network output w.r.t. quantifiable image-level pathological *concepts* in patches. Although such explanations are pathologist-friendly, image-level *concepts* are neither fit nor meaningful for real-world large histology images that contain many localized concepts. Furthermore, attention-based learning [40], and multimodal mapping between image and diagnostic report [71] are devised to localize network attention. However, the pixel-wise and patch-based processing in all the aforementioned techniques ignore biological entities' notion; thus, they are not easily understood by pathologists. Separately, the earlier stated entity graph-based processing provides an intuitive platform for pathologists. However, research on explainability and visualization using entity graphs has been scarce: CGC-Net [74] analyzes cluster assignment of nodes in CG to group them according to their appearance and tissue types. CGExplainer [31] introduces a post-hoc graph-pruning explainer to identify decisive cells and interactions. Robust spatial filtering [61] utilizes an attention-based GNN and node occlusion to highlight cell contributions. No previous work has comprehensively analyzed and quantified graph explainers in computational pathology while expressing explanations in a pathologist-understandable form to the best of our knowledge. This gap between the existing and desired explainability of DL outputs in digital pathology motivates our work herein.
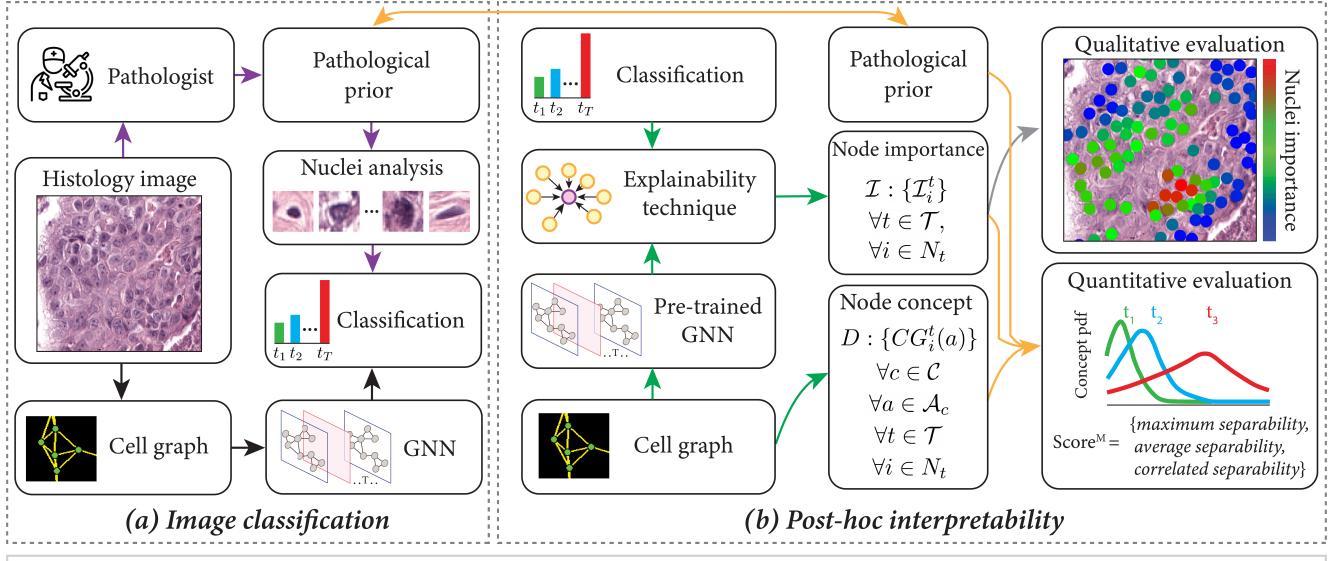
Figure 2. Overview of the proposed framework. (a) presents pathologist, and entity-based (cell-graph + GNN) diagnosis of a histology image. (b) presents nuclei-level pathologically relevant *concept* measure $D$, a post-hoc graph explainability technique to derive nuclei-level importance $\mathcal{I}$ for *concepts* $\mathcal{C}$, measurable *attributes* $\mathcal{A}_c$, and classes $\mathcal{T}$. $D$, $\mathcal{I}$ and prior pathological knowledge defining *concepts'* relevance are utilized to propose a novel set of quantitative metrics to evaluate the explainer quality in pathologist-understandable terms.

# 3. Method

In this section, we present entity graph processing, explainability methods, and our proposed evaluation metrics. First, we transform a histology region-of-interest (RoI) into a *biological entity graph*. Second, we introduce a "black-box" GNN that maps the *entity graph* to a corresponding class label. Third, we employ a post-hoc graph explainer to generate explanations. Finally, we perform a qualitative and quantitative assessments of the generated explanations. An overview of the methodology is shown in Figure 2.

## 3.1. Entity graph notations

We define an attributed undirected entity graph $G := (V, E, H)$ as a set of nodes $V$, edges $E$, and node attributes $H \in \mathbb{R}^{|V| \times d}$. $d$ denotes the number of attributes per node, and $|.|$ denotes set cardinality. The graph topology is defined by a symmetric graph adjacency, $A \in \mathbb{R}^{|V| \times |V|}$, where $A_{u,v} = 1$ if $e_{uv} \in E$. We denote the neighborhood of a node $v \in V$ as $\mathcal{N}(v) := \{u \in V \mid v \in V, \ e_{uv} \in E \}$. We denote a set of graphs as $\mathcal{G}$.

## 3.2. Entity graph construction

Our methodology begins with transforming RoIs into entity graphs. It ensures the method's inputs are pathologically interpretable, as the inputs consist of biologically-defined objects that pathologists can directly *relate-to* and *reason-with*. Thus, image-to-graph conversion moves from *uninterpretable* to *interpretable* input space. In this work,

we consider cells as entities, thereby transforming RoIs into cell-graphs (CGs). A CG nodes and edges capture the morphology of cells and cellular interactions. A CG topology acquires both tissue micro and macro-environment, which is crucial for characterizing cancer subtypes.

First, we detect nuclei in a RoI at $40\times$ magnification using Hover-Net [20], a nuclei segmentation algorithm pre-trained on MoNuSeg [39]. We process patches of size $72\times72$ around the nuclei by ResNet34 [27] pre-trained on ImageNet [14] to produce nuclei visual attributes. We further concatenate nuclei spatial attributes, *i.e.* nuclei centroids min-max normalized by RoI dimension. The nuclei and their attributes (visual and spatial) define the nodes and node attributes of the CG, respectively. Following prior work [46], we construct the CG topology by employing thresholded $k$-Nearest Neighbors algorithm. We set $k = 5$, and prune the edges longer than 50 pixels (12.5 $\mu$m). The CG-topology encodes how likely two nearby nuclei will interact [17]. A CG example is presented in Figure 1.

## 3.3. Entity graph learning

Given $\mathcal{G}$, the set of CGs, the aim is to infer the corresponding cancer subtypes. We use GNNs [52, 13, 37, 24, 64, 68, 19], the conceptual analogous of 2D convolution for graph-structured data, to classify the CGs. A GNN layer follows two steps: for each node $v \in V$, (i) *aggregation step*: the states of neighboring nodes, $\mathcal{N}(v)$, are aggregated via a differentiable and permutation-invariant operator to produce $a(v) \in \mathbb{R}^d$, then, (ii) *update step*: the state of $v$

is updated by combining the current node state $h(v) \in \mathbb{R}^d$ and the aggregated message $a(v)$ via another differentiable operator. After $L$ iterations, *i.e.* the number of GNN layers, a *readout step* is employed to merge all the node states via a differentiable and permutation-invariant function to result in a fixed-size graph embedding. Finally, the graph embeddings are processed by a classifier to predict the class label.

In this work, we use a flavor of Graph Isomorphism Network (GIN) [66], that uses *mean* and a *multi-layer perceptron* (MLP) in the *aggregation* and *update* step respectively. Formally, we define a layer as,

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

where $h(v)$ denotes features of node $v$, and $l \in \{1, ..., L\}$. Our GNN consists of 3-GIN layers, with each layer including a 2-layer MLP. The dimension of latent node embeddings is fixed to 64 for all layers. We use *mean* operation in *readout step*, and feed the graph embedding to a 2-layer MLP classifier. The GNN is trained end-to-end by minimizing cross-entropy loss between predicted logits and target cancer subtypes. We emphasize that the entity-based processing follows a pathologist's diagnostic procedure that identifies diagnostically relevant nuclei and analyzes cellular morphology and topology in a RoI, as shown in Figure 2.

### 3.4. Post-hoc graph explainer

We generate an explanation per entity graph by employing post-hoc graph explainers. The explanations allow to evaluate the pathological relevance of black-box neural network reasoning. Specifically, we aim to evaluate the agreement between the pathologically relevant set of nuclei in a RoI, and the explainer identified set of important nuclei, *i.e.* nuclei driving the prediction, in corresponding CG. In this work, we consider three types of graph explainers for explaining CGs, which follow similar operational setting, *i.e.* (i) input data are attributed graphs, (ii) a GNN is trained *a priori* to classify the input data, and (iii) each data point can be inferred independently to produce an explanation. We present the graph explainers in the following sections and their detailed mathematical formulations in the Appendix.

**GRAPHLRP:** Layerwise relevance propagation (LRP) [6] propagates the output logits backward in the network using a set of propagation rules to quantify the positive contribution of input pixels for a certain prediction. Specifically, LRP assigns an importance score to each neuron such that the output logit relevance is preserved across layers. While initially developed for explaining fully-connected layers, LRP can be extended to GNN by treating the GNN *aggregation step* as a fully connected layer that projects the graph adjacency matrix on the node attributes as in [53]. LRP outputs per-node importance.

**GRAPHGRAD-CAM:** GRAD-CAM [54] is a feature attribution approach designed for explaining CNNs operating on images. It produces class activation explanation following two steps. First, it assigns weights to each channel of a convolutional layer $l$ by computing the gradient of the targeted output logit w.r.to each channel in layer $l$. Second, importance of the input elements are computed by the weighted combination of the forward activations at each channel in layer $l$. The extension to GNN is straightforward [48], and only requires to compute the gradient of the predicted logits w.r.to a GNN layer. Following prior work [48], we take the average of node-level importance-maps obtained from all the GNN layers $l \in \{1, ..., L\}$ to produce smooth per-node importance.

**GRAPHGRAD-CAM++:** GRAD-CAM++ [11] is an increment on GRAD-CAM by including spatial contributions into the channel-wise weight computation of a convolutional layer. The extension allows weighting the contribution by each spatial location at a layer for improved spatial localization. The spatial locations in a convolutional layer are analogous to the size of the graph in a GNN layer. With this additional consideration, we propose an extension of GRAD-CAM++ to graph-structured data.

**GNNEXPLAINER:** GNNEXPLAINER [67, 31] is a graph pruning approach that aims to find a compact sub-graph $G_s \subset G$ such that mutual information between $G_s$ and GNN prediction of $G$ is maximized. Sub-graph $G_s$ is regarded as the explanation for the input graph $G$. GNNEXPLAINER can be seen as a feature attribution technique with binarized node importances. To address the combinatorial nature of finding $G_s$, GNNEXPLAINER formulates it as an optimization problem that learns a mask to activate or deactivate parts of the graph. [31] reformulates the initial approach in [67] to learn a mask over the nodes instead of edges. The approach in [31] is better suited for pathology as the nodes, *i.e.* biological entities, are more intuitive and substantial for disease diagnosis than heuristically-defined edges. The optimization for an entity graph results in per-node importance.

### 3.5. Quantitative metrics for graph explainability

In the presence of several graph explainers producing distinct explanations for an input, it is imperative to discern the explainer that produces the most pathologically-aligned explanation. Considering the limitations of existing qualitative and quantitative measures presented in Section 1, we propose a novel set of quantitative metrics based on class separability statistics using pathologically relevant *concepts*. Intuitively, a good explainer should emphasize the relevant *concepts* that maximize the class separation. Details of the metric evaluations are presented as follows.

**Input:** A graph explainer outputs an explanation, *i.e.* node-level *importance* $\mathcal{I}$, for an input CG. To quantify a
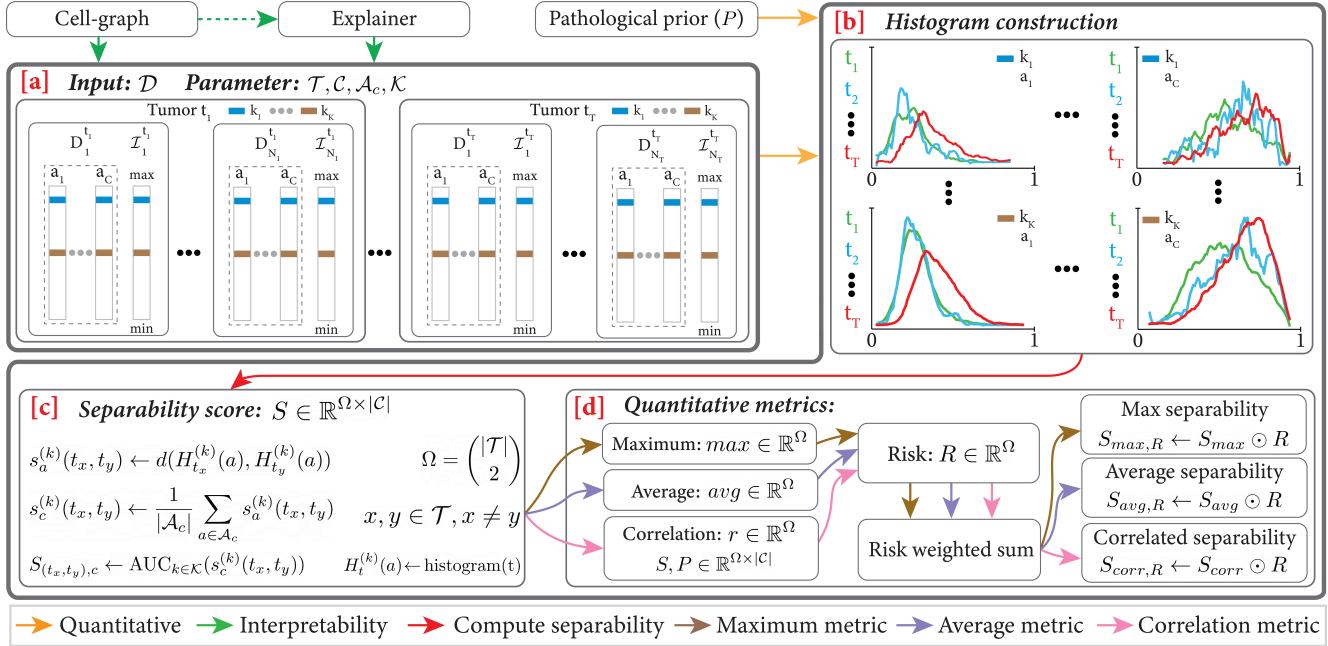
Figure 3. Overview of proposed quantitative assessment. (a) presents input dataset $\mathcal{D}$, and parameters *concepts* $\mathcal{C}$, measurable *attributes* $\mathcal{A}_c$, classes $\mathcal{T}$, and importance thresholds $\mathcal{K}$. For simplicity $|\mathcal{A}_c| = 1, \forall c \in \mathcal{C}$ in this figure. (b) shows histogram probability densities for $\forall a \in \mathcal{A}_c, \forall k \in \mathcal{K}, \forall t \in \mathcal{T}$. (c) displays the algorithm for computing class separability score $S$. (d) presents the algorithm for computing the proposed class separability-based risk-weighted quantitative metrics.

*concept* $c \in \mathcal{C}$, $\mathcal{C}$ denoting the set of *concepts*, we measure nuclear *attributes* $a \in \mathcal{A}_c$ for each nucleus in CG, *e.g.*, for $c =$ *nuclear shape*, we measure $\mathcal{A}_c = \{$*perimeter, roughness, eccentricity, circularity*$\}$. We create a dataset $\mathcal{D} = \bigcup_{t \in \mathcal{T}} \mathcal{D}_t$, $\mathcal{T}$ denoting the set of cancer subtypes. We define $\mathcal{D}_t := \{(D_i^t, \mathcal{I}_i^t) | i = 1, \ldots, N_t\} \forall t \in \mathcal{T}$, where $N_t$ is the number of CGs for tumor type $t$. $\mathcal{I}_i^t$ and $D_i^t$ are, respectively, the sorted importance matrix for a CG indexed by $i$ and corresponding node-level attribute matrix. To perform inter-concept comparisons, we conduct *attribute*-wise normalization across all $D_i^t \forall t, i$. In order to compare different explainers, we conduct CG-wise normalization of $\mathcal{I}$. The structure of input dataset $\mathcal{D}$ is presented in Figure 2(a).

Note that the notion of important nuclei vary (1) per-CG since the number of nodes vary across CGs, and (2) per-explainer. Hence, selecting a *fixed* number of important nuclei per-CG and per-explainer is not meaningful. To overcome this issue, we assess different number of important nuclei $k \in \mathcal{K}$, selected based on node importances, per-CG and per-explainer. In the following sections we will show how to aggregate the results for a given explainer.

**Histogram construction:** Given the input dataset $\mathcal{D}$, and parameters $\mathcal{K}, \mathcal{C}, \mathcal{A}_c, \mathcal{T}$, we apply threshold $k \in \mathcal{K}$ on $\mathcal{I}_i^t, \forall t \in \mathcal{T}, \forall i \in N_t$ to select CG-wise most important nuclei. The cancer subtype-wise selected set of nuclei data from $\mathcal{D}$ are used to construct histograms $H_t^{(k)}(a), \forall a \in \mathcal{A}_c, \forall c \in \mathcal{C}$ and $\forall t \in \mathcal{T}$. For histogram $H_t^{(k)}(a)$, bin-edges are

decided by quantizing the complete range of *attribute* $a$, *i.e.* $\mathcal{D}(a)$, by a fixed step size. We convert each $H_t^{(k)}(a)$ into a probability density function. Similarly, sets of histograms are constructed by applying different thresholds $k \in \mathcal{K}$. Sample histograms are shown in Figure 2(b).

**Separability Score ($S$):** Given two classes $t_x, t_y \in \mathcal{T}$ and corresponding probability density functions $H_{t_x}^{(k)}(a)$ and $H_{t_y}^{(k)}(a)$, we compute *class separability* $s_a^{(k)}(t_x, t_y)$ based on optimal transport as the Wasserstein distance between the two density functions. We average $s_a^{(k)}(t_x, t_y)$ over all $a \in \mathcal{A}_c$ to obtain a score $s_c^{(k)}(t_x, t_y)$ for *concept* $c$ and threshold $k$. Finally, we compute the area-under-the-curve (AUC) over the threshold range $\mathcal{K}$ to get the aggregated class separability $S_{(t_x, t_y), c}$ for a *concept* $c$. The class separability score indicates the significance of *concept* $c$ for the purpose of separating $t_x$ and $t_y$. Thus, separability scores can be used to compare different *concepts* and to identify relevant ones for differentiating $t_x$ and $t_y$. A pseudo-algorithm is presented in Algorithm 1, and illustrated in Figure 2(c). A separability matrix $S \in \mathbb{R}^{\Omega \times |\mathcal{C}|}$ is built by computing class separability scores for all pair-wise classes, *i.e.* $\forall (t_x, t_y) \in \Omega := \binom{|\mathcal{T}|}{2}$ and $\forall c \in \mathcal{C}$.

**Statistics of Separability Score:** Since explainability is not uniquely defined, we include multiple metrics highlighting different facets. We compute three separability statistics $\forall (t_x, t_y) \in \Omega$ using $S$ as given in Equation (2), *i.e.* (1) *max-*

*imum*: the utmost separability, (2) *average*: the expected separability. These two metrics encode (model+explainer)'s focus, *i.e.* "how much the black-box model implicitly uses the *concepts* for class separability?" (3) *correlation*: encodes the agreement between (model+explainer)'s focus and pathological prior $P$. $P \in \mathbb{R}^{\Omega \times |\mathcal{C}|}$ signifies the relevance $\forall c \in \mathcal{C}$ for differentiating $(t_x, t_y) \in \Omega$, *e.g.* nuclear *size* is highly relevant for classifying benign and malignant tumor as important nuclei in malignant are larger than important nuclei in benign.

$$s_{\max}(t_x, t_y) = \max_{c \in \mathcal{C}} S_{(t_x, t_y), c}$$
$$s_{\mathrm{avg}}(t_x, t_y) = \frac{1}{|C|} \sum_{c \in \mathcal{C}} S_{(t_x, t_y), c} \qquad (2)$$
$$s_{\mathrm{corr}}(t_x, t_y) = \rho(S_{(t_x, t_y), c=1,..,|\mathcal{C}|}, P_{(t_x, t_y), c=1,..,|\mathcal{C}|})$$

where $\rho$ denotes Pearson correlation. $s_{\max}$, $s_{avg} \in [0, \infty)$ show separation between unnormalized class-histograms; and $s_{corr} \in [-1, 1]$ shows agreement between $S$ and $P$. We build $S_{\max}$, $S_{\mathrm{avg}}$ and $S_{\mathrm{corr}}$ by computing Equation (2) $\forall (t_x, t_y) \in \Omega$. Metrics' complementary may lead to relevant *concepts* different to pathological understanding.

**Risk:** We *conceptually* introduce the notion of risk as a weight to indicate the cost of misclassifying a sample of class $t_x$, erroneously as class $t_y$ [62, 26]. Indeed, misclassifying a malignant tumor as a benign tumor is riskier than misclassifying it as an atypical tumor. Thus, we construct a risk vector $R \in \mathbb{R}^{\Omega}$. In this work, each entry in $R$ defines the symmetric risk of differentiating $t_x$ from $t_y$ measured as the number of class-hops needed to evolve from $t_x$ to $t_y$.

**Metrics:** Finally, we propose three quantitative metrics based on class separability to assess an explainer quality. The metrics are computed as the risk weighted sum of the statistics of separability scores, *i.e.*, (1) *maximum separability* $S_{\max, R} := S_{\max} \odot R$, (2) *average separability* $S_{\mathrm{avg}, R} := S_{\mathrm{avg}} \odot R$, (3) *correlated separability* $S_{\mathrm{corr}, R} := S_{\mathrm{corr}} \odot R$, where $\odot$ defines the Hadamard product. The first two metrics are pathologist-independent, and the third metric requires expert pathologists to impart the domain knowledge in the form of pathological prior $P$. Such prior can be defined individually by a pathologist or collectively by consensus of several pathologists, and it is independent of the algorithm generated explanations.

# 4. Results

This section describes the analysis of CG explainability for breast cancer subtyping. We evaluate three types of graph explainers and quantitatively analyze the explainer quality using the proposed class separability metrics.

## 4.1. Dataset

We experiment on BReAst Cancer Subtyping (BRACS), a large collection of breast tumor RoIs [46]. BRACS con-

---

**Algorithm 1:** Class separability computation.

**Input:** $\mathcal{D} = \{(D_i^t, \mathcal{I}_i^t)\}, t \in \mathcal{T}, i \in N_t$
**Parameter:** $\mathcal{T}, \mathcal{C}, \mathcal{A}_c, \mathcal{K}$
**Result:** $S \in \mathbb{R}^{\binom{|\mathcal{T}|}{2} \times |\mathcal{C}|}$

**for** $c$ in $\mathcal{C}$ **do** // go over concepts
  **for** $k$ in $\mathcal{K}$ **do** // go over nuclei thresh
    **for** $a$ in $\mathcal{A}_c$ **do** // go over attributes
      **for** $t$ in $\mathcal{T}$ **do** // go over classes
        var $\leftarrow D_i^t(a)[: k]$ // sorted $I_i^t$
        $H_t^{(k)}(a) \leftarrow$ histogram(var)
      **for** $(t_x, t_y)$ in $\binom{|\mathcal{T}|}{2}$ **do** // go over class pairs
        $s_a^{(k)}(t_x, t_y) \leftarrow d(H_{t_x}^{(k)}(a), H_{t_y}^{(k)}(a))$
    $s_c^{(k)}(t_x, t_y) \leftarrow \frac{1}{|\mathcal{A}_c|} \sum_{a \in \mathcal{A}_c} s_a^{(k)}(t_x, t_y)$
  $S_{(t_x, t_y), c} \leftarrow \mathrm{AUC}_{k \in \mathcal{K}}(s_c^{(k)}(t_x, t_y))$

---

sists of 4391 RoIs at $40\times$ resolution from 325 H&E stained breast carcinoma whole-slides. The RoIs are annotated by the consensus of three pathologists as, (1) Benign (B): normal, benign and usual ductal hyperplasia, (2) Atypical (A): flat epithelial atypia and atypical ductal hyperplasia, and (3) Malignant (M): ductal carcinoma *in situ* and invasive. The RoIs consist of an average #pixels=$3.9 \pm 4.3$ million, and average #nuclei=$1468 \pm 1642$, and are stain normalized using [60]. The train, validation, and test splits are created at the whole-slide level, including 3163, 602, and 626 RoIs.

## 4.2. Training

We conducted our experiments using PyTorch [45] and the Deep Graph Library (DGL) [65]. The GNN architecture for CG classification is presented in Section 3.3. The CG classifier was trained for 100 epochs using Adam optimizer [36], $10^{-3}$ learning rate and 16 batch size. The best CG-classifier achieved $74.2\%$ weighted F1-score on the test set for the three-class classification. Average time for processing a $1\mathrm{K} \times 1\mathrm{K}$ RoI on a NVIDIA P100 GPU is 2s for CG generation and 0.01s for GNN inference.

## 4.3. Qualitative assessment

Figure 4 presents explanations, *i.e.* nuclei importance maps, from four studied graph explainers. We observe that GRAPHGRAD-CAM and GRAPHGRAD-CAM++ produce similar importance maps. The GNNEXPLAINER generates almost binarized nuclei importances. Interestingly, the gradient and pruning-based techniques consistently highlight similar regions. Indeed, the approaches focus on relevant epithelial region and unfocus on stromal nuclei and lymphocytes outside the glands. Differently, GRAPHLRP produces less interpretable maps through high spatial localiza-
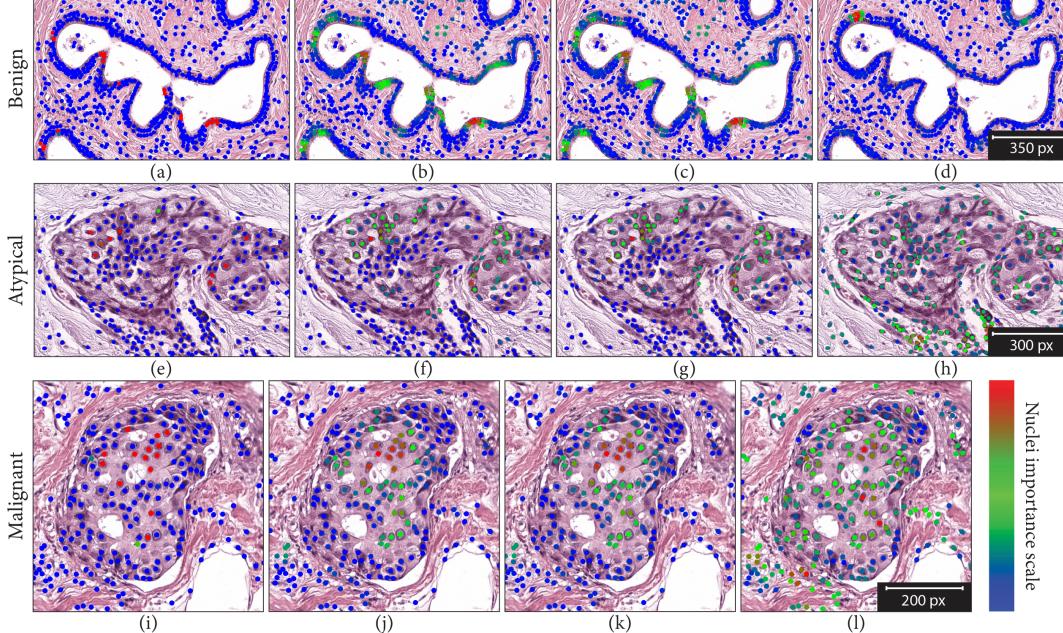
Figure 4. Qualitative results. The rows represent the cancer subtypes, *i.e.* Benign, Atypical and Malignant, and the columns represent the graph explainability techniques, *i.e.* GNNEXPLAINER, GRAPHGRAD-CAM, GRAPHGRAD-CAM++, and GRAPHLRP. Nuclei-level importance ranges from blue (the least important) to red (the most important).

tion (Figure 4(d)) or less spatial localization (Figure 4(h,l)). Qualitative visual assessment of Figure 4 conclude that, (1) *fidelity* preserving explainers result differently based on the underlying mechanism, (2) high *fidelity* does not guarantee straightforward pathologist-understandable explanations, (3) qualitative assessment cannot rigorously compare explainers' quality, and (4) large-scale tedious pathological evaluation is inevitable to rank the explainers.

### 4.4. Quantitative results

For cancer subtyping, relevant *concepts* are nuclear morphology and topology [49, 33, 44, 2]. Here, we focus on nuclear morphology, *i.e.* $\mathcal{C} = \{size, shape, shape variation, density, chromaticity\}$. Table 2 lists the *attributes* $\mathcal{A}_c, \forall c \in \mathcal{C}$. In our experiments, we select $\mathcal{K} = \{5, 10, ..., 50\}$ nuclei per CG. We further introduce a RANDOM explainer via *random* nuclei selection strategy per CG to assess a lower bound per quantitative metric. Table 1 presents the statistics of pair-wise class separability and aggregated separability w/ and w/o risk to assess the studied explainers quantitatively. Also, for each class pair $(t_x, t_y)$, we compute classification accuracy by using the CGs of type $t_x, t_y$.

Noticeably, GNNEXPLAINER achieves the best *maximum* and *average separability* for majority of pair-wise classes. GRAPHGRAD-CAM++ and GRAPHGRAD-CAM followed GNNEXPLAINER except for (B vs. A), where GRAPHLRP outperforms them. All explainers outperform RANDOM which conveys that the quality of the explainers' explanations are better than random. Notably,

GRAPHGRAD-CAM and GRAPHGRAD-CAM++ quantitatively perform very similarly, which is consistent with our qualitative analysis in Figure 4. Interestingly, a positive correlation is observed between pair-wise class accuracies and *average separability* for the explainers, *i.e.* better classification leads to better *concept* separability, and thus produces better explanations. Further, the observation does not hold for RANDOM generated explanations, which possesses undifferentiable *average concept* separability.

To obtain pathological prior to compute *correlation separability*, we consulted three pathologists to rank the *concepts* in terms of their relevance for discriminating each pair of classes. For instance, given an atypical RoI, we asked how important is nuclear *shape* to classify the RoI as *not* benign and *not* malignant. Acquired *concept* ranks for each class pair are *min-max* normalized to output prior matrix $P$. We observe that GNNEXPLAINER, GRAPHGRAD-CAM and GRAPHGRAD-CAM++ have positive *correlated separability* for (B vs. M), (A vs. M), and nearly zero values for (B vs. A). It shows that the explanations for (B vs. M) and (A vs. M) bear similar relevance of *concepts* as the pathologists, and focus on a different relevance of *concepts* for (B vs. A). GRAPHGRAD-CAM++ has the best overall agreement at the *concept*-level with the pathologists, followed by GRAPHGRAD-CAM and GNNEXPLAINER. RANDOM agrees significantly worse than the three explainers, and GRAPHLRP has the least agreement. Table 2 provides more insights by highlighting the per-*concept* metrics of GNNEXPLAINER. Nuclear *size* is the most relevant *con-*

| Tasks ($\Omega$) | | B vs. A | B vs. M | A vs. M | B vs. A vs. M | | | |
|---|---|---|---|---|---|---|---|---|
| Accuracy (in %) | | 77.19 | 90.29 | 80.42 | 74.92 | | | |
| Explainer | | Metric $\forall (t_x, t_y) \in \Omega$ ($\uparrow$) | | | Agg. Metric w/o Risk ($\uparrow$) | | Agg. Metric w/ Risk ($\uparrow$) | |
| GNNEXPLAINER | $s_{\max}(t_x,t_y)$ | **3.26** | **6.24** | **3.48** | $S_{\max}$ | **12.98** | $S_{\max,R}$ | **19.22** |
| GRAPHGRAD-CAM | | 1.24 | 4.41 | 3.36 | | 9.01 | | 13.42 |
| GRAPHGRAD-CAM++ | | 1.27 | <u>4.42</u> | <u>3.40</u> | | <u>9.09</u> | | <u>13.51</u> |
| GRAPHLRP | | <u>2.33</u> | 2.46 | 1.28 | | 6.07 | | 8.53 |
| RANDOM | | 1.02 | 1.26 | 1.11 | | 3.39 | | 4.65 |
| GNNEXPLAINER | $s_{\mathrm{avg}}(t_x,t_y)$ | **1.54** | **2.78** | 1.93 | $S_{\mathrm{avg}}$ | **6.25** | $S_{\mathrm{avg},R}$ | **9.03** |
| GRAPHGRAD-CAM | | 1.15 | 2.57 | <u>2.08</u> | | 5.80 | | 8.37 |
| GRAPHGRAD-CAM++ | | 1.18 | <u>2.58</u> | **2.09** | | <u>5.85</u> | | <u>8.43</u> |
| GRAPHLRP | | <u>1.38</u> | 1.59 | 1.47 | | 4.44 | | 6.03 |
| RANDOM | | 1.05 | 1.00 | 0.95 | | 3.00 | | 4.00 |
| GNNEXPLAINER | $s_{\mathrm{corr}}(t_x,t_y)$ | $-0.02$ | 0.36 | 0.38 | $S_{\mathrm{corr}}$ | 0.72 | $S_{\mathrm{corr},R}$ | 1.08 |
| GRAPHGRAD-CAM | | <u>$-0.01$</u> | <u>0.57</u> | <u>0.58</u> | | <u>1.14</u> | | <u>1.71</u> |
| GRAPHGRAD-CAM++ | | **$-0.01$** | **0.58** | **0.59** | | **1.16** | | **1.74** |
| GRAPHLRP | | $-0.15$ | $-0.49$ | $-0.23$ | | $-0.87$ | | $-1.36$ |
| RANDOM | | $-0.37$ | $-0.31$ | $-0.18$ | | $-0.86$ | | $-1.17$ |

Table 1. Quantitative assessment of graph explainers: GNNEXPLAINER, GRAPHGRAD-CAM, GRAPHGRAD-CAM++ and GRAPHLRP, using proposed *maximum, average*, and *correlated separability* metrics. Results are provided for each pair-wise breast subtyping tasks, and are aggregated w/o and w/ risk weighting, *i.e.* $S_{\max}$ and $S_{\max,R}$. The first and second best values are indicated in **bold** and <u>underline</u>.

| Concept (Attributes) / Tasks ($\Omega$) | B vs. A | B vs. M | A vs. M | w/o risk ($\uparrow$) | w/ risk ($\uparrow$) |
|---|---|---|---|---|---|
| Size (area) | **3.26** | **6.24** | **3.47** | **12.97** | **19.21** |
| Shape (perimeter, roughness, eccentricity, circularity) | 1.27 | 2.23 | 1.60 | 5.10 | 7.34 |
| Shape variation (shape factor) | 0.69 | 2.30 | 1.99 | 4.97 | 7.28 |
| Density (mean density, std density) | 1.01 | 0.80 | 0.52 | 2.33 | 3.14 |
| Chromaticity (GLCM contrast, homogeneity, ASM, entropy, variance) | <u>1.44</u> | <u>2.31</u> | <u>2.07</u> | <u>5.82</u> | <u>8.13</u> |
| *Average separability* ($\uparrow$) | 1.54 | 2.78 | 1.93 | 6.25 | 9.03 |

Table 2. Quantification of *concepts* for pair-wise and aggregated class separability in GNNEXPLAINER. The first and second best values are indicated in **bold** and <u>underline</u>. The per-*concept attributes* are presented in the first column.

*cept*, followed by *chromaticity* and *shape variation*. Comparatively nuclear *density* is the least relevant *concept*.

## 5. Conclusion

In this work, we presented an approach for explaining black-box DL solutions in computational pathology. We advocated for biological entity-based analysis instead of conventional pixel-wise analysis, thus providing an intuitive space for pathological understanding. We employed four graph explainability techniques, *i.e.* graph pruning (GNNEXPLAINER), gradient-based saliency (GRAPHGRAD-CAM, GRAPHGRAD-CAM++) and layerwise relevance propagation (GRAPHLRP), to explain "black-box" GNNs processing the entity graphs. We proposed a novel set of user-independent quantitative metrics expressing pathologically-understandable *concepts* to evaluate the graph explainers, which relaxes the exhaustive qualitative assessment by expert pathologists. Our analysis concludes that the explainer bearing the best class separability in terms of *concepts* is GNNEXPLAINER, followed by GRAPHGRAD-CAM++ and GRAPHGRAD-CAM. GRAPHLRP is the worst explainer in this category while outperforming a randomly created explanation. We observed that the explainer quality is directly proportional to the GNN's classification performance for a pair of classes. Furthermore, GRAPHGRAD-CAM++ produces explanations that best agrees with the pathologists in terms of *concept* relevance, and objectively highlights the relevant set of *concepts*. Considering the expansion of entity graph-based processing, such as radiology, computation biology, satellite and natural images, graph explainability and their quantitative evaluation is crucial. The proposed method encompassing domain-specific user-understandable terminologies can potentially be of great use in this direction. It is a meta-method that is applicable to other domains and tasks by incorporating relevant entities and corresponding *concepts*. For instance, with entity-graph nodes denoting car/body parts in Stanford Cars [38]/ Human poses [3], and expert knowledge available on car-model/ activity, our method can infer relevant parts by quantifying their agreement with experts.

# References

[1] M. Adnan, S. Kalra, and H.R. Tizhoosh. Representation Learning of Histopathology Images using Graph Neural Networks. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2020. 2

[2] K.H. Allison, M.H. Rendi, S. Peacock, T. Morgan, J.G. Elmore, and D.L. Weaver. Histologic Features associated with Diagnostic Agreement in Atypical Ductal Hyperplasia of the Breast: Illustrative Cases from the B-Path Study. *Histopathology*, 69(6):1028–1046, 2016. 7

[3] M. Andriluka, L. Pishchulin, P. Gehler, and S. Bernt. 2d human pose estimation: New benchmark and state of the art analysis. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3686–3693, 2014. 8

[4] A.B. Arrieta, N. Diaz-Rodriguez, J.D. Ser, A. Bennetot, S. Tabik, A. Barbado, S. Garcia, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila, and F. Herrera. Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges toward Responsible AI. *Information Fusion*, 58:82–115, 2020. 2

[5] B. Aygunes, S. Aksoya, R.G. Cinbis, K. Kosemehmetoglu, S. Onder, and A. Uner. Graph Convolutional Networks for Region of Interest Classification in Breast Histopathology. In *SPIE Medical Imaging: Digital Pathology*, 2020. 2

[6] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.R. Müller, and W. Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE*, 10(7), 2015. 1, 2, 4, 12

[7] F. Baldassarre and H. Azizpour. Explainability Techniques for Graph Convolutional Networks. *International Conference on Machine Learning Workshops*, 2019. 1

[8] K. Bera, K.A. Schalper, D.L. Rimm, V. Velcheti, and A. Madabhushi. Artificial intelligence in digital pathology — new tools for diagnosis and precision oncology. *Nature Reviews Clinical Oncology*, 16:703–715, 2019. 1

[9] A. Binder, M. Bockmayr, M. Hagele, S. Wienert, D. Heim, K. Hellweg, A. Stenzinger, L. Parlow, J. Budczies, B. Goeppert, D. Treue, M. Kotani, M. Ishii, M. Dietel, A. Hocke, C. Denkert, K.O. Mller, and F. Klauschen. Towards computational fluorescence microscopy: Machine learning-based integrated prediction of morphological and molecular tumor profiles. *arXiv:1805.11178*, 2018. 1

[10] K. Bruno, M.O. Andrea, P.M. Allen, M.N. Catherine, A.S. Matthew, T. Lorenzo, A.S. Arief, and H. Saeed. Looking under the hood Deep neural network visualization to interpret whole slide Image analysis outcomes for colorectal polyps. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017. 1, 2

[11] A. Chattopadhay, A. Sarkar, P. Howlader, and V.N. Balasubramanian. Grad-CAM++: Generalized gradient-based visual explanations for deep convolutional networks. In *IEEE Winter Conference on Applications of Computer Vision*, volume 2018-Janua, pages 839–847, 2018. 1, 2, 4, 13

[12] R.J. Chen, M.Y. Lu, J. Wang, D.F.K. Williamson, S.J. Rodig, N.I. Lindeman, and F. Mahmood. Pathomic Fusion: An Integrated Framework for Fusing Histopathology and Genomic Features for Cancer Diagnosis and Prognosis. *IEEE Transactions on Medical Imaging*, 2020. 2

[13] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering. In *Advances in Neural Information Processing Systems*, 2016. 3

[14] J. Deng, W. Dong, R. Socher, L.J. Li, K. Li, and F.F. Li. ImageNet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 3

[15] A. Dhurandhar, V. Iyengar, R. Luss, and K. Shanmugam. A Formal Framework to Characterize Interpretability of Procedures. In *International Conference on Machine Learning*, 2017. 2

[16] F. Doshi-Velez and B. Kim. Towards a rigorous science of interpretable machine learning. In *arXiv:1702.08608*, 2017. 2

[17] K. Francis and B.O. Palsson. Effective intercellular communication distances are determined by the relative time constants for cyto/chemokine secretion and diffusion. *Proceedings of the National Academy of Sciences*, 94(23):12258–12262, 1997. 3

[18] S. Gadiya, D. Anand, and A. Sethi. Histographs: Graphs in histopathology. *SPIE Medical Imaging: Digital Pathology*, 11320, 2020. 2

[19] J. Gilmer, S.S. Schoenholz, P.F. Riley, O. Vinyals, and G.E. Dahl. Neural Message Passing for Quantum Chemistry. *International Conference on Machine Learning*, 2017. 3

[20] S. Graham, Q.D. Vu, S.A. Raza, A. Azam, Y.H. Tsang, J.T. Kwak, and N. Rajpoot. Hover-Net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images. *Medical Image Analysis*, 58, 2019. 3

[21] M. Graziani, V. Andrearczyk, S. Marchand-Maillet, and H. Müller. Concept attribution: Explaining CNN decisions to physicians. In *Computers in Biology and Medicine*, volume 123, 2020. 1, 2

[22] C. Gunduz, B. Yener, and S.H. Gultekin. The cell graphs of cancer. *Bioinformatics*, 20(1):145–151, 2004. 1, 2

[23] M. Hägele, P. Seegerer, S. Lapuschkin, M. Bockmayr, W. Samek, F. Klauschen, K.R. Müller, and A. Binder. Resolving challenges in deep learning-based analyses of histopathological images using explanation methods. In *Nature Scientific Reports*, volume 10, 2020. 1, 2

[24] W.L. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs. In *Advances in neural information processing systems*, pages 1024–1034, 2017. 3

[25] R.M. Haralick, K. Shanmugam, and I. Dinstein. Texural features for image classification. *IEEE transaction on systems, man and cybernatics*, 3(6):610–621, 1973. 14

[26] H. He and Y. Ma. *Imbalanced Learning: Foundations, Algorithms, and Applications*. Wiley-IEEE Press, 1st edition, 2013. 6

[27] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 3

[28] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. In *Advances in Neural Information Processing Systems*, 2015. 13

[29] R.R. Hoffman, S.T. Mueller, G. Klein, and J. Litman. Metrics for explainable AI: Challenges and prospects. In *arXiv:1812.04608*, 2018. 2

[30] A. Holzinger, B. Malle, P. Kieseberg, P.M. Roth, H. Müller, R. Reihs, and K. Zatloukal. Towards the Augmented Pathologist: Challenges of Explainable-AI in Digital Pathology. In *arXiv:1712.06657*, 2017. 1, 2

[31] G. Jaume, P. Pati, A.F. Rodriguez, F. Florinda, G. Scognamiglio, A.M. Anniciello, J.P. Thiran, O. Goksel, and M. Gabrani. Towards Explainable Graph Representations in Digital Pathology. In *International Conference on Machine Learning Workshops*, 2020. 2, 4, 13

[32] S. Javed, A. Mahmood, M.M. Fraz, N.A. Koohbanani, K. Benes, Y.W. Tsang, K. Hewitt, D. Epstein, D. Snead, and N. Rajpoot. Cellular community detection for tissue phenotyping in colorectal cancer histology images. *Medical Image Analysis*, 63, 2020. 2

[33] A. Kashyap, M. Jain, S. Shukla, and M. Andley. Role of Nuclear Morphometry in Breast Cancer and its Correlation with Cytomorphological Grading of Breast Cancer: A Study of 64 Cases. *Journal of Cytology*, 35(1):41–45, 2018. 7

[34] B. Kim, M. Wattenberg, J. Gilmer, C. Cai, J. Wexler, F. Viegas, and R. Sayres. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). *International Conference on Machine Learning*, page 2673–2682, 2018. 1

[35] P. Kindermans, K. T. Schutt, M. Alber, K.R. Muller, and S. Dahne. PatternNet and PatternLRP - improving the interpretability of neural networks. *arXiv:1705.05598*, 2015. 1

[36] D.P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. In *International Conference on Learning Representations*, 2015. 6

[37] T. Kipf and M. Welling. Semi supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, pages 1–14, 2017. 1, 3

[38] J. Krause, M. Stark, J. Deng, and L. Fei-Fei. 3d object representations for fine-grained categorization. In *IEEE International Conference on Computer Vision Workshops*, pages 554–561, 2013. 8

[39] N. Kumar, R. Verma, S. Sharma, S. Bhargava, A. Vahadane, and A. Sethi. A dataset and a technique for generalized nuclear segmentation for computational pathology. *IEEE Transactions on Medical Imaging*, 36(7):1550–1560, 2017. 3

[40] M.Y. Lu, D.F.K. Williamson, T.Y. Chen, R.J. Chen, M. Barbieri, and F. Mahmood. Data Efficient and Weakly Supervised Computational Pathology on Whole Slide Images. *arXiv:2004.09666*, 2020. 1, 2

[41] S. Mohseni, N. Zarei, and E.D. Ragan. A multidisciplinary survey and framework for design and evaluation of explainable AI systems. In *arXiv:1811.11839*, 2018. 2

[42] G. Montavon, S. Bach, A. Binder, W. Samek, and K.R. Muller. Explaining NonLinear Classification Decisions with Deep Taylor Decomposition. *Pattern Recognition*, 65:211–222, 2015. 1, 2, 12

[43] A. Nguyen and M. Rodriguez Martinez. On quantitative aspects of model interpretability. In *arXiv:2007.07584*, 2020. 2

[44] L. Nguyen, A.B. Tosun, J.L. Fine, D.L. Taylor, and S.C. Chennubhotla. Achitectural patterns for differential diagnosis of proliferative breast lesions from histopathological images. In *IEEE International Symposium on Biomedical Imaging*, pages 152–155, 2017. 7

[45] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019. 6

[46] P. Pati, G. Jaume, L.A. Fernandes, A. Foncubierta, F. Feroce, A.M. Anniciello, G. Scognamiglio, N. Brancati, D. Riccio, M.D. Bonito, G.D. Pietro, G. Botti, O. Goksel, J.P. Thiran, M. Frucci, and M. Gabrani. HACT-Net: A Hierarchical Cell-to-Tissue Graph Neural Network for Histopathological Image Classification. In *Uncertainty for Safe Utilization of Machine Learning in Medical Imaging, and Graphs in Biomedical Image Analysis*, pages 208–219, 2020. 1, 2, 3, 6

[47] P. Pati, G. Jaume, A. Foncubierta, F. Feroce, A.M. Anniciello, G. Scognamiglio, N. Brancati, M. Fiche, E. Dubruc, D. Riccio, M.D. Bonito, G.D. Pietro, G. Botti, J.P. Thiran, M. Frucci, O. Goksel, and M. Gabrani. Hierarchical Graph Representations in Digital Pathology. In *Uncertainty for Safe Utilization of Machine Learning in Medical Imaging, and Graphs in Biomedical Image Analysis*, 2021. 2

[48] P.E. Pope, S. Kolouri, M. Rostami, C.E. Martin, and H. Hoffmann. Explainability methods for graph convolutional neural networks. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 10764–10773, 2019. 1, 2, 4, 12

[49] N. Rajbongshi, K. Bora, D.C. Nath, A.K. Das, and L.B. Mahanta. Analysis of Morphological Features of Benign and Malignant Breast Cell Extracted From FNAC Microscopic Image Using the Pearsonian System of Curves. *Journal of Cytology*, 35(2):99–104, 2018. 7

[50] M.T. Ribeiro, S. Singh, and C. Guestrin. Why should i you? Explaining the predictions of any classifier. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144, 2016. 2

[51] W. Samek, A. Binder, G. Montavon, S. Lapuschkin, and K.R. Muller. Evaluating the visualization of what a deep neural network has learned. In *IEEE Transactions on Neural Networks and Learning Systems*, volume 28, pages 2660–2673, 2017. 2

[52] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, pages 61–80, 2009. 3

[53] R. Schwarzenberg, M. Huebner, D. Harbecke, C. Alt, and L. Hennig. Layerwise relevance visualization in convolutional text graph classifiers. *EMNLP Workshop*, pages 58–62, 2019. 2, 4, 12

[54] R.R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, and D. Batra. Grad-CAM : Visual Explanations from Deep Networks. In *International Conference on Computer Vision*, pages 618–626, 2017. 1, 2, 4, 12

[55] A. Serag, A. Ion-Margineanu, H. Qureshi, R. McMillan, M.S. Martin, J. Diamond, P. O'Reilly, and P. Hamilton. Translational AI and Deep Learning in Diagnostic Pathology. *Frontiers in Medicine*, 2019. 1

[56] H. Sharma, N. Zerbe, C. Boger, S. Wienert, O. Hellwich, and P. Hufnagl. A Comparative Study of Cell Nuclei Attributed Relational Graphs for Knowledge Description and Categorization in Histopathological Gastric Cancer Whole Slide Images. *IEEE Symposium on Computer-Based Medical Systems*, pages 61–66, 2017. 2

[57] H. Sharma, N. Zerbe, D. Heim, S. Wienert, S. Lohmann, O. Hellwich, and P. Hufnagl. Cell nuclei attributed relational graphs for efficient representation and classification of gastric cancer in digital histopathology. *SPIE Medical Imaging: Digital Pathology*, 9791, 2016. 2

[58] H. Sharma, N. Zerbe, S. Lohmann, K. Kayser, O. Hellwich, and P. Hufnagl. A review of graph-based methods for image analysis in digital histopathology. *Diagnostic Pathology*, 2015. 2

[59] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. *arXiv:1312.6034*, 2013. 1

[60] M. Stanisavljevic, A. Anghel, N. Papandreou, S. Andani, P. Pati, J.H. Rüschoff, 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 Workshops*, pages 424–436. 2019. 6

[61] M. Sureka, A. Patil, D. Anand, and A. Sethi. Visualization for histopathology images using graph convolutional neural networks. In *arXiv:2006.09464*, 2020. 2

[62] N. Thai-Nghe, Z. Gantner, and L. Schmidt-Thieme. Cost-sensitive learning methods for imbalanced data. In *The 2010 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2010. 6

[63] H.R. Tizhoosh and L. Pantanowitz. Artificial Intelligence and Digital Pathology: Challenges and Opportunities. *Journal of Pathology Informatics*, 38(9), 2018. 1, 2

[64] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph Attention Networks. *International Conference on Learning Representations*, 2018. 3

[65] M. Wang, L. Yu, D. Zheng, Q. Gan, Y. Gai, Z. Ye, M. Li, J. Zhou, Q. Huang, C. Ma, Z. Huang, Q. Guo, H. Zhang, H. Lin, J. Zhao, J. Li, A.J. Smola, and Z. Zhang. Deep Graph Library: Towards Efficient and Scalable Deep Learning on Graphs. *CoRR*, 2019. 6

[66] K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How Powerful are Graph Neural Networks? In *International Conference on Learning Representations*, 2018. 1, 4

[67] R. Ying, D. Bourgeois, J. You, M. Zitnik, and J. Leskovec. GNNExplainer: Generating Explanations for Graph Neural Networks. In *Advances in Neural Information Processing Systems*, 2019. 1, 2, 4, 13

[68] R. Ying, J. You, C. Morris, X. Ren, W.L. Hamilton, and J. Leskovec. Hierarchical Graph Representation Learning with Differentiable Pooling, 2018. 3

[69] J. Yosinski, J. Clune, A. M. Nguyen, T. J. Fuchs, and H. Lipson. Understanding neural networks through deep visualization. *International Conference on Machine Learning Workshops*, 2015. 1

[70] M.D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. *European Conference on Computer Vision*, 2014. 1

[71] Z. Zhang, P. Chen, M. McGough, F. Xing, and C. Wang. Pathologist-level interpretable whole-slide cancer diagnosis with deep learning. *Nature Machine Intelligence*, pages 236–245, 2019. 2

[72] Y. Zhao, F. Yang, Y. Fang, H. Liu, N. Zhou, J. Zhang, J. Sun, S. Yang, B. Menze, X. Fan, and J. Yao. Predicting lymph node metastasis using histopathological images based on multiple instance learning with deep graph convolution. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2

[73] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning Deep Features for Discriminative Localization. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2921–2929, 2016. 12

[74] Y. Zhou, S. Graham, N.A. Koohbanani, M. Shaban, P.A. Heng, and N. Rajpoot. CGC-Net: Cell Graph Convolutional Network for Grading of Colorectal Cancer Histology Images. *International Conference on Computer Vision Workshops*, 2019. 1, 2

[75] L. M. Zintgraf, T. S. Cohen, T. Adel, and M. Welling. Visualizing deep neural network decisions: Prediction difference analysis. *International Conference on Learning Representations*, 2017. 1

## Post-hoc explainers

In this section, we present the details of the considered graph explainability techniques (explainers) in this work: GRAPHLRP, GRAPHGRAD-CAM, GRAPHGRAD-CAM++, GNNEXPLAINER, and RANDOM.

### Notation

We define an attributed undirected entity graph $G := (V, E, H)$ as a set of nodes $V$, edges $E$, and node attributes $H \in \mathbb{R}^{|V| \times d}$. $d$ denotes the number of attributes per node, and $|.|$ denotes set cardinality. We denote an edge between nodes $u$ and $v$ as $e_{uv} \in E$. The graph topology is defined by a symmetric graph adjacency, $A \in \mathbb{R}^{|V| \times |V|}$, where $A_{uv} = 1$ if $e_{uv} \in E$. $H_{n,k}$ expresses the $k$-th attribute of the $n$-th node. The forward prediction of a cell-graph $G_{CG}$ is denoted as, $y = \mathcal{M}(G_{CG})$, where $\mathcal{M}$ is a pre-trained GNN, and $y \in \mathbb{R}^{|\mathcal{T}|}$ are the output logits. Notation $y(t)$, $t \in \mathcal{T}$ denotes the output logit of the $t$-th class. We refer to the logit of the predicted class as $y_{\max} = \max_{t \in \mathcal{T}} y(t)$, and the predicted class as $t_{\max} = \arg\max_{t \in \mathcal{T}} y(t)$.

### Layerwise relevance propagation: GRAPHLRP

Layerwise Relevance Propagation (LRP) [6] is a feature attribution based post-hoc explainer. LRP explains an output logit by determining the individual contribution of each input element to the logit value. An output logit, defined as the output *relevance* for a given class, is layerwise back-propagated until the input to compute the positive or negative impact of the input elements on the output logit. LRP, initially proposed for fully connected layers (LRP-FC), works as follows. Given a pre-trained fully connected layer $W \in \mathbb{R}^{z_1 \times z_2}$ between layer 1 and layer 2, where $z_1$ and $z_2$ are the number of neurons in layer 1 and layer 2, respectively, we compute the contributions of a neuron $i$, $i \in \{1, ..., z_1\}$ using the propagation rules in [42]. In this work, we are interested in identifying the input elements *positively* contributing to the prediction. To this end, we use the $z^+$ propagation rule that back-propagates the *positive* neuron contribution from layer 2 to layer 1 as:

$$R_i = \sum_{j}^{z_2} \frac{f_i |w_{ij}|}{\sum_{k}^{z_1} f_k |w_{kj}|} R_j \qquad \text{(LRP-FC)}$$

where $|w_{ij}|$ is the absolute value of the weight between $i$-th and $j$-th neuron in layer 1 and 2, respectively. $f_i$ denotes the activation of the $i$-th neuron in layer $l$.

The extension from LRP-FC to LRP for graph isomorphism network (GIN) layers (GRAPHLRP) is achieved by following the observations in [53]. First, the *aggregate step* in GNN corresponds to projecting the graph's adjacency matrix on the node attribute space. For simplicity, assuming a 1-layer MLP as an update function, the GIN layer with

*mean* aggregator can be re-written in its global form as:

$$H^{(l+1)} = \sigma\Big(W^{(l)}(I + \tilde{A})H^{(l)}\Big) \qquad (3)$$

where $\tilde{A}$ is the degree-normalized graph adjacency matrix, *i.e.* $\tilde{A}_{ij} = \frac{1}{|\mathcal{N}(i)|} A_{ij}$. $\sigma$ is the ReLU activation function. Second, this representation allows us to treat the term $(I + \tilde{A})$ as a regular, fully connected layer. We can then apply the $z^+$ propagation rule with weights $w_{ij}$ defined as:

$$w_{ij} = 1 \quad \text{if } i = j \qquad (4)$$

$$w_{ij} = \frac{1}{|\mathcal{N}(i)|} \quad \text{if } e_{ij} \in E \qquad (5)$$

$$w_{ij} = 0 \quad \text{otherwise} \qquad (6)$$

LRP outputs an importance score for each node $i$ in the input graph.

### Saliency-based: GRAPHGRAD-CAM

Grad-CAM [54] is a feature attribution post-hoc explainer that identifies salient regions of the input driving the neural network prediction. It assigns importance to each element of the input to produce Class Activation Map [73]. While originally developed for explaining CNNs operating on images, GRAD-CAM can be extended to GNNs operating on graphs [48].

GRAPHGRAD-CAM processes in two steps. First, it assigns an importance score to each channel of a graph convolutional layer. The importance of channel $k$ in layer $l$ is computed by looking at the gradient of the predicted output logit $y_{\max}$ w.r.t. the node attributes at layer $l$ of the GNN. Formally it is expressed as:

$$w_k^{(l)} = \frac{1}{|V|} \sum_{n=1}^{|V|} \frac{\partial y_{\max}}{\partial H_{n,k}^{(l)}} \qquad (7)$$

In the second step, a node-wise importance score is computed using the forward node feature activations $H^{(l)}$ as:

$$L(l,v) = \text{ReLU}\Big( \sum_{k}^{d^{(l)}} w_k^{(l)} H_{n,k}^{(l)} \Big) \quad \text{(GRAPHGRAD-CAM)}$$

where $L(l, v)$ denotes the importance of node $v \in V$ in layer $l$, and $d^{(l)}$ denotes the number of node attributes at layer $l$. Since we are only interested in the positive node contributions, *i.e.* nodes that positively influence the class prediction, we apply a ReLU activation to the node importances. Following prior work [48], we take the average node importance scores obtained over all the GNN layers $l \in \{1, ..., L\}$ to obtain smoother node importance scores.

**Saliency-based:** GRAPHGRAD-CAM++

GRAPHGRAD-CAM++ extends GRAD-CAM++ [11] to graph structured data. It improves the node importance localization by introducing node-wise contributions to channel importance scoring in Equation 7. Specifically, the modification is presented as,

$$w_k^{(l)} = \frac{1}{|V|} \sum_{n=1}^{|V|} \alpha_{n,k}^{(l)} \frac{\partial y_{max}}{\partial H_{n,k}^{(l)}} \tag{8}$$

where $\alpha_{n,k}^{(l)}$ are node-wise weights expressed for each attribute $k$ at layer $l$. The derivation of a closed-form solution for $\alpha_{n,k}^{(l)}$ is analogous to the derivation in [11], where the size of graph, *i.e.* number of nodes, replaces the spatial dimensions of a channel as:

$$\alpha_{n,k}^{(l)} = \frac{\frac{\partial^2 y_{max}}{(\partial H_{n,k}^{(l)})^2}}{2 \frac{\partial^2 y_{max}}{(\partial H_{n,k}^{(l)})^2} + \sum_{n=1}^{|V|} H_{n,k}^{(l)} \left( \frac{\partial^3 y_{max}}{(\partial H_{n,k}^{(l)})^3} \right)} \tag{9}$$

The subsequent node importance computation in GRAPHGRAD-CAM++ is same as GRAPHGRAD-CAM.

**Graph pruning:** GNNEXPLAINER

The GNNEXPLAINER [67, 31] is a graph pruning based post-hoc explainer for explaining GNNs. GNNEXPLAINER is model-agnostic, *i.e.* it can be used with any flavor of GNN. Intuitively, GNNEXPLAINER tries to find the minimum sub-graph $G_s \subset G$ such that the model prediction $y = \mathcal{M}(G)$ is retained. The inferred sub-graph $G_s$ is then regarded as the *explanation* for $G$. This approach can be seen as a feature attribution method with *binarized* node importance scores, *i.e.* a node $v \in V$ has importance one if $v \in V_s$, and zero otherwise. Exhaustively searching $G_s$ in the space created by nodes $V$ and edges $E$ is infeasible due to the combinatorial nature of the task. Instead, GNNEXPLAINER formulates the task as an optimization problem that learns a mask to activate or deactivate parts of the graph. The initial formulation by [67], developed for explaining node classification tasks, learns a mask over the edges, *i.e.* over the adjacency matrix. Instead, we follow the prior work in [31] to learn a mask over the nodes. Indeed, as we are concerned with classifying $G$, the optimal explanation $G_s$ can be a disconnected graph. Furthermore, in cell graphs, the nodes representing biological entities are more intuitive and substantial for disease diagnosis than edges, that are heuristically-defined.

Formally, we seek to learn a mask $M_V$ such that the induced masked sub-graph $G_s$, (1) is as small as possible, (2) outputs a binary node importance, and (3) provides the same prediction as the original graph. These constraints can

be modeled by considering a loss function as:

$$\mathcal{L} = \mathcal{L}_{KD}(\hat{y}, y^{(m)}) + \alpha_{M_V} \sum_{i}^{|V|} \sigma(M_{V_i}^{(m)}) + \alpha_{\mathcal{H}} \mathcal{H}^e(\sigma(M_V^{(m)})) \tag{10}$$

where, $m$ is the optimization step and $\sigma$ is the sigmoid activation function. The first term is a knowledge-distillation loss $\mathcal{L}_{KD}$ between $\hat{y} = \mathcal{M}(G)$ and $y^{(m)} = \mathcal{M}(G_s)$ ensuring that $y^{(m)} \approx \hat{y}$. The second term aims to minimize the size of the mask $M_V$. The third term binarizes the mask by minimizing the element-wise entropy $\mathcal{H}^e$ of $M_V$. Following previous work [28], $\mathcal{L}_{KD}$ is built as a combination of distillation and cross-entropy loss,

$$\mathcal{L}_{KD} = \lambda \mathcal{L}_{CE} + (1 - \lambda)\mathcal{L}_{dist} \text{ where } \lambda = \frac{\mathcal{H}^e(y^{(m)})}{\mathcal{H}^e(\hat{y})} \tag{11}$$

where $\mathcal{L}_{CE}$ is the regular cross-entropy loss and $\mathcal{L}_{dist}$ is the distillation loss. When the element-wise entropy $\mathcal{H}^e(y^{(m)})$ increases, the term $\mathcal{L}_{CE}$ gets larger and reduces the probability of changing the prediction. Each term in Equation 10 is empirically weighed such that their contributions to $\mathcal{L}$ are comparable. We set $\alpha_{M_V} = 0.005$ and $\alpha_{\mathcal{H}} = 0.1$. We learn $M_V$ using Adam optimizer with a learning rate of 0.01. $\mathcal{L}$ is optimized for 1000 steps with an early stopping mechanism, which triggers if the class prediction using $G_s$ is changed. Therefore, $G_s$ and $G$ always predict the same class, *i.e.* $t_{max}^{(m)} = \hat{t}_{max} \, \forall m$.

**Random selection:** RANDOM

The RANDOM baseline is implemented using a *random* nuclei selection. The number of selected nuclei per RoI is given by the threshold value $k \in \mathcal{K}$.

## BRACS dataset

In this paper, the BRACS dataset is used to analyze CG explainability for breast cancer subtyping. The pixel-level and entity-level statistics of the dataset are presented in Table 3. Training, validation, and test splits are created at the whole-slide level for conducting the experiments. The details of the class-wise distribution of images in each split are presented in Table 3.

## Concepts and Attributes

In this paper, we focus on pathologically-understandable nuclear *concepts* $\mathcal{C}$ pertaining to nuclear morphology for breast cancer subtyping. To quantify each $c \in \mathcal{C}$, we use several measurable *attributes* $\mathcal{A}_c$. Table 4 presents the list of *concepts* and corresponding *attributes* used to perform the proposed quantitative analysis in this work. Also, Table 4 includes the class-wise expected criteria for each *concept*.

| | Metric | Benign | Atypical | Malignant | Total |
|---|---|---|---|---|---|
| Image | Number of images | 1741 | 1351 | 1299 | 4391 |
| | Number of pixels (in million) | 3.9±3.54 | 1.62±1.48 | 6.35±5.2 | 3.9±4.3 |
| | Max/Min pixel ratio | 180.1 | 75.3 | 128.6 | 235.6 |
| CG | Number of nodes | 1331±1134 | 635±510 | 2521±1934 | 1468±1642 |
| | Number of edges | 4674±4131 | 2309±2110 | 8591±7646 | 5102±6089 |
| | Max/Min node ratio | 312.5 | 416.7 | 312.5 | 434.8 |
| Image split | Train | 1231 | 1008 | 928 | 3163 |
| | Validation | 261 | 162 | 179 | 602 |
| | Test | 249 | 185 | 192 | 626 |

Table 3. Statistics of BRACS dataset.

| Concept ($\mathcal{C}$) | Attribute ($\mathcal{A}$) | Computation | Benign | Atypical | Malignant |
|---|---|---|---|---|---|
| Size | Area | $A(x)$ | Small | Small-Medium | Medium-Large |
| Shape | Perimeter | $P(x)$ | Smooth | Mild irregular | Irregular |
| | Roughness | $\frac{P_{\mathrm{ConvHull}}(x)}{P(x)}$ | | | |
| | Eccentricity | $\frac{a_{\mathrm{minor}}(x)}{a_{\mathrm{major}}(x)}$ | | | |
| | Circularity | $\frac{4\pi A(x))}{P(x)^2}$ | | | |
| Shape variation | Shape factor | $\frac{4\pi A(x)}{P^2_{\mathrm{ConvHull}}}$ | Monomorphic | Monomorphic | Pleomorphic |
| Spacing | Mean spacing | $\mathrm{mean}(d_y\|y \in \mathrm{kNN(x)})$ | Evenly crowded | Evenly spaced | Variable |
| | Std spacing | $\mathrm{std}(d_y\|y \in \mathrm{kNN(x)})$ | | | |
| Chromatin | GLCM dissimilarity | $\sum_i \sum_j \|i-j\|p(i,j)$ | Light euchromatic | Hyperchromatic | Vesicular |
| | GLCM contrast | $\sum_i \sum_j (i-j)^2 p(i,j)$ | | | |
| | GLCM homogenity | $\sum_i \sum_j \frac{p(i,j)}{1+(i-j)^2}$ | | | |
| | GLCM ASM | $\sum_i \sum_j p(i,j)^2$ | | | |
| | GLCM entropy | $-\sum_i \sum_j p(i,j)\log(p(i,j))$ | | | |
| | GLCM variance | $\sum_i \sum_j (i-\mu_i)^2 p(i,j)$ with $\mu_i = \sum_i \sum_j ip(i,j)$ | | | |

Table 4. Pathologically-understandable nuclear *concepts*, corresponding measurable *attributes*, and computations are shown in Columns 1, 2, 3, respectively. The expected *concept* behavior for three breast cancer subtypes is shown in Columns 4, 5, 6, respectively.

The *attributes* of the nuclei in a RoI are computed as presented in Table 4. It uses the RoI and corresponding nuclei segmentation map, denoted as $I_{\mathrm{seg}}$. Area of a nucleus $x$, denoted as $A(x)$, is defined as the number of pixels belonging to $x$ in $I_{\mathrm{seg}}$. $P(x)$, the perimeter of $x$, is measured as the contour length of $x$ in $I_{\mathrm{seg}}$. $P_{\mathrm{ConvHull}}(x)$, the convex hull perimeter of $x$, is defined as the contour length of convex hull induced by $x$ in $I_{\mathrm{seg}}$. The major and minor axis of $x$, noted as $a_{\mathrm{major}}(x)$ and $a_{\mathrm{minor}}(x)$, are the longest diameter of $x$ and the longest line segment perpendicular to $a_{\mathrm{major}}(x)$, respectively. The chromatin *attributes* are computed from the normalized gray level co-occurrence matrix (GLCM) [25], which captures the probability distribution of co-occurring gray values in $x$.

## Quantitative assessment

In this section, we analyze two key components of the proposed quantitative metrics: the histogram construction and class separability scores for threshold set $\mathcal{K}$. Furthermore, we relate the analysis to the class-wise expected criteria for each *concept* presented in Table 4.

## Histogram analysis

Histogram construction is a key component in the proposed quantitative metrics. Figure 5 presents per-class histograms for each explainer and the best *attribute* per *concept*. We set the importance threshold to $k = 25$, *i.e.* for each RoI, we select 25 nuclei with the highest node importance. The best *attribute* for a *concept* is the one with the highest average pair-wise class separability.

The row-wise observation exhibits that GNNEXPLAINER and GRAPHLRP provide, respectively, the maximum and the minimum pair-wise class separability. The histograms for a *concept* and for an explainer can be analyzed to assess the agreement between the selected important nuclei *concept*, and the expected *concept* behavior as presented in Table 4, for all the classes. For instance, nuclear *area* is expected to be higher for malignant RoIs than benign ones. The *area* histograms for GNNEXPLAINER, GRAPHGRAD-CAM and GRAPHGRAD-CAM++ indicate that the important nuclei set in malignant RoIs includes nuclei with higher area compared to benign RoIs. Similarly, the important nuclei in malignant RoIs are expected to be vesicular, *i.e.* high texture entropy, compared to light euchromatic, *i.e.* moderate texture entropy, in benign RoIs. The *chromaticity* histograms for GNNEXPLAINER, GRAPHGRAD-CAM and GRAPHGRAD-CAM++ display this behavior. Additionally, the histogram analysis can reveal the important *concepts* and important *attributes*. For instance, nuclear *density* proves to be the least important *concept* for differentiating the classes.

## Separability score for threshold set $\mathcal{K}$

Multiple importance thresholds $\mathcal{K}$ are required to address the varying notion of important nuclei across different cell graphs and different explainers. Figure 6 presents the behavior of pair-wise class separability for using various $k \in \mathcal{K} = \{5, 10, ..., 50\}$. For simplicity, we present the behavior for the best *attribute* per *concept*. In general, the pair-wise class separability is observed to decrease with decreasing $k$. Intuitively, decreasing $k$ results in including more unimportant nuclei into the evaluation, thereby gradually decreasing the class separability.

The degree of agreement between the difference in the expected behavior per *concept* and the pair-wise class separability in Figure 6, for all pair-wise classifications and various $k \in \mathcal{K}$ can be used to assess the explainer's quality. For instance, according to Table 4, the difference in the expected nuclear *size* can be considered as benign–atypical $<$ benign–malignant, and atypical–malignant $<$ benign–malignant. GNNEXPLAINER, GRAPHGRAD-CAM and GRAPHGRAD-CAM++ display these behaviors $\forall k \in \mathcal{K}$. GNNEXPLAINER provides the highest class separability in each pair-wise classification, thus proving to be the best ex-

plainer pertaining to *size concept*. Detailed inspection of Figure 6 shows that all the differences in the expected behavior, per *concept* for all pair-wise classifications, is inline with the *concept*-wise expected behavior in Table 4, $\forall c \in \mathcal{C}$ and $\forall k \in \mathcal{K}$. Overall, GNNEXPLAINER is seen to be the best explainer as it agrees to the majority of the expected differences $\forall c \in \mathcal{C}$ for all pair-wise classifications, while providing high-class separability. Furthermore, *size* proves to be the most important *concept* that provides the maximum class separability across all pair-wise classifications.

## Qualitative assessment

Figure 7 and Figure 8 present CG explanations produced by GNNEXPLAINER, GRAPHGRAD-CAM, GRAPHGRAD-CAM++ and GRAPHLRP for RoIs across benign, atypical and malignant breast tumors. It can be observed that GNNEXPLAINER learns to binarize the explanations, thereby producing the most compact explanations by retaining the most important nuclei set of nuclei with high importance. However, GRAPHGRAD-CAM and GRAPHGRAD-CAM++ produce explanations with more distributed nuclei importance than GNNEXPLAINER. GRAPHLRP produces the largest explanations by retaining most of the nuclei in the CGs.
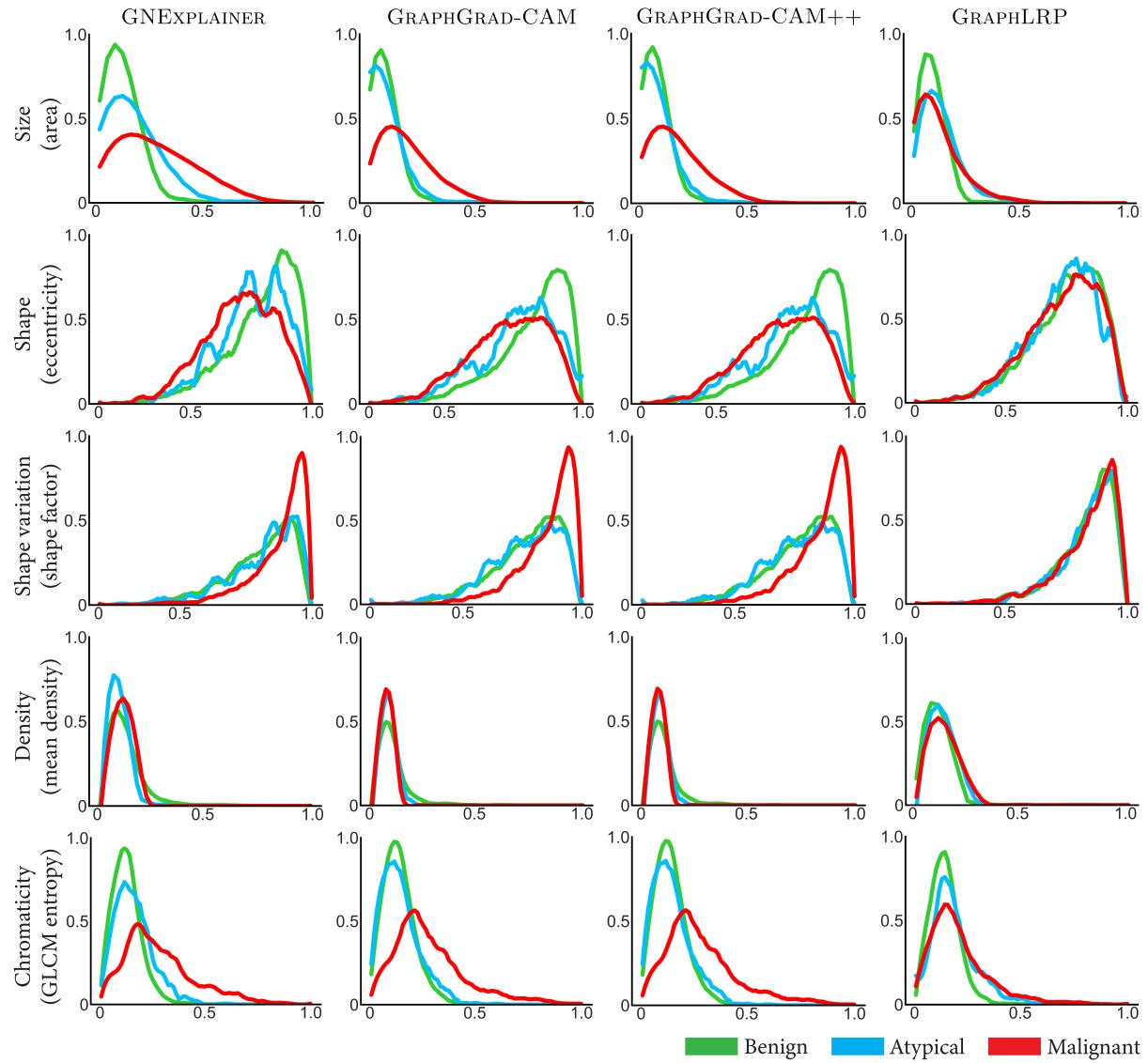
Figure 5. Per-class histograms for different *concepts* across different graph explainers. For simplicity, histograms are presented for the best *attribute* per *concept* at fixed importance threshold $k = 25$.
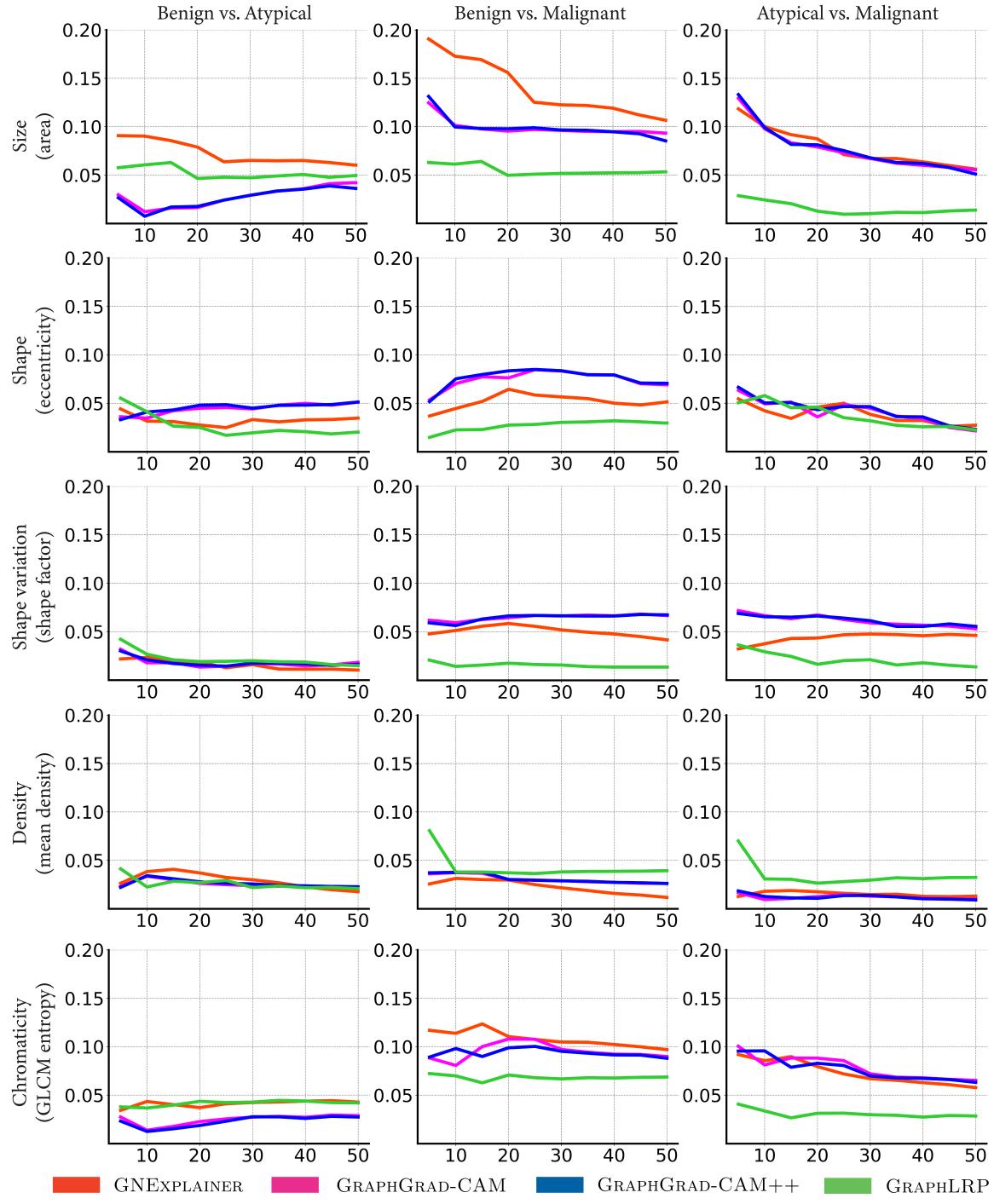
Figure 6. Visualizing the variation of pair-wise class separability score (Y-axis) w.r.t. various nuclei importance thresholds in $\mathcal{K}$ (X-axis). The analysis is provided for different graph explainers, and for the best *attribute* per *concept*.
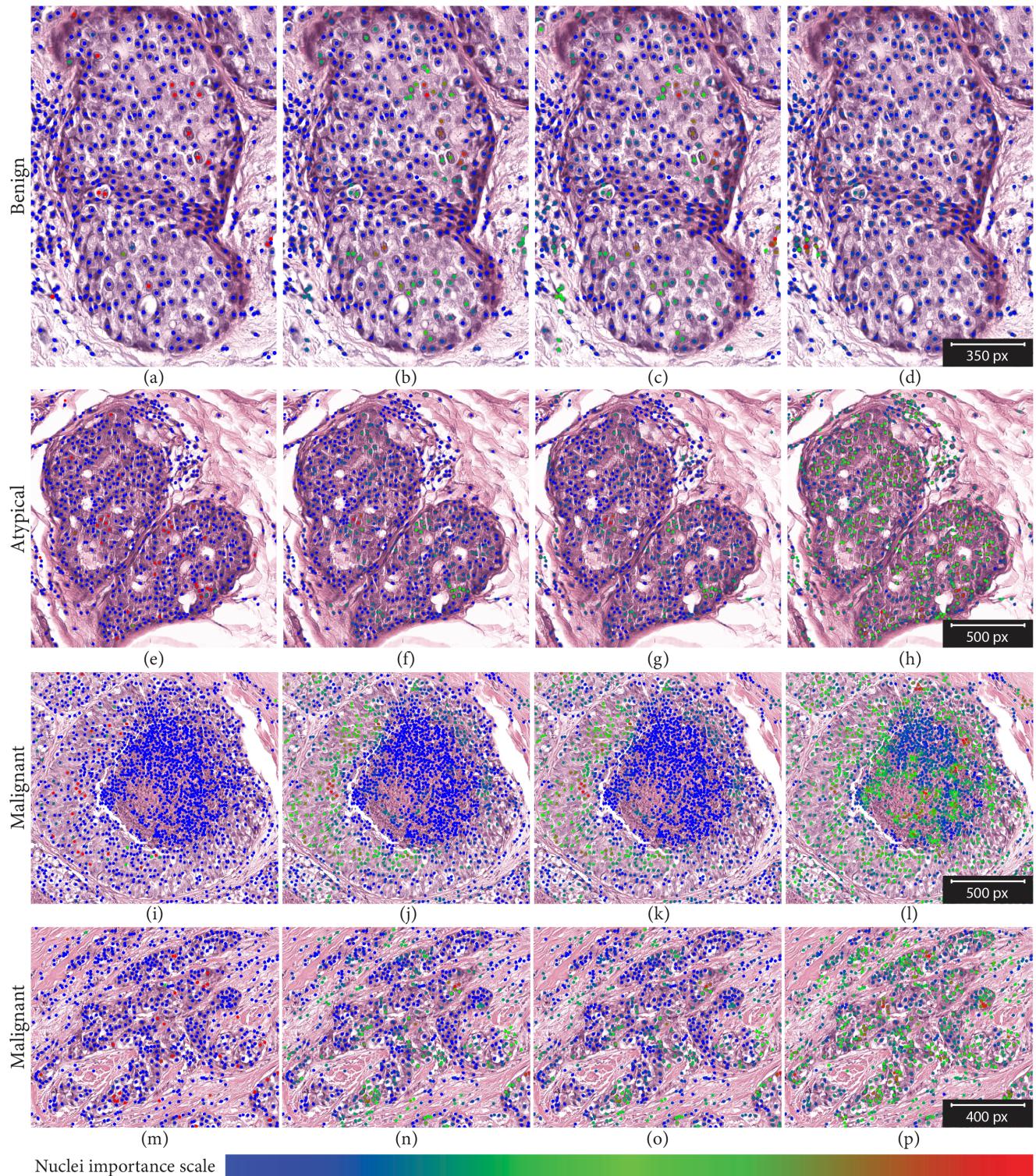
Figure 7. Qualitative results. The rows represent breast cancer subtypes, and columns represent graph explainers, *i.e.* GNNEXPLAINER, GRAPHGRAD-CAM, GRAPHGRAD-CAM++, and GRAPHLRP. Nuclei level importance ranges from blue (the least important) to red (the highest important).
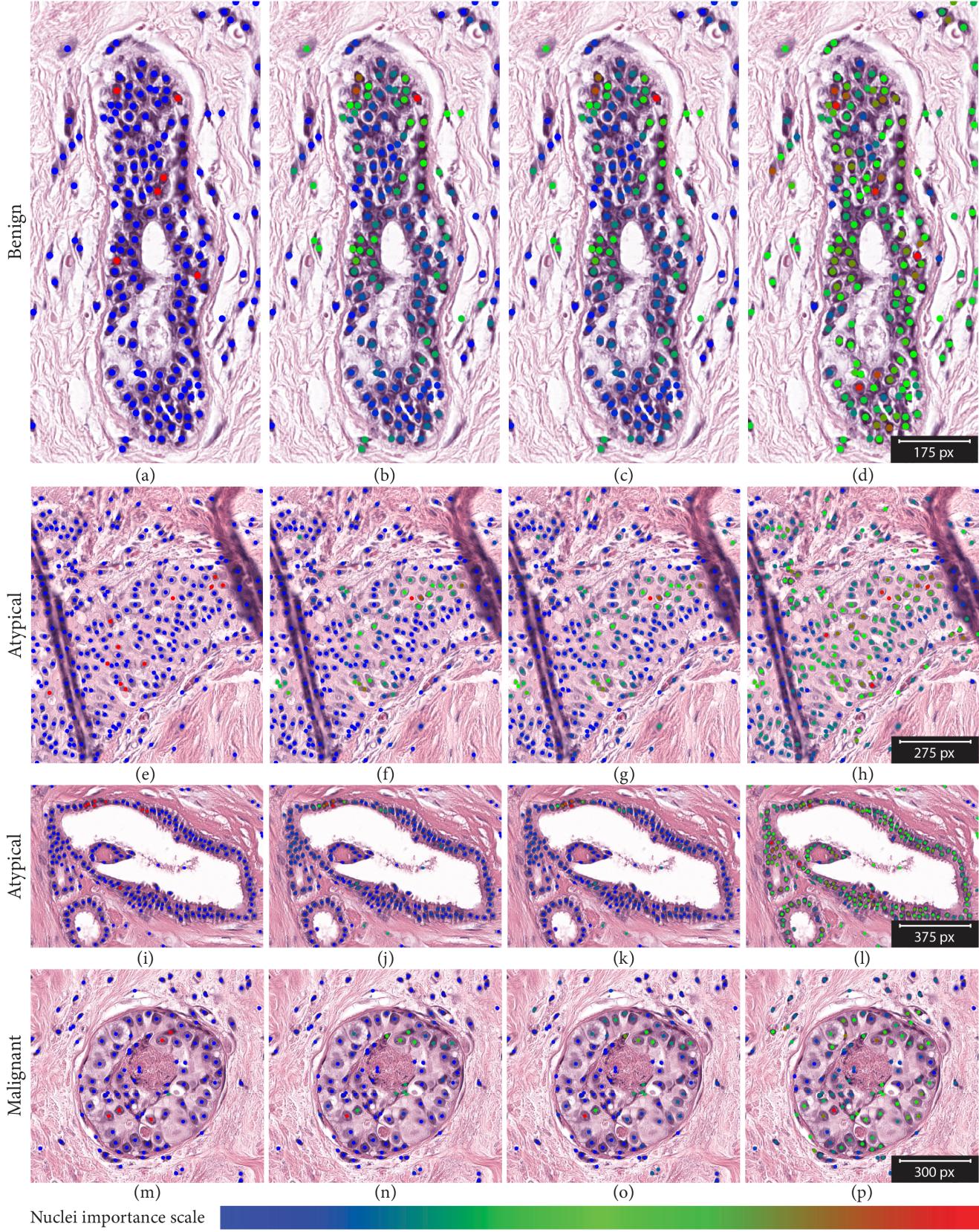
Figure 8. Qualitative results. The rows represent breast cancer subtypes, and columns represent graph explainers, *i.e.* GNNEXPLAINER, GRAPHGRAD-CAM, GRAPHGRAD-CAM++, and GRAPHLRP. Nuclei level importance ranges from blue (the least important) to red (the highest important).