

Shedding light on the black box of a neural network used to detect prostate cancer in whole slide images by occlusion-based explainability

Matej Gallo^{a,*}, Vojtěch Krajčanský^{a,1}, Rudolf Nenutil^b, Petr Holub^c, Tomáš Brázdil^a

^a Faculty of Informatics, Masaryk University, Botanická 68a, 602 00 Brno, Czech Republic

^b Department of Pathology, Masaryk Memorial Cancer Institute, Žlutý kopec 7, 656 53 Brno, Czech Republic

^c Institute of Computer Science, Masaryk University, Šumavská 416/15, 602 00 Brno, Czech Republic

ARTICLE INFO

Keywords:

Machine learning
Digital histopathology
Explainable AI
Artificial intelligence
Occlusion sensitivity analysis
Prostate cancer

ABSTRACT

Diagnostic histopathology faces increasing demands due to aging populations and expanding healthcare programs. Semi-automated diagnostic systems employing deep learning methods are one approach to alleviate this pressure. The learning models for histopathology are inherently complex and opaque from the user's perspective. Hence different methods have been developed to interpret their behavior. However, relatively limited attention has been devoted to the connection between interpretation methods and the knowledge of experienced pathologists. The main contribution of this paper is a method for comparing morphological patterns used by expert pathologists to detect cancer with the patterns identified as important for inference of learning models. Given the patch-based nature of processing large-scale histopathological imaging, we have been able to show statistically that the VGG16 model could utilize all the structures that are observable by the pathologist, given the patch size and scan resolution. The results show that the neural network approach to recognizing prostatic cancer is similar to that of a pathologist at medium optical resolution. The saliency maps identified several prevailing histomorphological features characterizing carcinoma, e.g., single-layered epithelium, small lumina, and hyperchromatic nuclei with halo. A convincing finding was the recognition of their mimickers in non-neoplastic tissue. The method can also identify differences, i.e., standard patterns not used by the learning models and new patterns not yet used by pathologists. Saliency maps provide added value for automated digital pathology to analyze and fine-tune deep learning systems and improve trust in computer-based decisions.

1. Introduction

The increasing lifespan in developed countries inevitably leads to higher incidences of cancer due to the aging population. Expanding cancer screening programs and personalized medicine increases healthcare systems' workload, including diagnostic specialties such as radiology and histopathology. This effect is partly compensated by progress in digitization, providing more efficient processing, archiving,

and retrieval of medical records. The use of digitized medical images represents the next evolutionary step [1]. Radiology is more advanced in this respect, already routinely utilizing Picture Archiving and Communicating Systems. In contrast, comparable pathology systems utilizing whole slide images (WSIs) are currently being approved for diagnostic use [2] and introduced into routine workflows [3].

The availability of digitized WSI, and large scans of histopathological samples, typically sized from gigapixels to tens of gigapixels, provides

Abbreviations: AOPC, Area Over the Perturbed Curve; ASAP, Automated Slide Analysis Platform; AUC, Area Under Curve; DAB, 3,3-Diaminobenzidine; DeconvNet, Deconvolution Network; DTD, Deep Taylor Decomposition; EG, Expected Gradients; HER, Effective Heat Ratios; FFPE, Formalin-fixed Paraffin-embedded; FN, False Negative; FP, False Positive; FROC, Free-response Receiver Operation Characteristics; GB, Guided Backpropagation; GNN, Graph Neural Network; H&E, Hematoxylin and Eosin; I*G, Input*Gradient; IG, Integrated Gradients; LIME, Local Interpretable Model-Agnostic Explanations; LRP-z, Layer-Wise Relevance Propagation with Z-rule; LRP-ε, Layer-Wise Relevance Propagation with epsilon rule; MIL, Multiple-instance Learning; OSA, Occlusion Sensitivity Analysis; PDA, Prediction Difference Analysis; RISE, Randomized Input Sampling for Explanations; ROAR, RemOve And Retrain; SHAP, Shapley Additive Explanations; TCAV, Testing with Concept Activation Vectors; TMA, Tissue Microarray; TN, True Negative; TP, True Positive; WSI, Whole Slide Image; xPattern, Explained Pattern; xPOI, Explanation Point of Interest.

* Corresponding author.

E-mail address: 422328@mail.muni.cz (M. Gallo).

¹ Authors have contributed equally to this work.

<https://doi.org/10.1016/j.nbt.2023.09.008>

Received 16 March 2023; Received in revised form 29 August 2023; Accepted 30 September 2023

Available online 2 October 2023

1871-6784/© 2023 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

the basis for developing more advanced approaches. This includes using neural networks for image analysis to facilitate diagnostics and prognostics or to assist pathologists in reducing their workload in routine tasks [4–6]. The detection of prostate cancer in core biopsies is one example of such an application. This represents a typical and often tedious part of a pathologist's daily routine, where samples from patients identified as prostate-specific antigen-positive during screening are mostly negative. Nevertheless, all slides containing multiple tissue sections must be carefully evaluated for small cancer foci. Unsurprisingly, some of the first deep learning applications in WSI analysis were dedicated to this problem, providing impressive results [7–9].

The deep learning approach is sometimes criticized for not providing insights into its internal mechanisms, which raises issues of trust by the pathologists and may hinder acceptance by regulatory authorities [10, 11]. In addition, unraveling the precise features used by the system could provide additional information that would be useful for pathologists and for further training in the machine learning process itself (see e. g., [12], where deep learning revealed new information not directly attributable to known histomorphological characteristics). As a result, a large amount of work has already been devoted to the development of interpretable machine-learning methods. For example, some of these methods use probability heatmaps to identify the histopathological features associated with adverse prognosis in glioma [13]. Interpretable AI has already been successfully applied in other areas of healthcare. An exhaustive survey [14,15] shows applications in Alzheimer's disease detection, Parkinson's disease detection, COVID-19 detection, pneumonia diagnosis, or ophthalmic disease detection. Still lacking is a systematic mapping of outputs from the methods interpreting a given model to the knowledge of pathologists.

The main contribution of this paper is a method for comparing morphological patterns used by expert pathologists to detect cancer with the patterns identified as important for inference of learning models. The method starts with building a catalog of relevant morphological features pathologists use to recognize cancerous tissue. Intuitively, we measure how close the reasoning of the learning model is to that of pathologists. More specifically, we have estimated the proportion of morphological features important for the model that also belong to the catalog to all morphological features important for the model. The method allows us to decide whether a statistically significant portion of patterns is employed by both pathologists and learning models, thus enhancing the trustworthiness of the learning models. It can also identify differences, i.e. standard patterns not used by the learning models and new patterns not yet used by pathologists.

As a use case, we focus on prostate carcinoma. From the pathologist's point of view, the morphological criteria of malignancy in the prostate are well-established and relatively reproducible [16]. To demonstrate the method, we have trained a model based on VGG16 [17] using a dataset of cases from Masaryk Memorial Cancer Institute, Brno. The model is applied to the WSIs patch-wise, with a patch size of 512 px × 512 px at the resolution of 0.468 μm/px. The model predicts 256 px × 256 px center of each patch; hence, the patches overlap to cover a tissue on the WSI completely. The model shows the state-of-the-art performance of 98% AUC for patch-wise prediction and 100% slide-level AUC score. This shows that VGG16 is capable of high-quality patch-based segmentation on data from a single source.

To identify important morphological features, an Occlusion Sensitivity Analysis (OSA) [18] is utilized for WSI, where the explanations are stitched from explanations of individual patches which overlap in the case of our use case model. The advantage of the occlusion is that it is easily explainable even to the users of learning models without a deep understanding of their inner workings: Parts of the input images are systematically hidden from the model, and we observe how the predictions of the model change. OSA results are evaluated both by an expert as well as by automated metrics such as Causal Insertion and Causal Deletion [19], Area over Perturbed Curve [20], Sensitivity- n [21], and Effective Heat Ratios [22]. Using the automated metrics, OSA

are also compared with other well-known, more advanced methods for generating saliency maps. We evaluate against Input*Gradient (I*G) [23], Guided Backpropagation (GB) [24], Deep Taylor Decomposition (DTD) [25], LRP- ϵ [26], DeconvNet [18], and Integrated Gradients (IG) [27].

We analyze the model's behavior using the method described above and conclude that most morphological features designated as important for our trained VGG16 model by OSA are also important for pathologists. The analysis method shows that the VGG16 model effectively utilizes all the features that the pathologist can recognize, given patch size and resolution.

1.1. Related work

Processing entire raw WSIs spanning hundreds of thousands of pixels in height and width with three or four channels would be memory inefficient. Nearly all methods utilize a patch-based approach [28–45] in which a WSI is converted into a set of equal-sized patches that are only the fraction of a size of the original WSI. The WSI is then represented as either a set of patches [27–31,34–37,39–44] or further processed to emphasize spatial relationships between patches and represented as a graph [34,39].

Stained tissue often varies greatly in color due to the type of scanner being used or due to a chemical preparation of a slide. Some authors have reported better results when these variations are minimized using stain normalization [42,46–50].

The scarcity of well-annotated datasets has resulted in a shift from supervised learning utilizing fully annotated datasets to self-supervised learning [33,36,37] utilizing partially annotated datasets and weakly-supervised multiple-instance learning (MIL) [28,43,44] utilizing only WSI-level information about the presence or absence of tumors.

1.1.1. Explainability methods

Recently, more emphasis has been placed on how models reach their decisions. Some models contain a built-in mechanism as part of their architectures called attention to guide the analysis of input data [32,35, 38–40,51]. These can be examined to determine the most influential parts of the input features. Different techniques had to be developed for other models to assign importance scores to input.

Saliency maps are the most straightforward way of presenting such scores for images. Each pixel is assigned a score, called attribution, quantifying its contribution to the final prediction. Gradient-based methods (I*G [23], IG [27], Expected Gradients (EG) [52], Grad-CAM [53], SmoothGrad [54], GB [24]) use gradients to calculate the importance scores; perturbation-based methods (Shapley Additive Explanations (SHAP) [55–58], Local Interpretable Model-Agnostic Explanations (LIME) [59], CXPLAIN [60], Randomized Input Sampling for Explanations (RISE) [19], Prediction Difference Analysis (PDA) [61], OSA [18], Anchors [62]) produce saliency maps by perturbing the input image, usually by deactivating—e.g., zeroing—pixels [63,64]; propagation-based methods (DTD [25], Layer-Wise Relevance Propagation (LRP) [26]) use rules to redistribute the output score among the input features; and other explanation methods utilizing concepts (Testing with Concept Activation Vectors (TCAV) [65]), deconvolution (DeconvNet [18]), activation comparison (DeepLIFT [66]), special neural networks (Neural Additive Models (NAMs) [67]).

For graph neural networks, modifications exist of previously mentioned algorithms [68] such as GNN-LRP [69] and GraphLIME [70], as well as novel methods specifically designed to work with relational or graph structures (RelEx [71], GNNExplainer [72,73], XGNN [74], Sub-GraphX [75], SE-GNN [76], and ProtGNN [77]). Explanations based on related patches can be presented in the form of subgraphs.

1.1.2. Explainability evaluation

Evaluation of generated explanations is challenging [78], mainly because they are subjective, and their quality depends on factors such as

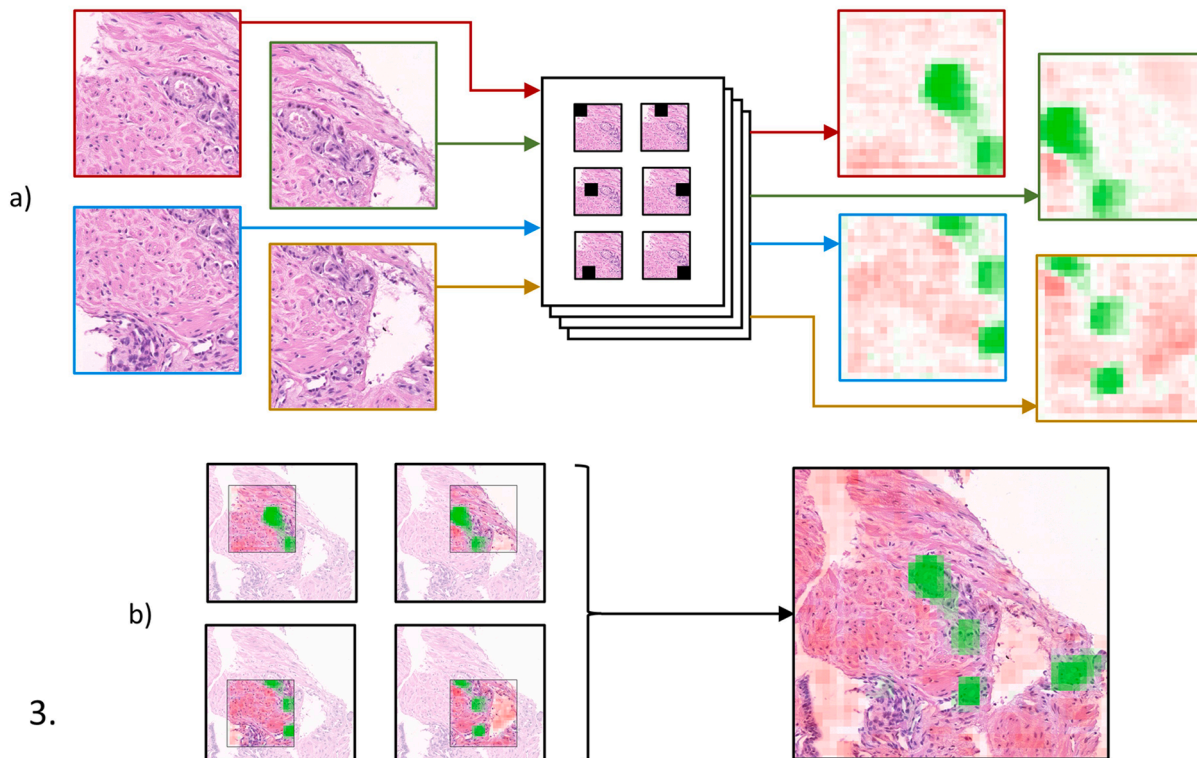
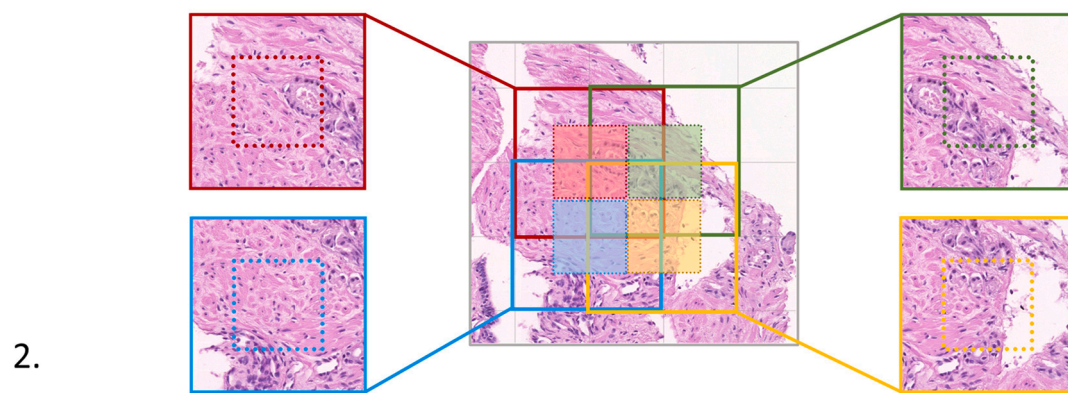
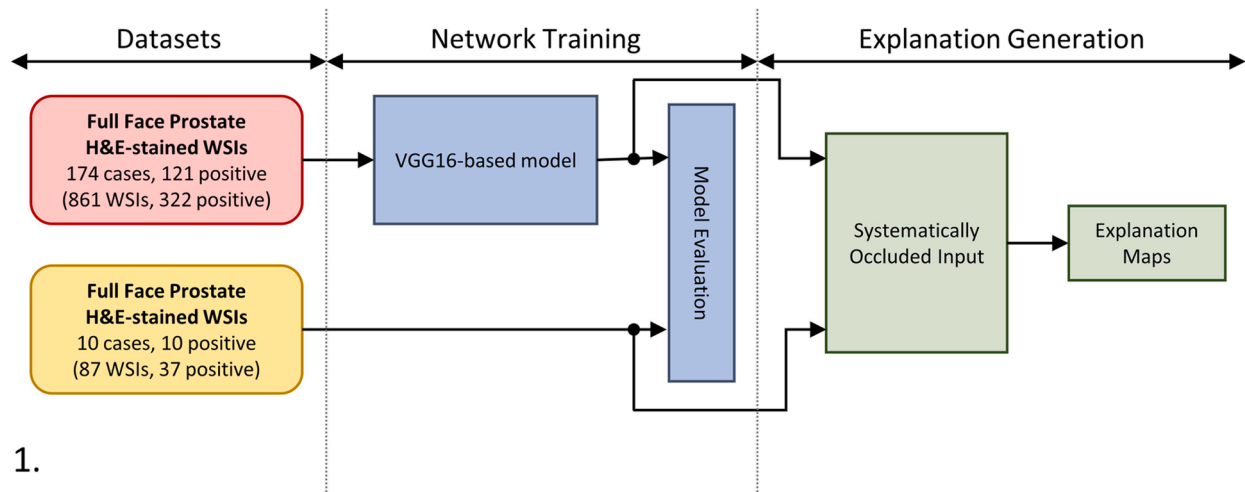


Fig. 1. Network training and patch saliency map generation. Section 1 shows the overall training and generation process; Section 2 illustrates how input patches overlap as a consequence of including surrounding contextual information (the whole patch) around the label-relevant central region (denoted by dotted squares); Section 3a) illustrates the systematic occlusion of individual patches and resulting occlusion patch saliency map for individual patches; Section 3b) shows how the final occlusion saliency map on a WSI is obtained by averaging the overlapping occlusion patch saliency maps.

the user's experience and goals. The authors of [79] presents an entire typology of explanations based on an explained system, the method by which it is being explained, and the relationship between the two. A good explanation should meet several criteria, such as interpretability (easy to understand) and fidelity/faithfulness (accurately describe explained model). Due to the subjective nature of the explanations, a common approach is to present the explanations to domain experts and subsequently gather feedback via questionnaires or interviews. A survey conducted in [80] suggests caution when selecting the form in which the AI results are presented to the user. Authors demonstrate how different forms of presentation, e.g., via saliency maps or counter-examples, may impact an AI system's perceived usefulness and trustworthiness. While our work also evaluates the explanations from the point of view of the pathologist, it is not done by means of interviews or questionnaires.

However, several automated methods exist to measure the soundness of explanations [81,82]. Many of these are based on systematic attribution-based input perturbations and subsequent observation of performance degradation (Causal deletions and Causal insertions [19], Area over the Perturbed Curve (AOPC) [20], Sensitivity- n [21], and RemOve And Retrain (ROAR) [83]). Effective Heat Ratios (EHR) [22] measure the overlap between the saliency maps and ground truth annotations.

Several studies have been conducted analyzing the overall impact of AI-assisted diagnosis. In [84], the authors demonstrate an AI solution that performs at a level equal to pathologists. In some cases, such as detecting Gleason pattern four types of prostate cancer, it even surpasses the pathologists' detection rate. They also demonstrated that pathologists, with the help of an AI, achieved higher labeling consistency and a significant reduction of time spent on each slide.

2. Material and methods

2.1. Material

To train and test deep learning models, we use the dataset of WSIs stained with hematoxylin/eosin (containing 3–5 tissue core sections each) that are part of the digital archive at the Department of Pathology, Masaryk Memorial Cancer Institute, Brno. They were scanned using a Panoramic® MIDI scanner (3DHistech, Budapest, Hungary) with a 20x objective lens at a 0.172 $\mu\text{m}/\text{pixel}$ resolution. The WSIs were stored in MIRAX format as uncompressed PNG images. Each WSI is a large image of dimensions 105,185 px \times 221,772 px. The dataset consists of the following:

1. Training WSIs, obtained from 157 consecutive core biopsies (104 patients with carcinoma, 53 negative. The distribution of WHO grade groups was as follows: 1: 38, 2: 31, 3: 16, 4: 9, 5: 10); in total, 264 WSIs with cancer and 436 without cancer.
2. Test WSIs, obtained from 10 patients with cancer, selected from additional consecutive cases to represent different Gleason patterns and types of infiltration; in total, 37 WSIs with cancer and 50 without cancer. The distribution of WHO grade groups was as follows: 1: 5, 2: 1, 3: 1, 4: 1, 5: 2.

Details are given in supplementary data (Supplementary Table S1.1).

The WSIs were checked in the Automated Slide Analysis Platform (ASAP) (<https://computationalpathologygroup.github.io/ASAP/>) [85], and all biopsy cores containing carcinoma were manually annotated in ASAP for further analysis. Annotations were performed as polygons containing carcinoma areas. The whole process is presented as a flow-chart in Fig. 1.

2.1.1. Dataset access

The dataset is available as raw files stored in Mirax MRXS format (<https://openslide.org/formats/mirax/>) compatible with the OpenSlide library [86]. Annotations used for the evaluation stage are available as

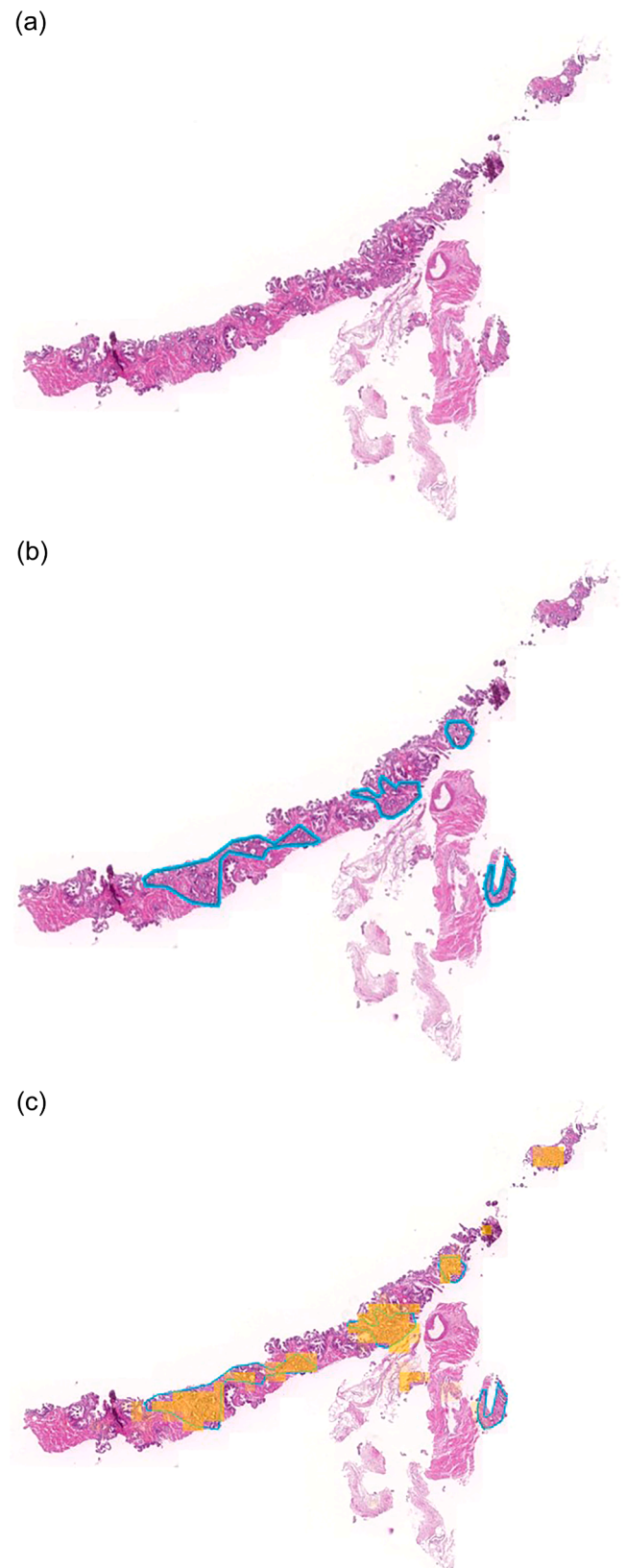


Fig. 2. Example of the WSI workflow. a) Original WSI; b) WSI manually annotated by the pathologist; c) WSI for validation, comparing model inference with manual annotation.

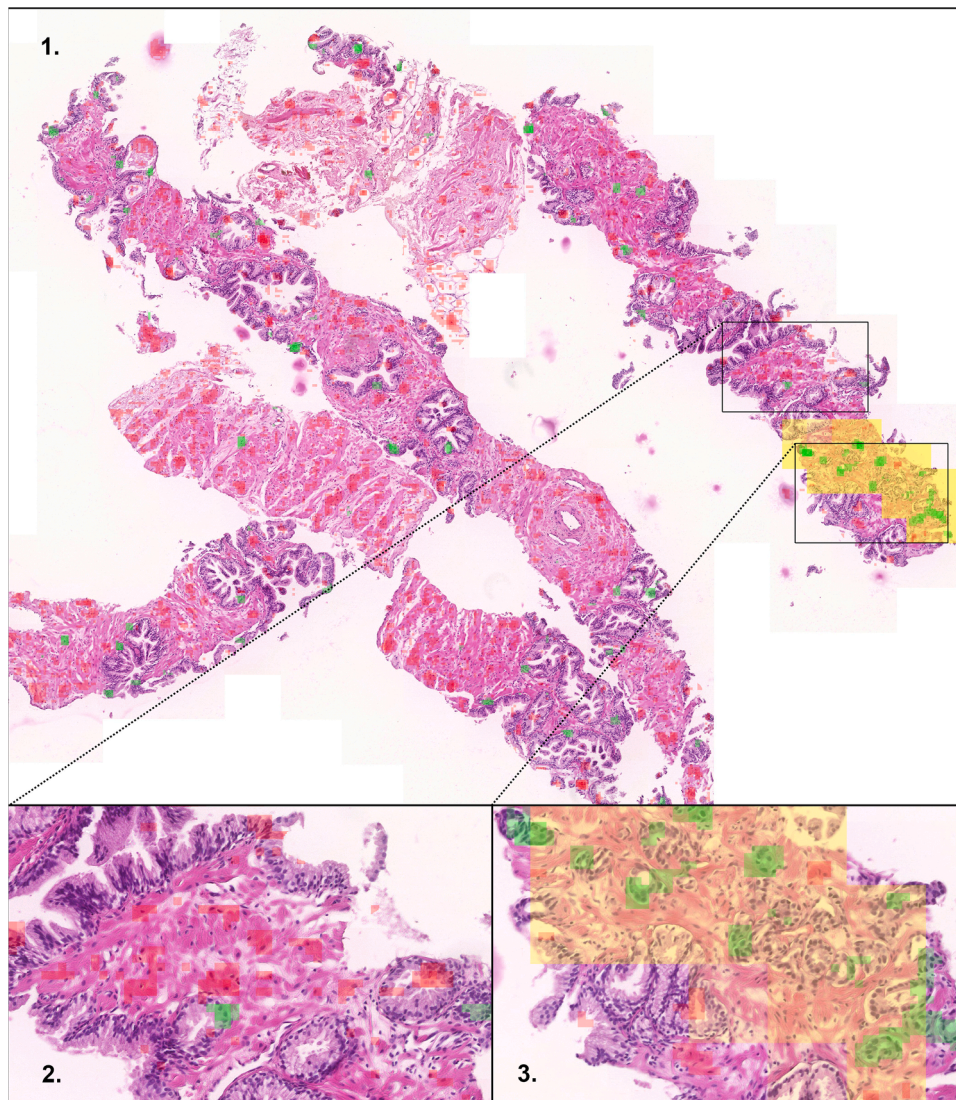


Fig. 3. Network predictions and occlusion saliency map overlays. The prediction layer (yellow) labels the focus of Gleason pattern 3 carcinoma. The explanatory layer (green and red) labels those parts of the image critical for estimating carcinoma probability using OSA. Regions with positive attributions, i.e., supporting the classification of a patch as malignant, are green; while those with negative attribution, i.e., suppressing such a classification, are red. The details below show parts of a WSI with carcinoma (right) and without carcinoma (left).

XML files compatible with ASAP [85]. The dataset is pseudonymized, and access can be requested via BBMRI-ERIC European Research Infrastructure by following its access policy (<https://www.bbMRI-eric.eu/services/access-policies>); the request should be placed via BBMRI-ERIC to Masaryk Memorial Cancer Institute. (Potential requesters can use https://directory.bbMRI-eric.eu/#/collection/bbMRI-eric:ID:CZ_MMCI:collection:LTS and add that to the selected collections and proceed to request samples via BBMRI-ERIC Negotiator platform.).

2.2. Machine learning models

The training set consists of WSI annotated with a pixel-level segmentation of cancerous tissue (see Fig. 2). Note that even though the labeling is pixel-level, the precision of tumor tissue delineation cannot be as precise due to the inherent difficulty of finding exact borders of cancerous tissue. We train a neural network to do such segmentation more coarsely.

We proceed according to the well-established approach used in several papers on deep learning and WSIs, e.g., [28–32,38–45,87]. A given WSI is cut into overlapping patches of size 512 px × 512 px with

stride 256 px (see Fig. 1, step 2). We concentrate on the following problem.

Problem statement: Classify patches according to the presence/absence of cancer in their central square area of size 256 px × 256 px.

More precisely, patches serve as inputs for a binary classifier that decides whether a given patch's central area intersects the cancerous areas of the tissue. Outputs for all patches of a given WSI can be organized into a coarse heatmap segmenting the tumorous tissue in the WSI with a precision of 256 px (see Fig. 2). See Section S1.1 of supplementary data for more detailed description of our model training.

2.3. Explainability analysis

2.3.1. Explainability using occlusion sensitivity analysis

To demonstrate our method for evaluation of explainability, a simple method of OSA is employed, which highlights regions of the WSI having a large impact on the output value of the model. The analysis results in an *occlusion saliency map* with the same dimensions as the input WSI. For each point, it specifies its attribution, i.e., numerically, how significant its impact is on the model's output.

OSA is applied on a per patch basis, giving an occlusion patch saliency map for each patch. The resulting patch maps are combined into a single occlusion saliency map for the whole WSI. We assume that each patch P has its own local coordinates relative to coordinates in the WSI. Concretely, assume that P is a patch whose upper left corner lies at position (u, v) in the WSI. Now considering $i, j \in \{0, \dots, 511\}$ the local position (i, j) in the patch P corresponds to the "global" position $(u+i, v+j)$ in the WSI.

2.3.2. Generating occlusion patch saliency maps

Let us consider an input patch P of size $512 \text{ px} \times 512 \text{ px}$, as described in the previous section. Intuitively, OSA is based on systematically covering (occluding) square regions of P by setting their values to zero and recording the changes in the output value of the model. The idea is summarized in Fig. 1, step 3a.

Concretely, we occlude square regions of size $55 \text{ px} \times 55 \text{ px}$ with a stride of 25 px. That is, for every pair of indices $i, j \in \{0, \dots, 19\}$, we consider an input patch P_{ij} obtained from P by zeroing out the square region of size $55 \text{ px} \times 55 \text{ px}$ with the upper left corner at the position $(i * 25, j * 25)$ in P .

Now consider a function F computed by our model, i.e., given the input patch P , the value $F(P)$ is the probability of a tumor present in P . Assume that $F(P) = \sigma(f(P))$ where σ is a logistic sigmoid (our output activation function) and f returns the logit output of the last layer.

The effect of perturbing P to P_{ij} can be measured using the difference $f(P) - f(P_{ij})$. However, note that using just this difference would result in wildly different scales of values for different patches. Conversely, taking the difference $F(P) - F(P_{ij})$ might result in too small differences due to saturation at the sigmoid. The aim is to combine the occlusion patch saliency maps into a single occlusion saliency map on the WSI, so the differences need to be normalized into (roughly) the same range.

Our solution is to apply the logistic sigmoid to the difference. Formally, for all patches P , define the occlusion patch saliency map O^P by

$$O_{ij}^P = \sigma(f(P) - f(P_{ij})) \text{ for } i, j \in \{0, \dots, 19\}$$

By upsampling the occlusion patch saliency map O^P , using the nearest-neighbor upsampling method, we obtain an occlusion patch saliency map M^P of size $512 \text{ px} \times 512 \text{ px}$ for each patch P . Given $k, l \in \{0, \dots, 255\}$, we write M_{kl}^P to denote the value of the map at the position (k, l) .

Finally, each occlusion patch saliency map is linearly scaled from $[0, 1]$ to $[-1, 1]$ range (using the transform $2(M^P - \frac{1}{2})$) so that negative values correspond to evidence against cancer, and positive values correspond to evidence for cancer. Zero represents no effect on the output.

2.3.3. Generating occlusion saliency maps for WSI

The occlusion saliency map for the whole WSI is obtained by stitching the occlusion patch saliency maps M^P for individual patches P (Fig. 3). The problem is that the patches overlap, so each pixel in the WSI belongs to up to 9 different patches. This problem is solved simply using averaging. For illustration, see Fig. 1, step 3b.

The final occlusion saliency map M is defined as follows. Consider a point at the position (u, v) in the input WSI. This point is covered by patches P^1, \dots, P^n . For each patch, P^i , the point at position (u, v) in the WSI corresponds to a local position (k_i, l_i) in P^i . We define M_{uv} by

$$M_{uv} = \frac{1}{n} \sum_{i=1}^n M_{k_i l_i}^{P^i}$$

As observed in our experiments, this averaging has a reasonably strong smoothing effect giving more homogeneous regions of similar values corresponding to known tissue patterns (see **Saliency map evaluation**).

2.3.4. Other explainability methods compared

The following methods were compared visually and by automated metrics presented earlier against OSA: Input*Gradient (I*G), Guided Backpropagation (GB), Deep Taylor Decomposition (DTD), Layer-Wise Relevance Propagation (LRP), Deconvolution (DeconvNet), and Integrated Gradients (IG).

The I*G [23] method calculates the derivative of the output given the input and multiplies the resulting gradients with the input image. Similarly, the GB [24] method computes the gradients of the output given the input using gradient backpropagation. However, only the positive gradients can pass through the ReLU activations while the gradients are propagated through the network. DTD [25] and LRP- ϵ [26] both redistribute relevancy on a layer-by-layer basis to the input features. DeconvNet [18] is based on applying transposed convolution to the feature maps and upsampling the results to the input resolution. IG [27] first samples several points in the input image space lying on the line segment connecting the input image x and a fixed reference image x^R (completely black patch in our case) and calculates the gradients of model output for each sampled point. The final feature attributions are obtained by summing the gradients.

We use concrete implementations of the above methods from the package iNNvestigate [88]. The saliency maps for all methods have been obtained: the same test set and the trained model described earlier were used. The analyzer from the iNNvestigate package retrieved patch saliency maps for each patch. Similarly to the approach described in **Explainability using occlusion sensitivity analysis**, the WSI-level saliency maps were stitched from patch-level saliency maps and the overlapping areas averaged. For each method, the saliency maps were scaled to the $[-1, 1]$ range where -1 represents strong evidence against the patch containing the pixel being classified as cancer, and 1 represents strong evidence in favor of classifying the patch as cancer.

2.4. Manual evaluation of saliency maps

Our main contribution is an evaluation procedure measuring the quality of explanations by comparing the saliency map with the labeling of relevant morphological features by pathologists.

2.4.1. Explanation points of interest

To evaluate a given saliency map, we identify its explanation points of interest (xPOIs). Roughly speaking, xPOI is a point in the saliency map surrounded by sufficiently large positive/negative attributions. The xPOIs identify locations in the input WSI with a large positive/negative impact on the model output. Note that the higher attributions around xPOIs are demanded, the fewer xPOIs we get. Hence, there is a need to strike a balance between how high the attributions have to be in order for the explanations to be reliable and how well the xPOIs cover the input WSI.

In the case of the *occlusion saliency map*, the xPOIs are defined as follows: A point in the map is an xPOI iff it is the center of a square region of dimensions $15 \text{ px} \times 15 \text{ px}$, where the absolute difference between the mean positive and mean negative attributions is greater than 0.55. The threshold 0.55 has been selected experimentally to achieve a reasonable coverage of patches with evenly spaced xPOIs (see **Sampling xPOIs**).

The difference between the mean positive and mean negative attributions can be either positive or negative, which gives *positive xPOI* and *negative xPOI*, respectively.

2.4.2. Classification of xPOIs

xPOIs are classified in two ways: based on the surrounding biologically relevant morphological features and whether they lie within a cancerous epithelium.

Pathologists recognize tumorous tissue using various forms of morphological features [3]. We have collected many of these features in

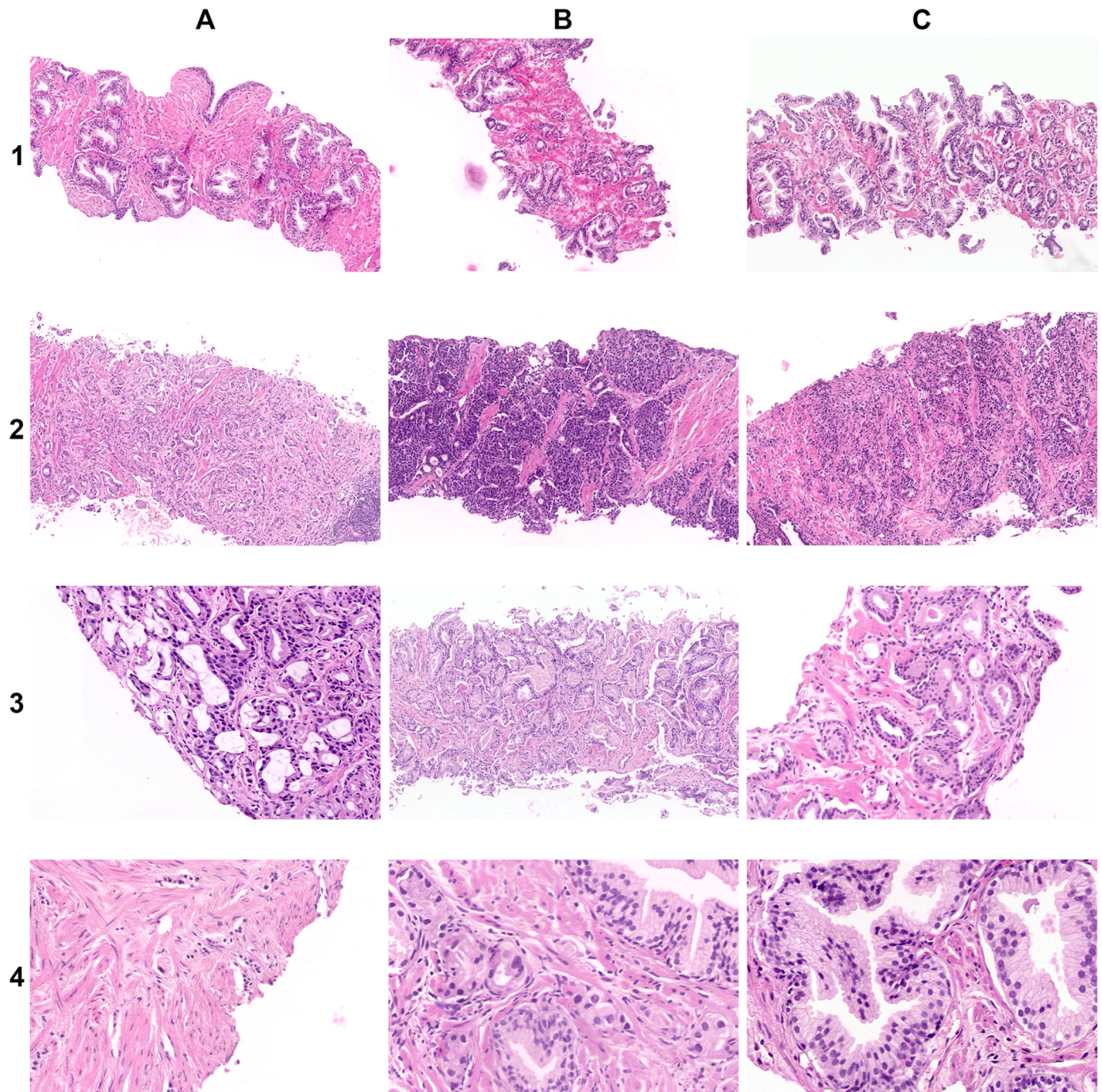


Fig. 4. Catalogue of selected morphological features, characterizing prostate cancer and non-cancerous tissue in core biopsies. The non-cancerous tissue (A1, C4) usually has a regular lobular architecture with isolated glands and a relatively balanced epithelium/stroma ratio or stromal predominance (A1). The glandular epithelium is two-layered. The cells have uniform smaller nuclei, and the nucleoli are not visible (C4). The luminal cells are highly polarized with abundant cytoplasm (C4). Conversely, carcinoma (A2 to A4, B1 to B4, C1 to C3) is characterized by distorted gland architecture. In Gleason grade 3 (B2, C1, C3) the small caliber and relatively uniform glands consist of single-layered epithelium (A3, C3), which may have periglandular clefts (C3) and can infiltrate in between normal glands (C1, B4). Tumor glands have rigid, sharp lumina (C1, A3, C3), which may contain blue mucin (A3), or crystalloids (C3). Gleason grade 4 (A2, A3, B4, C2) exhibits poorly formed fused cribriform or glomeruloid glands (A2), solid sheets, cords, medium or large nests with rosettes (B2, C2), and high nuclear density (B2). Gleason grade 5 (A4, C2, B3) exhibits infiltrative single cells and small cell groups (A4) or a large amount of necrotic debris within glands (B3). In most cancers, the cells are cuboidal to low cylindrical with modest cytoplasm (A3, C3) and have enlarged hyperchromatic nuclei (B4).

[Table 8](#) (pro-cancer features) and [Table 9](#) (non-cancer features). Most of these morphological features can be recognized based on their composition of simple patterns, such as single chains of nuclei (indicating single-layer epithelium) or small holes (indicating small lumina).

Now visual analysis of occlusion saliency maps around *some* xPOIs revealed explainability attributions forming continuous regions of size

about 50 μm . The corresponding tissue regions highlighted by these attributions contained morphological structures matching the simple patterns ([Fig. 5](#)). We have identified seven prominent patterns, four typically forming pro-cancer morphological features and three typical for non-cancer tissue. These seven simple patterns are termed the *explained patterns* or, in brief, *xPatterns* ([Fig. 6](#)).

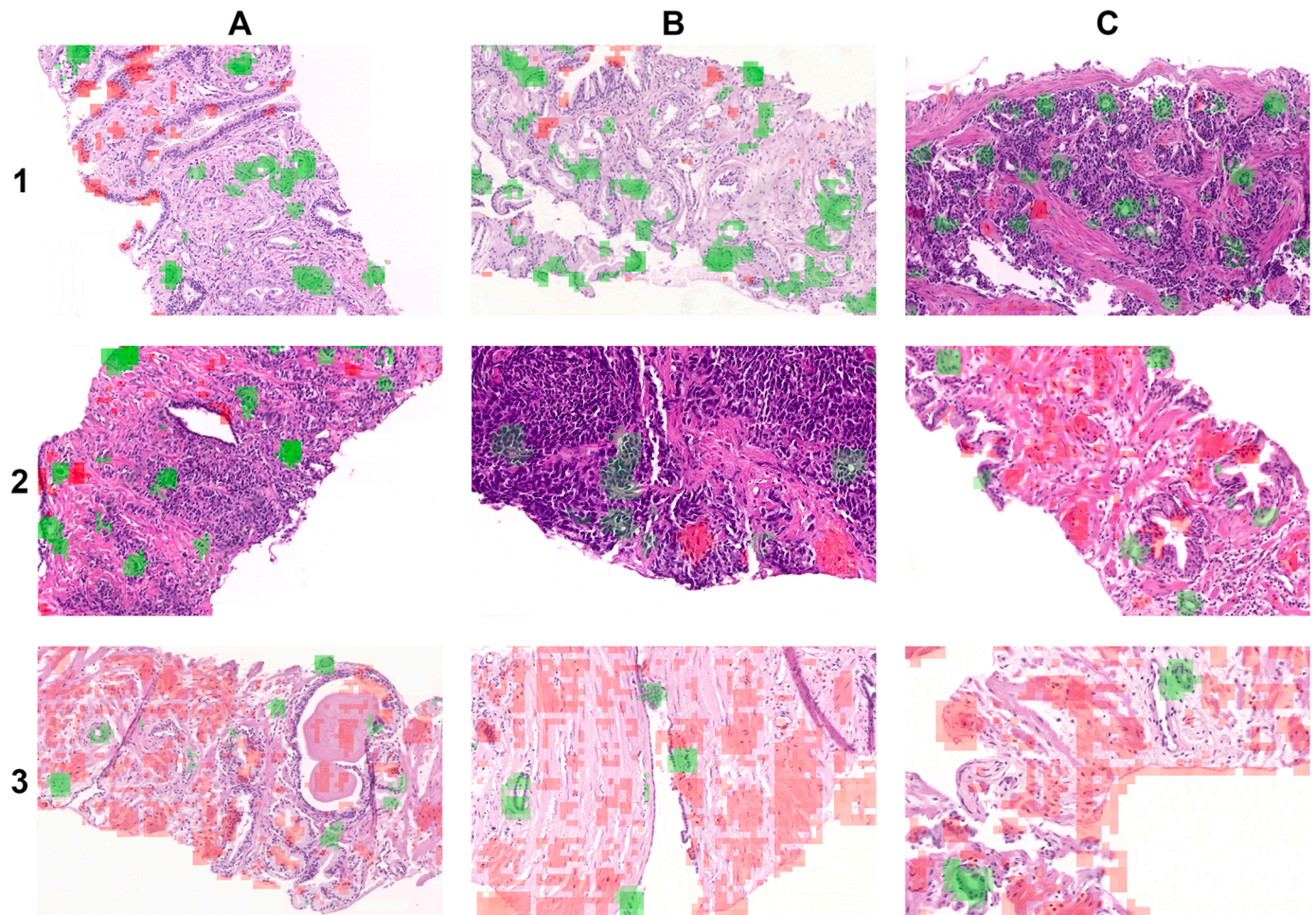


Fig. 5. Overview of occlusion saliency maps in prostate cancer and non-cancerous tissue. The explanatory layer comprises labeling attributions of diameter around 50 μm . The visual analysis of these attributions reveals they are located above some repeating image patterns related to morphological features, e.g., small round holes, high nuclear density, single chain of nuclei, and hyperchromatic nuclei with perinuclear halo (in green, favoring cancer) or areas of low nuclear density with eosinophilic background, two-layered chains of nuclei or chain of nuclei with abundant slightly eosinophilic neighborhood (in red, favoring non-cancer). The feeling that the location of labels is not stochastic is supported by the recognition of some morphological mimickers of attributions favoring cancer in non-cancerous tissue (C2, A3, B3, C3). Note a small blood vessel with activated endothelium in A3 and C3. This structure is often labeled cancerous due to the simulation of a small round hole or single chain of nuclei patterns. Additional examples of these patterns in the broader tissue context can be found in [Supplementary Fig. S2.1](#).

An xPOI is said to lie within a given xPattern if it lies within an instance of the xPattern, i.e., within a specific morphological structure in the tissue that can be classified as the given xPattern. It is shown that most of the xPOIs lie within instances of xPatterns.

Four pro-cancer xPatterns are found:

1. Single chain of nuclei (single-layered epithelium)
2. Small round hole (small lumina)
3. High nuclear density (high cellular density)
4. Larger nucleus with perinuclear halo (hyperchromatic nuclei with halo)

Likewise, three non-cancer xPatterns are found:

1. Two-layered chain of nuclei (two-layered epithelium)
2. Areas of low nuclear density with eosinophilic background (stromal predominance)
3. Chain of nuclei with abundant slightly eosinophilic neighborhood (highly polarized epithelium)

Finally, each xPOI is classified based on its corresponding xPattern and its relative position to cancerous epithelium:

- Given a positive xPOI lying within the xPattern P_i , the xPOI is said to be
 - *true positive with P_i (TP_i)* if it lies within the cancer epithelium
 - *false positive with P_i (FP_i)* if it lies outside the cancer epithelium
- Given a negative xPOI lying within an xPattern N_i , the xPOI is said to be
 - *true negative with N_i (TN_i)* if it lies outside the cancer epithelium
 - *false negative with N_i (FN_i)* if it lies within the cancer epithelium
- We have found that no positive xPOI lies within any xPattern N_i , and no negative xPOI lies within any P_i .
- Given an xPOI that does not lie within any xPattern, such an xPOI is said to be undefined.

Thus, to evaluate the saliency map, the number of xPOIs in each category is counted: TP_i , FP_i , TN_i , FN_i , and undefined (see [Table 7](#)). The more xPOIs belong to TP_i and TN_i , the better. Note that xPOIs counted as FP_i and FN_i correspond to mimickers – xPOIs lying within xPatterns contradicting the tumor annotation ([Fig. 7](#)).

2.4.3. Sampling xPOIs

As delineating all occurrences of the xPatterns in all testing slides would be too burdensome for the pathologist, we have decided to use statistical evaluation and identify xPatterns only in randomly sampled

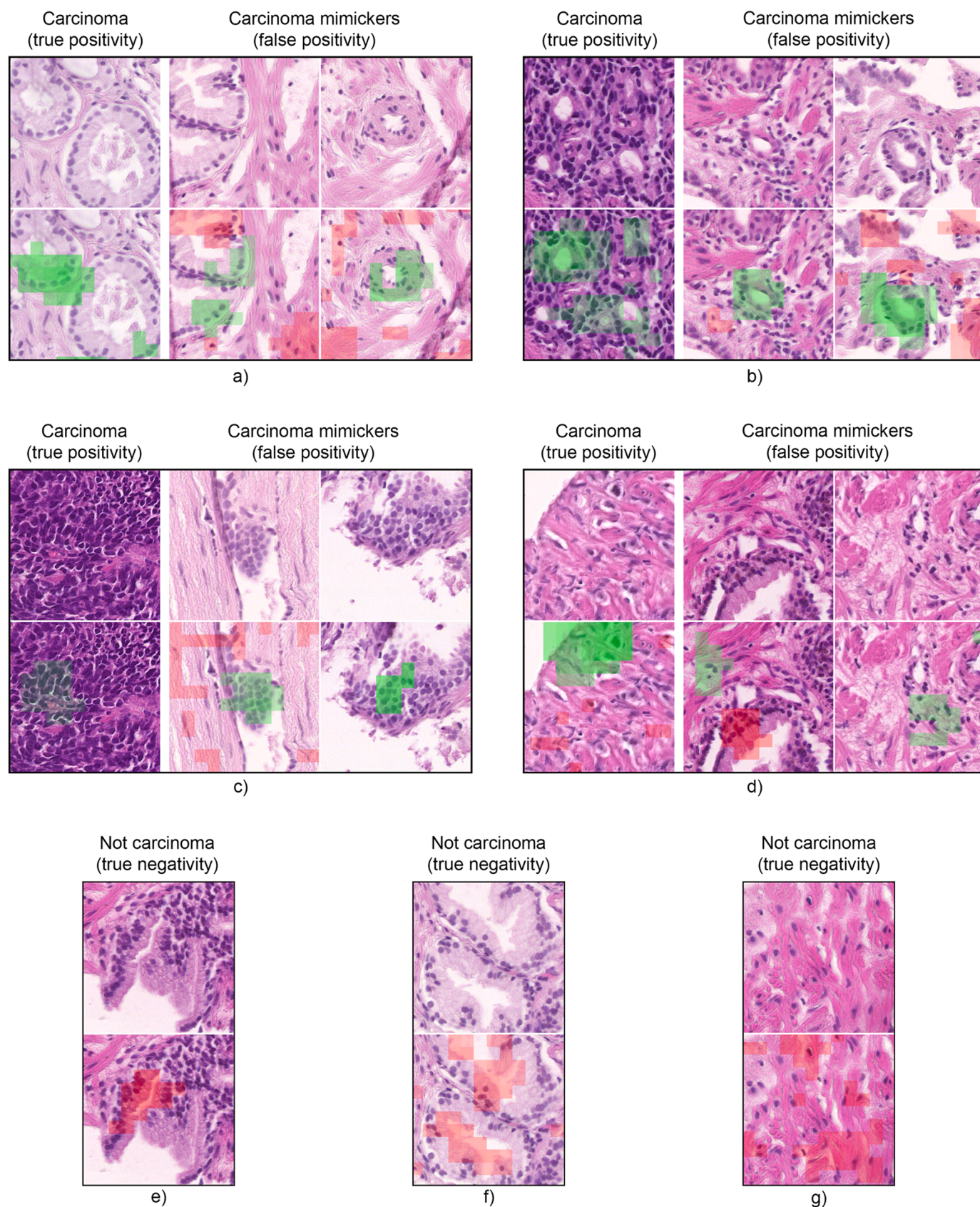


Fig. 6. Examples of identified xPatterns from occlusion saliency map xPOIs. a) single chain of nuclei, b) small round hole, c) high nuclear density, d) larger nucleus with perinuclear halo, e) two-layered chain of nuclei, f) chain of nuclei with abundant slightly eosinophilic neighborhood, g) areas of low nuclear density with eosinophilic background. The benign structures are vessels, reactive glands, or dense benign epithelium with slight nuclear enlargement. Detailed discussion of the identified xPatterns is in **Manual evaluation of occlusion saliency maps**.

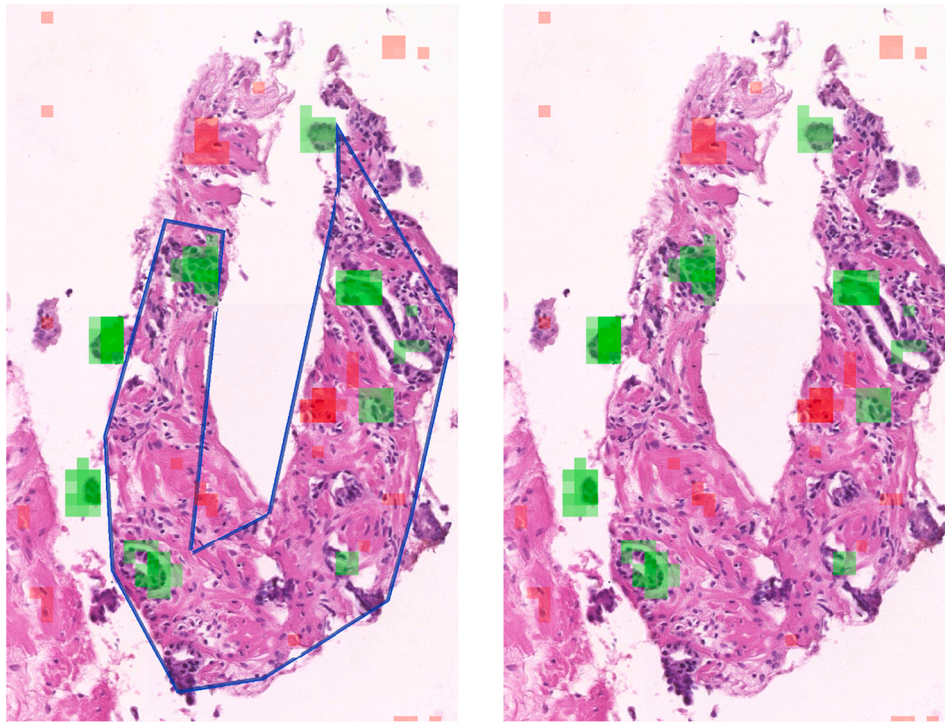


Fig. 7. Examples of FN misclassification (low network response). Part of Gleason pattern 4 + 5 carcinoma with sparse tumor islands infiltrating fibromuscular stroma. Although carcinoma foci are labeled correctly, the strong and frequent negative labels in stroma seem to push the decision to negativity. (Cancer predictions in yellow; manual annotations in blue; positive and negatively contributing regions in green and red, respectively.).

xPOIs. Random sampling tries to cover the whole tissue figure evenly, i. e., intuitively, avoiding clusters of xPOIs if other parts of the tissue are not covered yet. The sampled xPOIs are subsequently classified, and the classification is statistically evaluated.

The sampling of xPOIs is performed as follows. For each WSI, we consider intersection points of a coordinate grid of 280 px increment in both x and y dimensions. Considering only xPOIs on the intersections of the grid, 0.33% of intersections met the requirements of a valid xPOI, with 0.11% being positive and 0.22% negative. The sampling population consists of all intersection points that are also xPOIs from all 87 test WSIs (37 with cancer and 50 without cancer). During the actual sampling, the maximum number of samples is limited from a single WSI to 20. The mean sample count for a single slide is approximately 7.5, with extremes being 20 samples from one slide and one sample in the case of two slides.

2.5. Automated evaluation of saliency maps

To compare OSA method with the other explainability methods presented in **Other explainability methods compared**, the following metrics are employed measuring the faithfulness and clarity of saliency maps: Causal Deletion and Causal Insertion [19], Area over Perturbed Curve [20], Sensitivity- n [21], and Effective Heat Ratios [22].

Sensitivity- n is based on correlating the output deviation in response to perturbing n randomly selected input pixels with the sum of the pixels' attributions. Causal Insertion measures changes in the model's sensitivity in response to progressively adding pixels in the order of their attribution. The metric itself is the area under the curve plotting the sensitivity against the number of added pixels. Causal Deletion is the same, except that the pixels are removed. AOPC is similar to Causal Deletion, except that the result is an average of output deviations in response to the progressively perturbed input images. In the present case, the Effective Heat Ratios method gives the number of pixels with high enough attribution within the tumor annotation divided by the total number of pixels with high enough attribution.

Table 1

Patch-wise evaluation metrics on test set for Prostate and CAMELYON16 datasets.

Metric	Precision	Recall	AUC	Specificity
Camelyon	0.793	0.960	0.988	0.979
Prostate	0.823	0.925	0.981	0.952

Table 2

Comparison of patch-level performance on CAMELYON16 test set.

Method	Precision	Recall	AUC	Specificity
VGG16	0.793	0.960	0.988	0.979
SAMIL[82]	0.921	0.972	0.953	-
DeepGAT[83]	0.951	0.930	-	0.994
ResNet50[83]	0.888	0.819	-	0.987
DenseNet[83]	0.889	0.837	-	0.989
YOLOv4-GCPANet[84]	0.936	0.680	-	-
Spatial-ResNet34[85]	0.957	0.917	-	-

3. Results

3.1. Results for classification

Table 1 presents the patch-wise evaluation of the model on our prostate test set and a reference CAMELYON16 [89] test set. (Note that we do not evaluate tumor-level metrics (such as free-response receiver operating characteristic (FROC) [80]) as our main aim is the patch-level segmentation and optimizing FROC typically involves non-trivial post-processing of model outputs.) To obtain precision and recall, the output of the model is thresholded at 0.5. **Table 2** presents a comparison of the model to other contemporary architectures.

For completeness, the slide-level performance of the model was evaluated using a simple max-pooling strategy over the patch-wise predictions (each slide receives a score equal to the maximum score of

Table 3

Causal Deletion score (C^-) and Causal Insertion score (C^+) for the selected methods.

ExAI Method	I*G	GB	DTD	LRP- ϵ	OSA	DeconvNet	IG
C^- AUC	0.037	0.013	0.089	0.040	0.166	0.022	0.030
C^+ AUC	0.131	0.094	0.852	0.128	0.837	0.164	0.209

its patches). Interestingly, the model achieved a 100% slide-level area under curve (AUC) score on the Prostate test set and a 98% slide-level AUC score on the CAMELYON16 test set.

Note that even though VGG16 is a very simple and relatively old model, its results are still competitive, especially in situations with limited data and computational power. Our findings are in line with [64]. This means that for illustrating our model-agnostic evaluation of saliency maps, the VGG16 is sufficiently robust.

3.2. Saliency map evaluation

3.2.1. Automated saliency map evaluation

In this section, the different explainability methods introduced in **Other explainability methods compared** are compared using metrics introduced in **Automated evaluation of saliency maps**.

3.2.1.1. Causal insertion/deletion. Note that a high Causal Insertion score suggests important pixels have received high attributions. Conversely, for Causal Deletion, a low score indicates that the attributions express importance correctly. Due to the method relying on sensitivity, negative patches do not influence the final score and were excluded from the calculations during the experiments. Table 3 shows the results for the tested methods on the Prostate test set.

3.2.1.2. Area over the perturbed curve (AOPC). The results obtained on the Prostate test set are summarized in Table 4. Unlike the Causal Insertion/Deletion, negative patches are included in its calculation.

3.2.1.3. Sensitivity- n . For our evaluation, 16 equally spaced values of n from 0.1 to 1.0 are considered, representing the fraction of pixels in the patch to be removed. Results are summarized in Table 5.

3.2.1.4. Effective heat ratios. Similarly to Causal Insertion/Deletion, this metric cannot be used to assess the quality of saliency maps for negative slides as these slides have no ground truth annotations.

3.2.1.5. Discussion of the automated explainability evaluation. Most of the presented metrics utilize the occlusion principle. Visual inspection of the saliency maps reveals a strong similarity between DTD and OSA (Fig. 8). These two methods tend to (positively) highlight larger homogeneous regions roughly corresponding to morphological structures in the image. GB assigns positive attributions to spaces around nuclei, while the nuclei themselves are systematically assigned negative attributions. The remaining methods are much noisier; thus, deciding whether the underlying tissue region is positive or negative evidence is much more challenging.

As seen in the following section, OSA and DTD mainly highlight structures relevant to cancer detection; the spaces around nuclei highlighted by GB bear little to no relevance in the context of cancer

classification; and the gradient-based methods highlight large structures that are irrelevant to cancer classification (Supplementary Figure S4.6).

Causal Insertion (Table 3), Sensitivity- n (Table 5), and I (Table 6) all assigned higher scores to the less noisy methods – DTD and OSA.

Note the Causal Deletion and Insertion (Table 3) scores for GB. Also, note that the highlighted regions around nuclei by GB are the first to be changed by these methods. A detailed view of the results reveals that deletion in these areas greatly influences the model's output, while insertion has little impact. Hence, it is concluded that the highlighted areas do not contain sufficient evidence in isolation. This is further corroborated by the negative Sensitivity- n score (Table 5) assigned to GB. This result may likely be attributed to the GB assigning high positive scores mainly to white areas around the nuclei.

Of particular interest are the results using the AOPC metric (Table 4). Evaluation on positive patches only (AOPC⁺) yields similar results for all methods. Including negative patches (AOPC*) significantly affects all but two methods (DTD and OSA). Judging from scores obtained using only negative patches (AOPC⁻), it is believed these big changes may be attributed to the methods highlighting irrelevant areas, as mentioned earlier in this section. The I scores (Table 6) agree with this conjecture since I awards a higher score to saliency maps in which positive attributions are concentrated within annotated regions.

3.2.2. Manual evaluation of occlusion saliency maps

The results of the evaluation are presented in Table 7. Overall, of 646 sample xPOIs, 253 (39.1%) were evaluated as true positive (TP1–TP4), 107 (16.5%) as false positive (FP1–FP4), 273 (42.2%) as true negative (TN1–TN3), and 0 (0.0%) as false negative (FN1–FN3). The 13 (2%) sample xPOIs where it was unclear what information was highlighted by the saliency map were labeled as “Undefined”. A single chain of nuclei represents the predominant true positive morphological xPattern (132, i.e., 52.2% of true positive xPOIs and 20.4% of total xPOIs), followed by small round hole (22.5% of true positive xPOIs) and high nuclear density (19.0% of true positive xPOIs).

Note the change of morphological pattern distribution related to the Gleason score in carcinoma tissue (inside the border) (Table 7, underlined, in italics), showing a shift from a single chain of nuclei and a small round hole, xPatterns P1 and P2, to a high nuclear density and a larger nucleus with perinuclear halo, xPatterns P3 and P4, which represents some internal control of reading validity. There is no apparent predominance among the false positive patterns, but due to the similarity to true positives, they provide much more information regarding explainability. In true negatives, areas of low nuclear density with eosinophilic background, in most cases related to the predominant stromal component (TN2, 25.1% of total xPOIs, 57.2% of all negative sample xPOIs), is the prevailing xPattern followed by a two-layered chain of nuclei (TN1, 11.1% of total sample xPOIs, 25.4% of all negative sample xPOIs). The false negative samples are too rare to provide any reasonable explainability information, and no such case was present in sampled xPOIs. The undefined morphological patterns represent 2% of total sampled

Table 5

AUC of Pearson Correlation Coefficient (PCC) over different values of n on the entire dataset.

ExAI Method	I*G	GB	DTD	LRP- ϵ	OSA	DeconvNet	IG
PCC AUC	0.164	-0.305	0.703	0.161	0.566	0.125	0.396

Table 4

AOPC on the entire dataset (AOPC*), on positive patches only (AOPC⁺), and on negative patches only (AOPC⁻).

ExAI Method	I*G	GB	DTD	LRP- ϵ	OSA	DeconvNet	IG
AOPC*	-0.211	-0.211	0.029	-0.208	0.028	-0.256	-0.192
AOPC ⁺	0.121	0.201	0.194	0.122	0.177	0.206	0.129
AOPC ⁻	-0.272	-0.280	0.005	-0.269	0.009	-0.333	-0.252

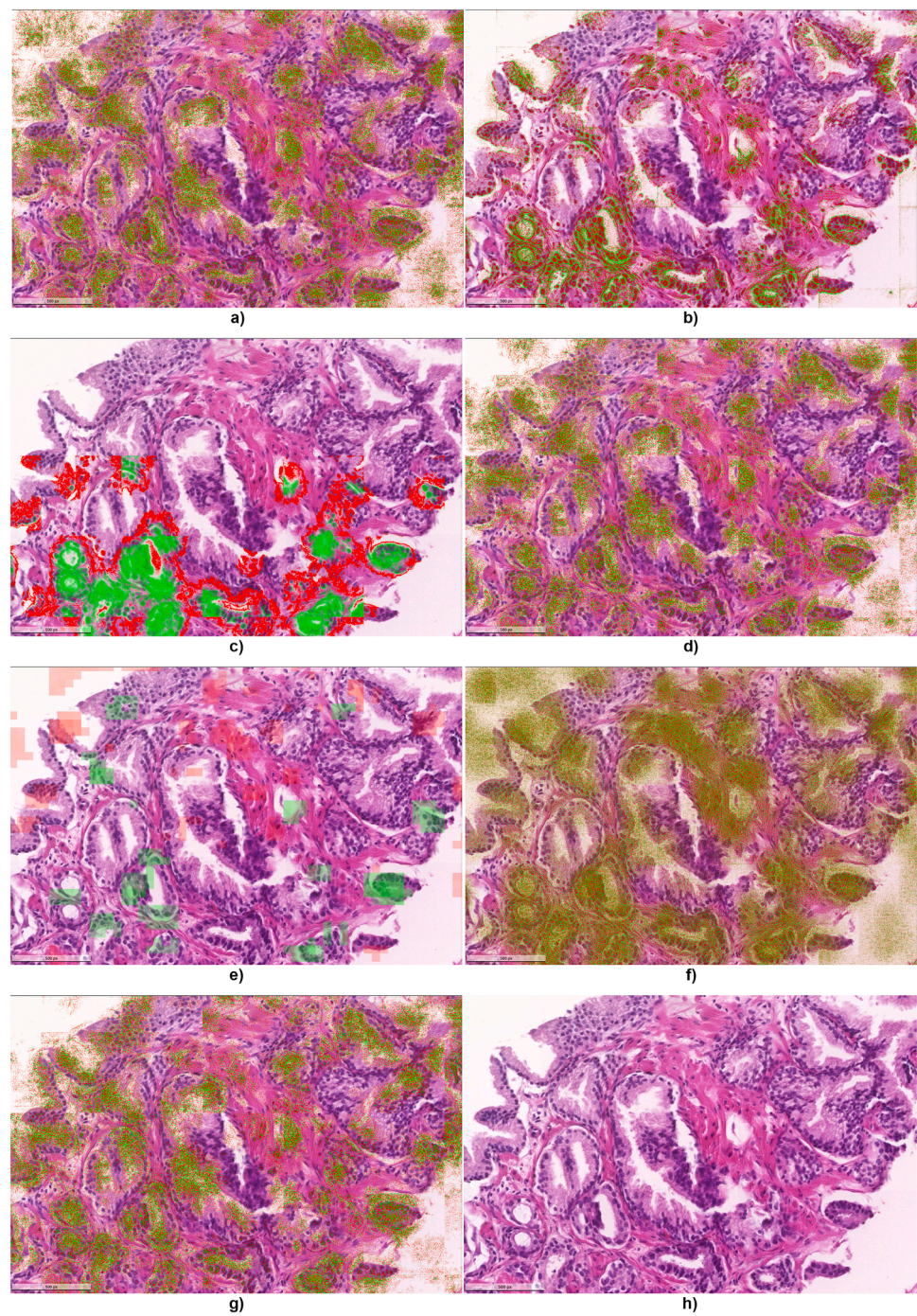


Fig. 8. Visual comparison of saliency maps produced by different methods. Gleason pattern 3 carcinoma with smaller glands infiltrating in between non-tumor glands. a) Input*Gradients, b) Guided Backpropagation, c) Deep Taylor Decomposition, d) LRP- ϵ , e) Occlusion Sensitivity Analysis, f) DeconvNet, g) Integrated Gradients, h) Original image. Channel-wise histogram equalization was applied on all methods except OSA to make the colors more visible.

Table 6

AUC of EHR ratios over quantiles on the entire dataset.

ExAI Method	I*G	GB	DTD	LRP- ϵ	OSA	DeconvNet	IG
EHR AUC	0.203	0.301	0.606	0.214	0.542	0.245	0.251

attributions across all categories.

Table 8 and Table 9 break down morphological features of cancer and non-cancer prostatic tissue into recurring patterns recognized by

OSA. Fig. 4 displays examples of morphological features, while Fig. 6 shows examples of recognized recurring patterns comprising the abovementioned features.

In subfigures a) through d) of Fig. 6 are presented the main morphological patterns underlying explainability attributions responsible for the classification of patches as malignant (first column, true positivity) and their mimickers in non-cancer tissue (another two right columns, false positivity). The explainability overlay (green) is provided in the bottom part. Note also some red labeled attributions representing true negativity. The “single chain of nuclei” can also be represented with an activated endothelium (subfigure a) third column). A small round

Distribution of morphological patterns under sampled explainability attributions in the test of WSIs. Detailed discussion of the results in the table is in Manual evaluation of occlusion saliency maps.

Morphological pattern under attribution	WSIs w/ carcinoma					Total (N = 37)	WSIs w/o carc. (N = 50)	Tot. %
	Gleason							
	3 + 3 (N = 14)	3 + 4 (N = 3)	4 + 3 (N = 11)	4 + 4 (N = 5)	4 + 5 (N = 4)			
Single chain of nuclei (TP1)	<u>52</u>	<u>7</u>	<u>65</u>	<u>6</u>	<u>2</u>	132	-	132 (20.4%)
Small round hole (TP2)	<u>12</u>	<u>2</u>	<u>24</u>	<u>11</u>	<u>8</u>	57	-	57 (8.8%)
High nuclear density (TP3)	<u>2</u>	<u>1</u>	<u>17</u>	<u>17</u>	<u>11</u>	48	-	48 (7.4%)
Larger nucleus with perinuclear halo (TP4)	<u>1</u>	<u>0</u>	<u>6</u>	<u>1</u>	<u>8</u>	16	-	16 (2.5%)
Undefined	<u>1</u>	<u>0</u>	<u>0</u>	<u>0</u>	<u>0</u>	1	-	1 (0.2%)
Single chain of nuclei (FP1)	5	0	0	1	0	6	29	35 (5.4%)
Small round hole (FP2)	4	0	0	0	0	4	35	39 (6.0%)
High nuclear density (FP3)	5	3	0	2	0	10	8	18 (2.8%)
Larger nucleus with perinuclear halo (FP4)	0	1	2	0	0	3	12	15 (2.3%)
Undefined	1	0	0	0	0	1	1	2 (0.3%)
Two-layered chain of nuclei (TN1)	13	2	11	11	6	43	29	72 (11.1%)
Areas of low nuclear density with eosinophilic background (TN2)	23	6	4	2	2	37	125	162 (25.1%)
Chain of nuclei with abundant slightly eosinophilic neighborhood (TN3)	5	0	1	1	2	9	30	39 (6.0%)
Undefined	4	1	0	2	0	7	3	10 (1.5%)

Out of 23 patterns, 14 can be detected by OSA, out of which 3 (denoted by an asterisk in the tables) are detected indirectly (distorted gland architecture, small uniform glands infiltrate, and regular lobular architecture with isolated glands or stroma predominant). 6 features are not recognized, 3 of which due to incompatible WSI resolution used for patches. The remaining three undetected patterns are blue mucin, eosinophilic amorphous secretions, and crystalloids, primarily represented by acellular areas.

Table 8

Morphological features of cancer prostatic tissue and their recognition by OSA. OoS – out of scale, SL – small lumina, HCD – high cellular density, SLE – single-layered epithelium, AA – acellular areas, HNH – hyperchromatic nuclei with halo. L – 20–40 × 10–5 mm FOV, M – 100–200 × 2–1 mm FOV, S – 400 × 0.5 mm FOV.

Feature type	Scale	Feature	Features (≈50 μm ø)	Found	xPatterns
Architectural	L	distorted gland architecture		No*	
		small uniform glands infiltrate in between normal glands	OoS	No*	
		poorly formed fused, cribriform or glomeruloid glands, high nuclear density in Gleason pattern 4	SL, HCD	Yes	small round hole, high nuclear density
		solid sheets, cords, medium or large nests with rosettes, comedo type necrosis	SL, HCD	Yes	small round hole, high nuclear density
		small caliber glands	SLE, SL	Yes	single chain of nuclei, small round hole
	M	crowded or compact glands clusters	HCD	Yes	high nuclear density
		blue mucin	AA	No	
		eosinophilic amorphous secretions	AA	No	
		crystalloids		No	
		rigid or sharp gland lumina, may have periglandular clefts	SLE, SL	Yes	single chain of nuclei, small round hole
Intraluminal		glands lack basal cells (single-layered epithelium in Gleason pattern 3)	SLE	Yes	single chain of nuclei
		infiltrative single cells in Gleason pattern 5	HNH	Yes	larger nucleus with perinuclear halo
					single chain of nuclei, small round hole
Cytoplasmic		cuboidal to low cylindrical cells with modest cytoplasm	SLE, SL	Yes	single chain of nuclei, small round hole
Nuclear		enlarged hyperchromatic nuclei	HNH	Yes	larger nucleus with perinuclear halo
	S	prominent enlarged nucleoli	OoS	N/A	
		often eosinophilic multiple nucleoli located in periphery	OoS	N/A	

histopathology, and potentially for identifying previously unrecognized morphological features related to histopathological diagnosis, prognosis, and prediction [90,91]. Finally, unraveling the large quantity of features within the network and exposing the key elements will help to promote trust in these and similar AI-based methods in pathology, enhancing the opportunities for incorporation into clinical use.

Table 9

Morphological features of non-cancer prostatic tissue and their recognition by OSA. OoS – out of scale, 2LE – two-layered epithelium, HPE – highly polarized epithelium. L – 20–40 × 10–5 mm FOV, M – 100–200 × 2–1 mm FOV, S – 400 × 0.5 mm FOV.

Feature type	Scale	Feature	Features (≈50 μm ø)	Found	xPatterns
Architectural	L	regular lobular architecture with isolated glands and relatively balanced epithelium/stroma ratio, or stroma predominant large glands	OoS	Yes*	areas of low nuclear density with eosinophilic background
	M	two-layered epithelium abundant cytoplasm in highly polarized cells	HPE, 2LE	Yes	two-layered chain of nuclei
Cytoplasmic			HPE	Yes	chain of nuclei with abundant slightly eosinophilic neighborhood
Nuclear	S	uniform smaller nuclei, nucleoli not visible	OoS	N/A	

Ethics approval

The project was approved by the Ethical Committee of Masaryk Memorial Cancer Institute, No. MOU 385 920.

Funding statement

This work has been supported by Czech Ministry of Health, (MMCI 00209805) and Czech Ministry of Education, Youth and Sports, (project BBMRI.cz, reg. no. LM2023033). The sponsors were not involved in the study design or execution.

Declaration of Competing Interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: Co-author, Petr HOLUB, is serving as an guest editor for the special issue of the journal.

Acknowledgments

Computational resources were provided by the e-INFRA CZ project (ID:90140), supported by the Ministry of Education, Youth and Sports of the Czech Republic.

Appendix A. Supporting information

Supplementary data associated with this article can be found in the online version at [doi:10.1016/j.nbt.2023.09.008](https://doi.org/10.1016/j.nbt.2023.09.008).

References

- [1] Holzinger A, Keiblinger K, Holub P, Zatloukal K, Müller H. AI for life: trends in artificial intelligence for biotechnology. *N Biotechnol* 2023;74:16–24. <https://doi.org/10.1016/j.nbt.2023.02.001>.
- [2] Evans AJ, Bauer TW, Bui MM, Cornish TC, Duncan H, Glassy EF, et al. US food and drug administration approval of whole slide imaging for primary diagnosis: a key milestone is reached and new questions are raised. *Arch Pathol Lab Med* 2018;142:1383–7. <https://doi.org/10.5858/arpa.2017-0496-CP>.

- [3] Stathonikos N, Nguyen TQ, Spoto CP, Verdaasdonk MAM, van Diest PJ. Being fully digital: perspective of a Dutch academic pathology laboratory. *Histopathology* 2019;75:621–35. <https://doi.org/10.1111/his.13953>.
- [4] Litjens GJS, Sánchez CI, Timofeeva N, Hermesen M, Nagtegaal ID, Kovacs I, et al. Deep learning as a tool for increased accuracy and efficiency of histopathological diagnosis. *Sci Rep* 2016;6:26286. <https://doi.org/10.1038/srep26286>.
- [5] Campanella G, Hanna MG, Geneslaw L, Mirafior AP, Silva VWK, Busam KJ, et al. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nat Med* 2019;25:1301–9. <https://doi.org/10.1038/s41591-019-0508-1>.
- [6] Esteve A, Chou K, Yeung S, Naik N, Madani A, Mottaghi A, et al. Deep learning-enabled medical computer vision. *NPJ Dig Med* 2021;4:5. <https://doi.org/10.1038/s41746-020-00376-2>.
- [7] Raciti P, Sue J, Ceballos R, Godrich R, Kunz J, Kapur S, et al. Novel artificial intelligence system increases the detection of prostate cancer in whole slide images of core needle biopsies. *Mod Pathol* 2020;33:2058–66. <https://doi.org/10.1038/s41379-020-0551-y>.
- [8] Ström P, Kartasalo K, Olsson H, Solorzano L, Delahunt B, Berney D, et al. Artificial intelligence for diagnosis and grading of prostate cancer in biopsies: a population-based, diagnostic study. *Lancet Oncol* 2020;21:222–32. [https://doi.org/10.1016/S1470-2045\(19\)30738-7](https://doi.org/10.1016/S1470-2045(19)30738-7).
- [9] Pantanowitz L, Quiroga-Garza GM, Bien L, Heled R, Laifienfeld D, Linhart C, et al. An artificial intelligence algorithm for prostate cancer diagnosis in whole slide images of core needle biopsies: a blinded clinical validation and deployment study. *Lancet Digit Health* 2020;2:e407–16. [https://doi.org/10.1016/S2589-7500\(20\)30159-X](https://doi.org/10.1016/S2589-7500(20)30159-X).
- [10] Abels E, Pantanowitz L, Aeffner F, Zarella MD, van der Laak J, Bui MM, et al. Computational pathology definitions, best practices, and recommendations for regulatory guidance: a white paper from the Digital Pathology Association. *J Pathol* 2019;249:286–94. <https://doi.org/10.1002/path.5331>.
- [11] van der Laak J, Litjens G, Ciompi F. Deep learning in histopathology: the path to the clinic. *Nat Med* 2021;27:775–84. <https://doi.org/10.1038/s41591-021-01343-4>.
- [12] Turkki R, Bychkov D, Lundin M, Isola J, Nordling S, Kovanen P, et al. Breast cancer outcome prediction with tumour tissue images and machine learning. *Breast Cancer Res Tr* 2019;177:41–52. <https://doi.org/10.1007/s10549-019-05281-1>.
- [13] Mobadersany P, Yousefi S, Amgad M, Gutman DA, Barnholtz-Sloan JS, Velázquez Vega JE, et al. Predicting cancer outcomes from histology and genomics using convolutional networks. *Proc Natl Acad Sci USA* 2018;115:E2970–9. <https://doi.org/10.1073/pnas.1717139115>.
- [14] Loh HW, Ooi CP, Seoni S, Barua PD, Molinari F, Acharya UR. Application of explainable artificial intelligence for healthcare: a systematic review of the last decade (2011–2022). *Comput Methods Prog Biomed* 2022;226:107161. <https://doi.org/10.1016/j.cmpb.2022.107161>.
- [15] Band SS, Yarahmadi A, Hsu C-C, Biyari M, Sookhak M, Ameri R, et al. Application of explainable artificial intelligence in medical health: a systematic review of interpretability methods. *Inform Med Unlocked* 2023;40:101286. <https://doi.org/10.1016/j.imu.2023.101286>.
- [16] Paner GP. *Acinar Adenocarcinoma*. In: Amin MB, Tickoo SK, editors. *Diagnostic Pathology: Genitourinary*. 3rd edition., Philadelphia, PA: Elsevier; 2022. p. 638–41.
- [17] Simonyan K., Zisserman A. Very deep convolutional networks for large-scale image recognition. 3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings 2015:1–14. <https://doi.org/10.48550/arXiv.1409.1556>.
- [18] Zeiler MD, Fergus R. Visualizing and understanding convolutional networks. In: Fleet D, Pajdla T, Schiele B, Tuytelaars T, editors. *Computer Vision – ECCV 2014*. Cham: Springer International Publishing; 2014. p. 818–33. https://doi.org/10.1007/978-3-319-10590-1_53.
- [19] Petsiuk V, Das A, Saenko K. Rise: Randomized input sampling for explanation of black-box models. *ArXiv Prepr ArXiv* 2018;180607421. <https://doi.org/10.48550/arXiv.1806.07421>.
- [20] Samek W, Binder A, Montavon G, Lapuschkin S, Müller K-R. Evaluating the visualization of what a deep neural network has learned. *IEEE Trans Neural Netw Learn Syst* 2016;28:2660–73. <https://doi.org/10.1109/TNNLS.2016.2599820>.
- [21] Ancona M, Ceolini E, Öztireli K, Gross M. Towards better understanding of gradient-based attribution methods for deep neural networks. *ArXiv Prepr ArXiv* 2017;171106104. <https://doi.org/10.48550/arXiv.1711.06104>.
- [22] Zhang Y, Khakzar A, Li Y, Farshad A, Kim ST, Navab N. Fine-grained neural network explanation by identifying input features with predictive information. *Adv Neural Inf Process Syst* 2021;34:20040–51.
- [23] Baehrens D, Schroeter T, Harmeling S, Kawanabe M, Hansen K, Müller K-R. How to explain individual classification decisions. *J Mach Learn Res* 2010;11:1803–31. <https://doi.org/10.5555/1756006.1859912>.
- [24] Springenberg JT, Dosovitskiy A, Brox T, Riedmiller M. Striving for simplicity: The all convolutional net. *ArXiv Prepr ArXiv* 2014;14126806. <https://doi.org/10.48550/arXiv.1412.6806>.
- [25] Montavon G, Lapuschkin S, Binder A, Samek W, Müller K-R. Explaining nonlinear classification decisions with deep Taylor decomposition. *Pattern Recognit* 2017;65: 211–22. <https://doi.org/10.1016/j.patcog.2016.11.008>.
- [26] Bach S, Binder A, Montavon G, Klauschen P, Müller K-R, Samek W. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS One* 2015;10:e0130140. <https://doi.org/10.1371/journal.pone.0130140>.
- [27] Sundararajan M., Taly A., Yan Q. Axiomatic attribution for deep networks. *International conference on machine learning*, PMLR; 2017. p. 3319–28. <https://doi.org/10.5555/3305890.3306024>.
- [28] Ilse M, Tomczak J, Welling M. Attention-based Deep Multiple Instance Learning. In: Dy J, Krause A, editors. *Proceedings of the 35th International Conference on Machine Learning*. 80. PMLR; 2018. p. 2127–36.
- [29] Roszkowiak L, Korzynska A, Siemion K, Zak J, Pijanowska D, Bosch R, et al. System for quantitative evaluation of DAB&H-stained breast cancer biopsy digital images (CHISEL). *Sci Rep* 2021;11:1–14. <https://doi.org/10.1038/s41598-021-88611-y>.
- [30] Van Rijthoven M, Balkenhol M, Silpa K, Van Der Laak J, Ciompi F. HookNet: multi-resolution convolutional neural networks for semantic segmentation in histopathology whole-slide images. *Med Image Anal* 2021;68:101890. <https://doi.org/10.1016/j.media.2020.101890>.
- [31] Pinckaers H, Bulten W, van der Laak J, Litjens G. Detection of prostate cancer in whole-slide images through end-to-end training with image-level labels. *IEEE Trans Med Imaging* 2021;40:1817–26. <https://doi.org/10.1109/TMI.2021.3066295>.
- [32] Ikromjanov K., Bhattacharjee S., Hwang Y.-B., Sumon R.I., Kim H.-C., Choi H.-K. Whole-Slide Image Analysis and Detection of Prostate Cancer using Vision Transformers. 2022 International Conference on Artificial Intelligence in Information and Communication (ICAIC); IEEE; 2022. p. 399–402. <https://doi.org/10.1109/ICAIC54071.2022.9722635>.
- [33] Chen T, Kornblith S, Swersky K, Norouzi M, Hinton GE. Big self-supervised models are strong semi-supervised learners. In: Larochelle H, Ranzato M, Hadsell R, Balcan MF, Lin H, editors. *Advances in Neural Information Processing Systems*. 33. Curran Associates, Inc.; 2020. p. 22243–55. <https://doi.org/10.5555/3495724.3497589>.
- [34] Zheng Y, Gindra RH, Green EJ, Burks EJ, Betke M, Beane JE, et al. A Graph-Transformer for Whole Slide Image Classification. *IEEE Trans Med Imaging* 2022; 41:3003–15. <https://doi.org/10.1109/TMI.2022.3176598>.
- [35] Thandiackal K., Chen B., Pati P., Jaume G., Williamson D.F.K., Gabrani M., et al. Differentiable Zooming for Multiple Instance Learning on Whole-Slide Images 2022. <https://doi.org/10.48550/ARXIV.2204.12454>.
- [36] Ciga O, Xu T, Martel AL. Self supervised contrastive learning for digital histopathology. *Mach Learn Appl* 2022;7:100198. <https://doi.org/10.1016/j.mlwa.2021.100198>.
- [37] Chhipa P.C., Upadhyay R., Pihlgren G.G., Saini R., Uchida S., Liwicki M. Magnification Prior: A Self-Supervised Method for Learning Representations on Breast Cancer Histopathological Images. *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2023. p. 2717–27. <https://doi.org/10.1109/WACV56688.2023.00274>.
- [38] Chakraborty S, Gupta R, Ma K, Govind D, Sarder P, Choi W-T, et al. Predicting the visual attention of pathologists evaluating whole slide images of cancer. *International Workshop on Medical Optical Imaging and Virtual Microscopy Image Analysis*. Springer; 2022. p. 11–21. https://doi.org/10.1007/978-3-031-16961-8_2.
- [39] Xiang J, Wang X, Wang X, Zhang J, Yang S, Yang W, et al. Automatic diagnosis and grading of Prostate Cancer with weakly supervised learning on whole slide images. *Comput Biol Med* 2022;106340. <https://doi.org/10.1016/j.combiomed.2022.106340>.
- [40] Zhou X, Tang C, Huang P, Mercaldo F, Santone A, Shao Y. LPCANet: Classification of Laryngeal Cancer Histopathological Images Using a CNN with Position Attention and Channel Attention Mechanisms. *Interdiscip Sci: Comput Life Sci* 2021;13: 666–82. <https://doi.org/10.1007/s12539-021-00452-5>.
- [41] Celik Y, Talo M, Yildirim O, Karabatak M, Acharya UR. Automated invasive ductal carcinoma detection based using deep transfer learning with whole-slide images. *Pattern Recognit Lett* 2020;133:232–9. <https://doi.org/10.1016/j.patrec.2020.03.011>.
- [42] Duran-Lopez L, Dominguez-Morales JP, Conde-Martin AF, Vicente-Diaz S, Linares-Barranco A. PROMETEO: a CNN-based computer-aided diagnosis system for WSI prostate cancer detection. *IEEE Access* 2020;8:128613–28. <https://doi.org/10.1109/ACCESS.2020.3008868>.
- [43] Chen C-L, Chen C-C, Yu W-H, Chen S-H, Chang Y-C, Hsu T-I, et al. An annotation-free whole-slide training approach to pathological classification of lung cancer types using deep learning. *Nat Commun* 2021;12:1–13. <https://doi.org/10.1038/s41467-021-21467-y>.
- [44] Lagree A, Shiner A, Alera MA, Fleschner L, Law E, Law B, et al. Assessment of digital pathology imaging biomarkers associated with breast cancer histologic grade. *Curr Oncol* 2021;28:4298–316. <https://doi.org/10.3390/curroncol28060366>.
- [45] Singhal N, Soni S, Bonthu S, Chattopadhyay N, Samanta P, Joshi U, et al. A deep learning system for prostate cancer diagnosis and grading in whole slide images of core needle biopsies. *Sci Rep* 2022;12:1–11. <https://doi.org/10.1038/s41598-022-07217-0>.
- [46] Ehteshami Bejnordi B, Litjens G, Timofeeva N, Otte-Holler I, Homeyer A, Karssemeijer N, et al. Stain Specific Standardization of Whole-Slide Histopathological Images. *IEEE Trans Med Imaging* 2016;35:404–15. <https://doi.org/10.1109/TMI.2015.2476509>.
- [47] Roy S, Kumar Jain A, Lal S, Kini J. A study about color normalization methods for histopathology images. *Micron* 2018;114:42–61. <https://doi.org/10.1016/j.micron.2018.07.005>.
- [48] Kang H, Luo D, Feng W, Zeng S, Quan T, Hu J, et al. StainNet: a fast and robust stain normalization network. *Front Med* 2021;8:746307. <https://doi.org/10.3389/fmed.2021.746307>.
- [49] Michiellini N, Caputo A, Scotto M, Mogetta A, Pennisi OAM, Molinari F, et al. Stain normalization in digital pathology: clinical multi-center evaluation of image

- quality. *J Pathol Inform* 2022;13:100145. <https://doi.org/10.1016/j.jpi.2022.100145>.
- [50] Zhao B, Han C, Pan X, Lin J, Yi Z, Liang C, et al. RestainNet: a self-supervised digital re-stainer for stain normalization. *Comput Electr Eng* 2022;103:108304. <https://doi.org/10.1016/j.compeleceng.2022.108304>.
- [51] Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, et al. An image is worth 16x16 words: transformers for image recognition at scale. *ArXiv Prepr ArXiv* 2020;201011929.
- [52] Erion G, Janizek JD, Sturmels P, Lundberg SM, Lee S-I. Improving performance of deep learning models with axiomatic attribution priors and expected gradients. *Nat Mach Intell* 2021;3:620–31. <https://doi.org/10.1038/s42256-021-00343-w>.
- [53] Selvaraju R.R., Cogswell M., Das A., Vedantam R., Parikh D., Batra D. Grad-CAM: Visual Explanations From Deep Networks via Gradient-Based Localization. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017. <https://doi.org/10.1109/ICCV.2017.74>.
- [54] Smilkov D, Thorat N, Kim B, Viégas FB, Wattenberg M. SmoothGrad: removing noise by adding noise. *CoRR* 2017. <https://doi.org/10.48550/arXiv.1706.03825>.
- [55] Lundberg SM, Lee S-I. A unified approach to interpreting model predictions. In: Guyon I, Luxburg UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, et al., editors. *Advances in Neural Information Processing Systems*, 30. Curran Associates, Inc.; 2017. <https://doi.org/10.5555/3295222.3295230>.
- [56] Frye C, Rowat C, Feige I. Asymmetric Shapley values: incorporating causal knowledge into model-agnostic explainability. *Adv Neural Inf Process Syst* 2020; 33:1229–39. <https://doi.org/10.5555/3495724.3495828>.
- [57] Wang J., Wiens J., Lundberg S. Shapley flow: A graph-based approach to interpreting model predictions. *International Conference on Artificial Intelligence and Statistics*, PMLR; 2021, p. 721–9. <https://doi.org/10.48550/arXiv.2010.14592>.
- [58] Biecek P. DALEX: explainers for complex predictive models in R. *J Mach Learn Res* 2018;19:3245–9. <https://doi.org/10.5555/3291125.3309646>.
- [59] Ribeiro M.T., Singh S., Guestrin C. “Why Should I Trust You?”: Explaining the Predictions of Any Classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, New York, NY, USA: Association for Computing Machinery; 2016, p. 1135–44. <https://doi.org/10.1145/2939672.2939778>.
- [60] Schwab P, Karlen W. CXPlain: causal explanations for model interpretation under uncertainty. In: Wallach H, Larochelle H, Beygelzimer A, Alché-Buc F, d', Fox E, Garnett R, editors. *Advances in Neural Information Processing Systems*, 32. Curran Associates, Inc.; 2019. <https://doi.org/10.5555/3454287.3455204>.
- [61] Zintgraf LM, Cohen TS, Adel T, Welling M. Visualizing deep neural network decisions: prediction difference analysis. *CoRR* 2017. <https://doi.org/10.48550/arXiv.1702.04595>.
- [62] Ribeiro M.T., Singh S., Guestrin C. Anchors: High-precision model-agnostic explanations. *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, 2018. <https://doi.org/10.1609/aaai.v32i1.11491>.
- [63] Holzinger A, Saranti A, Molnar C, Biecek P, Samek W. Explainable AI methods-a brief overview. *International Workshop on Extending Explainable AI Beyond Deep Models and Classifiers*. Springer; 2022. p. 13–38. https://doi.org/10.1007/978-3-031-04083-2_2.
- [64] Pocevičiūtė M, Eilertsen G, Lundström C. Survey of XAI in digital pathology. *Artificial Intelligence and Machine Learning for Digital Pathology*. Springer; 2020. p. 56–88.
- [65] Kim B., Wattenberg M., Gilmer J., Cai C., Wexler J., Viegas F., et al. Interpretability Beyond Feature Attribution: Quantitative Testing with Concept Activation Vectors (TCAV). In: {C}Dy J., Krause A., {C}, editors. *Proceedings of the 35th International Conference on Machine Learning*, vol. 80, PMLR; 2018, p. 2668–2677. <https://doi.org/10.48550/arXiv.1711.11279>.
- [66] Shrikumar A., Greenside P., Kundaje A. Learning important features through propagating activation differences. *International conference on machine learning*, PMLR; 2017, p. 3145–53. <https://doi.org/10.5555/3305890.3306006>.
- [67] Agarwal R, Melnick L, Frosst N, Zhang X, Lengerich B, Caruana R, et al. Neural additive models: interpretable machine learning with neural nets. In: Ranzato M, Beygelzimer A, Dauphin Y, Liang PS, Vaughan JW, editors. *Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc.; 2021. p. 4699–711. <https://doi.org/10.48550/arXiv.2004.13912>.
- [68] Krajina A, Brcic M, Kovac M, Sarcevic A. Explainable artificial intelligence: an updated perspective. *Review* 2022. <https://doi.org/10.23919/MIPRO55190.2022.9803681>.
- [69] Schnake T, Eberle O, Lederer J, Nakajima S, Schütt KT, Müller K-R, et al. XAI for graphs: explaining graph neural network predictions by identifying relevant walks. *CoRR* 2020. <https://doi.org/10.48550/arXiv.2006.03589>.
- [70] Huang Q, Yamada M, Tian Y, Singh D, Yin D, Chang Y. GraphLIME: local interpretable model explanations for graph neural networks. *CoRR* 2020. abs/2001.06216.
- [71] Zhang Y, Defazio D, Ramesh A. RelEx: A Model-Agnostic Relational Model Explainer. New York, NY, USA: Association for Computing Machinery; 2021. p. 1042–9. <https://doi.org/10.1145/3461702.3462562>.
- [72] Ying Z, Bourgeois D, You J, Zitnik M, Leskovec J. GNNExplainer: generating explanations for graph neural networks. In: Wallach H, Larochelle H, Beygelzimer A, Alché-Buc F d', Fox E, Garnett R, editors. *Advances in Neural Information Processing Systems*, vol. 32. Curran Associates, Inc.; 2019. <https://doi.org/10.5555/3454287.3455116>.
- [73] Pfeiffer B, Saranti A, Holzinger A. GNN-SubNet: disease subnetwork detection with explainable graph neural networks. *Bioinformatics* 2022;38:ii120–6. <https://doi.org/10.1093/bioinformatics/btac478>.
- [74] Yuan H., Tang J., Hu X., Ji S. XGNN: Towards Model-Level Explanations of Graph Neural Networks. *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, New York, NY, USA: Association for Computing Machinery; 2020, p. 430–8. <https://doi.org/10.1145/3394486.3403085>.
- [75] Yuan H., Yu H., Wang J., Li K., Ji S. On Explainability of Graph Neural Networks via Subgraph Explorations. In: {C}Meila M., Zhang T., {C}, editors. *Proceedings of the 38th International Conference on Machine Learning*, vol. 139, PMLR; 2021, p. 12241–12252. <https://doi.org/10.48550/arXiv.2102.05152>.
- [76] Dai E., Wang S. Towards self-explainable graph neural network. *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, p. 302–11. <https://doi.org/10.1145/3459637.3482306>.
- [77] Zhang Z., Liu Q., Wang H., Lu C., Lee C. Protgnn: Towards self-explaining graph neural networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, 2022, p. 9127–9135. <https://doi.org/10.48550/arXiv.2112.00911>.
- [78] Zhou J, Gandomi AH, Chen F, Holzinger A. Evaluating the quality of machine learning explanations: a survey on methods and metrics. *Electronics* 2021;10. <https://doi.org/10.3390/electronics10050593>.
- [79] Cabitza F, Campagner A, Malgieri G, Natali C, Schneeberger D, Stoeger K, et al. Quod erat demonstrandum? - towards a typology of the concept of explanation for the design of explainable AI. *Expert Syst Appl* 2023;213:118888. <https://doi.org/10.1016/j.eswa.2022.118888>.
- [80] Evans T, Retzlaff CO, Geißler C, Kargl M, Plass M, Müller H, et al. The explainability paradox: challenges for xAI in digital pathology. *Future Gener Comput Syst* 2022;133:281–96. <https://doi.org/10.1016/j.future.2022.03.009>.
- [81] Linardatos P, Papastefanopoulos V, Kotsiantis S. Explainable ai: a review of machine learning interpretability methods. *Entropy* 2020;23:18. <https://doi.org/10.3390/e23010018>.
- [82] Patrício C, Neves JC, Teixeira LF. Explainable deep learning methods in medical diagnosis: a survey. *ArXiv Prepr ArXiv* 2022;220504766.
- [83] Hooker S, Erhan D, Kindermans P-J, Kim B. A benchmark for interpretability methods in deep neural networks. *Adv Neural Inf Process Syst* 2019;32.
- [84] Jung M, Jin MS, Kim C, Lee C, Nikas IP, Park JH, et al. Artificial intelligence system shows performance at the level of uropathologists for the detection and grading of prostate cancer in core needle biopsy: an independent external validation study. *2022 35:10 Mod Pathol* 2022;35:1449–57. <https://doi.org/10.1038/s41379-022-01077-9>.
- [85] Litjens G. ASAP - Automated Slide Analysis Platform 2017. <https://computation.alpathologygroup.github.io/ASAP/>.
- [86] Goode A, Gilbert B, Harkes J, Jukic D, Satyanarayanan M. OpenSlide: a vendor-neutral software foundation for digital pathology. *J Pathol Inf* 2013;4:27. <https://doi.org/10.4103/2153-3539.119005>.
- [87] Liu Y., Gadepalli K., Norouzi M., Dahl G.E., Kohlberger T., Boyko A., et al. Detecting Cancer Metastases on Gigapixel Pathology Images. *ArXiv* 2017.
- [88] Alber M, Lapuschkin S, Seegerer P, Hägele M, Schütt KT, Montavon G, et al. iNNvestigate neural networks! *J Mach Learn Res* 2019;20:1–8.
- [89] CAMELYON16 - Grand Challenge. Grand-ChallengeOrg n.d. <https://camelyon16.grand-challenge.org/> (accessed August 10, 2023).
- [90] Eckardt J-N, Middeke JM, Riechert S, Schmittmann T, Sulaiman AS, Kramer M, et al. Deep learning detects acute myeloid leukemia and predicts NPM1 mutation status from bone marrow smears. *Leukemia* 2022;36:111–8. <https://doi.org/10.1038/s41375-021-01408-w>.
- [91] Govindarajan S, Swaminathan R. Differentiation of COVID-19 conditions in planar chest radiographs using optimized convolutional neural networks. *Appl Intel* 2021; 51:2764–75. <https://doi.org/10.1007/s10489-020-01941-8>.