RESEARCH



Equipping computational pathology systems with artifact processing pipelines: a showcase for computation and performance trade-offs

Neel Kanwal^{1*†}, Farbod Khoraminia^{2†}, Umay Kiraz^{3,4†}, Andrés Mosquera-Zamudio^{5†}, Carlos Monteagudo⁵, Emiel A. M. Janssen^{3,4}, Tahlita C. M. Zuiverloon², Chunming Rong¹ and Kjersti Engan^{1*}

Abstract

Background Histopathology is a gold standard for cancer diagnosis. It involves extracting tissue specimens from suspicious areas to prepare a glass slide for a microscopic examination. However, histological tissue processing procedures result in the introduction of artifacts, which are ultimately transferred to the digitized version of glass slides, known as whole slide images (WSIs). Artifacts are diagnostically irrelevant areas and may result in wrong predictions from deep learning (DL) algorithms. Therefore, detecting and excluding artifacts in the computational pathology (CPATH) system is essential for reliable automated diagnosis.

Methods In this paper, we propose a mixture of experts (MoE) scheme for detecting five notable artifacts, including damaged tissue, blur, folded tissue, air bubbles, and histologically irrelevant blood from WSIs. First, we train independent binary DL models as experts to capture particular artifact morphology. Then, we ensemble their predictions using a fusion mechanism. We apply probabilistic thresholding over the final probability distribution to improve the sensitivity of the MoE. We developed four DL pipelines to evaluate computational and performance trade-offs. These include two MoEs and two multiclass models of state-of-the-art deep convolutional neural networks (DCNNs) and vision transformers (ViTs). These DL pipelines are quantitatively and qualitatively evaluated on external and outof-distribution (OoD) data to assess generalizability and robustness for artifact detection application.

Results We extensively evaluated the proposed MoE and multiclass models. DCNNs-based MoE and ViTs-based MoE schemes outperformed simpler multiclass models and were tested on datasets from different hospitals and cancer types, where MoE using (MobileNet) DCNNs yielded the best results. The proposed MoE yields 86.15 % F1 and 97.93% sensitivity scores on unseen data, retaining less computational cost for inference than MoE using ViTs. This best performance of MoEs comes with relatively higher computational trade-offs than multiclass models. Furthermore, we apply post-processing to create an artifact segmentation mask, a potential artifact-free Rol map, a quality report,

¹Neel Kanwal, Farbod Khoraminia, Umay Kiraz and Andrés Mosquera-Zamudio contributed equally to this work.

*Correspondence: Neel Kanwal neel.kanwal@uis.no Kjersti Engan kjersti.engan@uis.no Full list of author information is available at the end of the article



© The Author(s) 2024. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit http://creativecommons.org/licenses/by/4.0/.

and an artifact-refined WSI for further computational analysis. During the qualitative evaluation, field experts assessed the predictive performance of MoEs over OoD WSIs. They rated artifact detection and artifact-free area preservation, where the highest agreement translated to a Cohen Kappa of 0.82, indicating substantial agreement for the overall diagnostic usability of the DCNN-based MoE scheme.

Conclusions The proposed artifact detection pipeline will not only ensure reliable CPATH predictions but may also provide quality control. In this work, the best-performing pipeline for artifact detection is MoE with DCNNs. Our detailed experiments show that there is always a trade-off between performance and computational complexity, and no straightforward DL solution equally suits all types of data and applications. The code and HistoArtifacts dataset can be found online at Github and Zenodo, respectively.

Keywords Computational pathology, Deep learning, Histological artifacts, Mixture of experts, Vision transformer, Whole slide images

Introduction

Cancer develops in organs when genetic mutations in normal cells trigger their transformation into tumor cells. This transformation may be triggered by frequent exposure to carcinogens, a class of substances (chemical, biological, or physical), or several other factors that have the potential to cause cancer [1]. Diagnosing cancer accurately and efficiently is critical for medical treatment and a reduced mortality rate, given its status as one of the deadliest diseases worldwide, with a projected estimate of 29 million deaths by 2040 [2, 3]. Histopathology is considered a gold standard for identifying cancerous cells, which involves examining tissue samples under a microscope using a histological glass slide [4]. However, this manual inspection and laboratory procedure is not without its pitfalls, as it is labor-intensive, subjective, and can be affected by inter- and intra-observer variability [5, 6]. Furthermore, the projected rise in cancer cases and the shortage of pathologists are significant issues that may lead to delayed diagnosis and treatment, resulting in a severe impact on clinical decision-making [7]. Therefore, streamlining the traditional diagnostic process through digitization and automation can provide timely diagnosis, improved treatment decisions, and efficacy [3]. Digital pathology (DP) has the potential to overcome these challenges by providing rapid diagnosis and smooth sharing of secondary opinions [8]. In fact, in the last decade, there has been a five-fold growth in DP research and development [9, 10]. This increase in the adoption of DP in clinical practice enables computation over the digitized version of histological slides, commonly called whole slide images (WSIs).

Computational pathology (CPATH) systems have the potential to unfold information embedded in WSIs by automated systems based on AI and image processing [10-12]. The seamless integration of CPATH with DP can enhance diagnostic or prognostic methodologies and save pathologists' time [6, 13]. However, artifacts that appear during the histological slide preparation are

ultimately transferred to the WSIs [14–16]. Artifacts are diagnostically irrelevant areas, and pathologists usually ignore these areas during manual inspection, but unfortunately, the presence of histological artifacts can hamper the performance of CPATH systems during automated diagnosis [10, 17]. Therefore, it is essential to equip the CPATH system with an *artifact detection pipeline* to exclude artifacts and ensure the flow of histologically relevant tissue for diagnostic or prognostic algorithms, as illustrated in Fig. 1. Thus, a CPATH system with artifact processing capacity will not only increase the likelihood of reliable and accurate predictions but also provide quality control (QC) for laboratory procedures, identifying weaknesses during the histotechnical stages (see review [10]) in acquiring WSIs.

In recent years, deep learning (DL) approaches have garnered more attention from the CPATH community due to their ability to extract hidden patterns in histological data [18–21]. Popular DL architectures such as deep convolution neural networks (DCNNs) and vision transformers (ViTs) have widely been used as state-ofthe-art (SOTA) to distinguish tissue patterns for different cancer types and perform image classification and segmentation tasks [16, 19, 22]. Some researches [23, 24] demonstrate that DCNNs perform better on small datasets, thanks to the inductive bias, which helps them to learn spatial relevance effectively. While other works [25-27] argue in favor of ViTs, showing that they are highly robust, attend to overall structural information, and are less biased towards textures. Nevertheless, both DL architectures may suffer from overfitting, poor generalization, and reproducibility issues, leading to overconfident predictions on new (external) data. To address these problems, ensembles of DL models (a.k.a. deep ensembles) have been used to overcome the weakness of an individual model [28-30]. Ensemble methods combine the prediction of independent models using averaging or majority voting. A mixture of experts (MoE) is an extended method that trains DL for a sub-task and



Fig. 1 An overview of computational pathology (CPATH) system equipped with artifact processing pipeline. Whole slide images (WSIs) are split into small sub-images (patches) to make them computationally tractable for deep learning (DL) models. These patches are fed to a mixture of experts (MoE) or multiclass models composed of state-of-the-art DL architectures to perform different CPATH classification tasks. Only patches with histological relevance can flow further for the downstream tasks. Finally, predictions are post-processed to produce different outcomes, such as a segmentation map, artifact report for quality control, region of interest mask, and artifact-free WSI for the diagnostic or prognostic algorithm to make a final clinical prediction

then combines the predictions dynamically to obtain a nuanced prediction. In short, the MoE approach consists of multiple DCNNs or ViTs, experts on each subclass, to achieve improved results. MoEs benefit in terms of reproducibility by reducing the variance of predictions but augmenting computational expense [31]. In contrast, the multiclass approach can be computationally efficient but does not involve the strength of multiple models, which are adaptive for looking into different aspects of data. Based on these arguments, the choice between DL approaches depends on application requirements. This raises a fundamental question: *how to build an effective artifact detection DL approach for CPATH systems with suitable trade-offs between computational complexity and performance*?

An effective DL approach for artifact detection applications (our case) might be created using MoEs, one DL model for each artifact class, or multiclass models with multiple output classes. In this paper, we propose the MoE-based DL approach, which uses a fusion mechanism to integrate predictions from experts and apply probabilistic thresholding to improve the sensitivity. We establish several DL pipelines using the MoE and multiclass models for detecting notable artifacts (i.e., damaged tissue, blur, folded tissue, air bubbles, and diagnostically irrelevant blood) from histological WSIs (see Fig. 1). Our DL pipelines produce four outcomes for the input WSI: i) Artifact segmentation map; ii) Artifact report for QC using six classes (five artifacts and artifact-free area); iii) Artifact-free mask with potential regions of interest (RoIs) with diagnostic relevance; and iv) Artifact-refined WSI for the diagnostic algorithm.

Our contributions to this work are summarized below:

- We develop four DL models (referred to as DL pipeline throughout the paper), with SOTA DCNNs (MobileNet [32]) and ViTs (ViT-tiny [33]), using MoE and a multiclass approach.
- We evaluate the computational complexity of the pipelines and systematically choose a learned probability threshold for maximizing the sensitivity of DL models in external validation.
- We conduct a qualitative and quantitative evaluation over external data (from different cancer types) and assess the efficiency of the proposed MoE scheme for detecting artifacts and QC.

The paper is structured as follows: "Related work" section presents recent studies involving DL approaches for computational pathology and related work for detecting artifacts. "Data materials" section provides data material descriptions. "Proposed method" section explains preprocessing for creating datasets, the proposed method, post-processing, evaluation metrics, and implementation details. "Experimental results and discussion" section discusses results for performance and computational complexity. "Conclusion" section concludes this work. Finally, "Limitations and future work" section discusses limitations and future directions for a smooth integration of artifact processing pipelines in CPATH systems.

Related work

Deep learning for computational pathology

Deep learning (DL) approaches have gained popularity in the CPATH community for different tasks [21, 34–36]. In recent years, several works [20, 37–40] have used popular DL architectures for diagnosis and prognostic algorithms. FDA-approved PAIGE [41] is an example of such a DLbased algorithm for prostate cancer. These works can be roughly divided into two branches, such as DCNN-based (MobileNet [32], DenseNet [42], ResNet [43], or GoogleNet [44], etc.) or ViT-based (ViT-Tiny([33], DINO [45], or SwinTransformer [46] etc.) approaches.

In the first branch, Srinidhi et al. [47] comprehensively reviewed different DL approaches for developing disease-specific classification algorithms using histological images. Riasatian et al. [48] applied transfer learning over DCNNs to classify various tumor types and accomplished remarkable results using three public histopathology datasets. Talo [49] demonstrated that pre-trained ResNet [43] and DenseNet [42] achieved better accuracy than traditional methods in the literature for classifying grayscale and color histopathological images. Similarly, Wang et al. [50] proposed a DCNN-based method based on GoogleNet [44] to locate tumors in breast and colon images using complex example-guided training for WSI analysis. Among other DCNN works, Meng et al. [29] compared several architectures for classification and segmentation problems on a cervical histopathology dataset. Their approach found the best results for precancerous lesions using ResNet-101 [43]. For the same task, MobileNet [32] was the fastest. Wang et al. [51] performed multi-class breast cancer classification in their two-stage dependency-based framework. A MobileNet [32] was used as a backbone to extract the features in the first stage. Then, the MobileNet [32] backbone was modified to perform sub-type classification for benign and malignant categories. Gandomkar et al. [38] deployed ResNet [43] for classifying breast histology images into benign or malignant and then identified them among several sub-types using a meta-classifier based on a decision tree.

Works in the second branch used ViTs, which have emerged as new SOTA, leveraging attention mechanisms to improve shape understanding and generalizablity [26] [27]. Stegmüller et al. [40] developed ViT-based ScoreNet for breast cancer classification. Their approach attended to some regions in the WSI for faster processing based on image semantics. Wessel et al. [39] used DINO [45] for predicting overall and disease-specific survival in renal cell carcinoma. Zidan et al. [46] introduced a ViT-based cascaded architecture for segmenting glands, nuclei, and stroma in colorectal cancer. Gao et al. [52] proposed instance-based ViT to capture global and local features for subtyping papillary renal carcinoma, achieving better performance over selected RoIs.

Unsurprisingly, in both branches, most of these DL algorithms were trained and tested on manually annotated clean data (with diagnostic relevance) and overlooked the impact of potential noise (histological artifacts) during the inference stage or unseen scenarios.

Schomig et al. [53], in their stress-testing study, showed that the accuracy of the prostate cancer DL algorithm was negatively affected by the presence of artifacts and resulted in more false positives. Even the presence of artifacts in the training data may result in poor learning by DL models, as they add irrelevant features to the data [10, 54]. Wright et al. [17] demonstrated that removing images with artifacts improved the accuracy of DL models by a significant margin. Laleh et al. [55] emphasized the need for robustness of DL-based CPATH systems against artifacts for their widespread clinical adaptability. Artifact processing pipeline that can detect, extract, and eliminate non-relevant patches from WSIs before running a diagnostic algorithm would avoid any detrimental effect on downstream image analysis [11, 17, 56]. Therefore, it is essential to equip CPATH systems with artifact detection ability, which is also the focus of this work, to obtain reliable predictions [17, 57, 58].

Detection of histological artifacts

Most researches focus on reducing color variations and image augmentations during the preprocessing phase in CPATH literature [59, 60]. Detection of artifacts is often an underrepresented aspect of WSI pre-processing [10]. Compared to color normalization techniques, there remains a scarcity of research detecting notable artifacts before feeding histologically relevant data to the diagnostic algorithms. While some works [17, 61–63] have relied on quickly identifying faulty WSIs by doing QC at low magnification. Avanki et al. [64] proposed a quality estimation method by combining blurriness, contrast, brightness, etc., to accept or discard WSI based on a reference. Bahlmann et al. [63] exploited texture features and stain absorption to separate diagnostically relevant and irrelevant regions. However, artifacts appearing in diagnostically relevant areas are likely to be missed. Apart from their limitations with lower magnification, they were validated based on specific staining and tissue types. Therefore, artifact detection methods need to be extended to higher magnification. Moreover, artifact detection methods that can identify specific artifacts are desirable for QC, as some artifacts, like a blur, can be avoided by re-scanning glass slides or de-blurring techniques.

Earlier works for artifact detection relied on traditional image processing and color-space transformation approaches. Gao et al. [65] detected blurry areas using 44 handcrafted (local statistics, brightness, etc) features. Hashimoto et al. [66] combined image sharpness and noise information to create a regression model for out-offocus detection. For folded tissue detection, Palokangas et al. [67] transformed red, green, and blue (RGB) images to hue, saturation, and intensity (HSI) to apply k-means clustering over the different saturation and intensity values. Bautista and Yagi [68] detected folds using RGB shift with fixed thresholding on luminance and saturation values to enhance color structure in thick (folded) areas. Kothari et al. [57] introduced a rank-sum approach that used connectivity descriptors and image features to discard folded tissues. Their approach used two adaptive thresholds on saturation and intensity ranges. Chadaj et al. [69] separated uninformative blood (hemorrhage) from blood vessels using cyan, magenta, yellow, and black (CYMK) color space and morphology. Mercan et al. [70] proposed a k-means method to classify diagnostically relevant vs. non-relevant patches using local binary patterns extracted from stains and L*a*b histograms. A detailed review of other artifact detection works can be found in Kanwal et al. [10]. Since color-based approaches can heavily underperform when exposed to data from different cohorts with stain variation, data-driven DL approaches are needed to resolve the challenges.

Among recent works using DL-based approaches, Albuquerque et al. [71] compared several DCNNs for detecting out-of-focus areas in their ordinal classification problem. Kohlberger et al. [72] proposed ConvFocus to quantify and localize blurry areas in WSI. Wetteland et al. [73, 74] proposed a segmentation model to find blood and damaged tissue in bladder cancer WSIs. Clymer et al. [75] developed a two-stage method to detect blood at low resolution using RetinaNet and later Xception CNN for subsequent classification. Babie et al. [76] used SOTA DCNNs with SVM, KNN, and decision tree classifiers to separate folded tissues from normal tissue in a binary fashion. HistoQC [77] was proposed to perform a content-based evaluation for detecting outliers in a cohort of WSIs using a combination of image metrics and supervised classifiers. However, it did not remove or detect each artifact in smaller regions on WSIs. Kanwal et al. [78] used several DCNNs to assess the impact of color normalization over blood and damaged tissue detection. In another work [16], they trained ViT-Tiny [33] for air bubble detection using knowledge distillation. All these works relied on training a single network to classify one or two artifacts against an artifact-free class. It is a well-known problem that DL models suffer from poor generalization, robustness, and overconfident predictions over out-of-distribution (OoD) data [79-81]. Thus, the high variance in the prediction of DL models needs to be addressed, especially when deployed in a critical application. A prominent DL technique, "deep ensembles", resolves these problems by training several baseline DL architectures and combining the resultant predictions to increase accuracy and OoD performance [5]. However, the success of the ensemble method relies on several factors, such as how baseline models are trained and integrated. The most widely used ensemble techniques include averaging and majority voting [31]. It is worth noting that a simple aggregation using averaging methods or majority voting is not a smart choice and is very sensitive to biased baseline models [31]. A mixture of experts (MoE) may address this shortcoming by combining base learners, who are experts on detecting particular artifact morphology. Unlike deep ensemble, where all models are trained on the same data, in MoE, each DL model is trained for a sub-task to master specific aspects of the data, resulting in improved robustness. This is the first work to provide a comprehensive DL-based artifact processing pipeline that takes the entire WSI, preprocess, infer, and post-process in an end-to-end fashion for artifact detection and QC applications.

Data materials

This section details the histological data used for training and validating DL models. The following in-house (private) datasets are used for the experiments.

Training and development data

We obtained 55 WSIs of bladder cancer resection biopsies from the Erasmus Medical Centre (EMC) in Rotterdam, The Netherlands. The glass slides were formalin-fixed, stained with Hematoxylin (purple) and Eosin (pink) (H&E) dyes, scanned with a Hamamatsu Nanozoomer 2.0HT at 40× and saved in *ndpi* format with a pixel size of 0.227 μ m × 0.227 μ m. These WSIs were properly anonymized to preserve patient privacy, and all ethical requirements were followed before the dataset was created. NK received training for the task and manually annotated five artifacts (blurry areas, folded tissues, blood (hemorrhage), air bubbles, and damaged tissue). The rest of the tissue was marked as an artifact-free region. Note that not all WSIs contained five artifacts present and they were not extensively labeled as distinct tissue types since this histological data is not used for any task other than artifact detection. However, each WSI had at least one annotation for RoI or the artifact region (i.e., blur, fold, etc). Later sections refer to this cohort as EMC_{dev} . A detailed description of the prepared dataset and its availability is mentioned in "Pre-processing" section.

External validation data

Since these external cohorts are not involved in training and development, they can be considered OoD data with different tissue types and staining. We have used the following cohorts for inference only to validate the generalizability and robustness of the proposed methods.

EMC cohort:

This cohort is a collection of bladder cancer WSIs from a multi-center cohort provided by Erasmus MC, Rotterdam, The Netherlands. These WSIs with *MRXS* format were prepared with H&E staining and scanned with a 3DHistech P100 scanner at 80× magnification. A few WSIs were selected based on the presence of artifacts to test their generalization ability. FK manually annotated these WSIs for five artifacts, in a similar fashion as mentioned earlier. We have used a 40× magnification level for inference as the models are trained at a similar level. We will refer to this dataset as *EMC_{inf}*, and it is a different cohort than the above-mentioned *EMC_{dev}*.

SUH cohort:

This cohort is a private breast cancer cohort of 258 surgical specimens. It contains H&E WSIs prepared from surgical specimens and collected by the Stavanger University Hospital (SUH) in Norway between 1978 and 2004. The WSIs are in *ndpi* format and scanned using the Hamamatsu NanoZoomer S60 at 40× magnification. An expert pathologist (UK) selected and manually annotated a few WSIs based on the severity of the presence of these artifacts. Only five artifacts were carefully annotated, and the rest were marked as artifact-free regions. We have used these WSIs to test DL pipelines over cancer types that differ from the ones they are trained on. We will refer to this dataset as *SUH_{inf}*.

INCLIVA cohort:

This cohort was prepared by the Department of Anatomical Pathology of the Hospital Clínico Universitario de Valencia, Spain. It is a collection between 1988 and 2020. The glass slides were prepared from skin cancer biopsies and were scanned with Roche's Ventana iScan HT at $40 \times$ magnification. WSIs were saved in *tiff* format. An expert dermatopathologist (AM) selected and annotated a few WSIs with artifacts from this cohort to validate the proposed pipeline over the external cohort. We will refer to this dataset as *INCLIVA*_{inf}.

Proposed method

This section describes the data pre-processing, the proposed method for MoE, post-processing, evaluation metrics, and details of the implementation of the DL pipelines.

Figure 1 gives a graphical overview of the proposed DL method for detecting histological artifacts in WSIs. We proceed with the artifact detection task in two steps. First, we train binary and multiclass models for patchwise classification. The binary models are trained to detect one particular artifact, i.e., blur against artifactfree. The multiclass models provide output with six classes (five artifacts and one artifact-free). In the second step, we used these trained binary models to create a sort of MoE for inference and post-processing the predictions. We combine predictions from each expert in MoE by fusing their outputs. We apply a probability threshold to maximize sensitivity for detecting notable artifacts and providing artifact-free WSI with diagnostic potential. We deploy multiclass models with probabilistic thresholding similar to MoE. A detailed description of the proposed method is given below.

Pre-processing

We have used the EMC_{dev} cohort to prepare the dataset. This included WSIs from this cohort, which were divided into 35/10/10 training, validation, and test WSIs to prevent data leakage.

Let a WSI at magnification level 40× (sometimes known as 400×) be denoted by $I_{WSI(i)}^{40\times}$ for specific $(i)^{th}$ WSI. Since $I_{WSI}^{40\times}$ are huge gigapixel images, it is not possible to process the entire WSI in compute memory at once. Most CPATH systems first tile or patch the WSI, or RoI, to make computation feasible before processing it further. The initial step in the patching procedure was to separate the foreground tissue from the background (white) areas irrelevant for image analysis. Foreground/ background separation is usually done with a low-resolution version of the image, which can later be interpolated to be used with the full-resolution image. We obtained tissue foreground by transforming the RGB (red, green, and blue) color space to HSV (hue, saturation, and value). Later, Otsu thresholding was performed on the value channel to separate the foreground-containing tissue

from the background. We set a uniform patch-coordinate sampling grid over the extracted foreground. Patches having at least 70% overlap with the annotation area (R) were retrieved after the extracted foreground was tiled across the grid with a non-overlapping stride, as depicted in Fig. 1. In short, if a patch had 70% overlap with the artifact mask, then it was labeled as a histological artifact and vice versa.

Assuming $\mathcal{T}: I_{\text{WSI}(i)\in \mathbb{R}}^{40\times} \to \{\mathbf{x}_j^i; j = 1 \cdots J\}$ denotes the patching process, which gives a set of J patches over \mathbb{R} . Here, $\mathbf{x}_j^i \in \mathbb{R}^{W \times H \times C}$ corresponds to patch j with coordinates (x_{ij}, y_{ij}) from WSI_i and H, W, and C represent the width, height, and channels of the patch, respectively. We refer to this prepared patch-based dataset from EMC_{dev} as $\mathcal{D} = (\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ containing N patches. Here \mathbf{x}_n is a vector and denotes n-th instance with $224 \times 224 \times 3$ pixels, and $y_n \in \{0, ..., k\}$ is a scalar, where k = 1 for binary and k = 5 for multiclass dataset formulation. For instance, in a multiclass dataset, '0' represents artifact-free class, and {1,2,3,4,5} correspond to the blood, blur, air bubbles, damaged tissue, and folded tissue classes, respectively. Table 1 shows the breakdown of patches in each subset of the dataset \mathcal{D} and Fig. 2 shows example instances for all classes obtained from $I_{\rm WSI}^{40\times}$. This training and development dataset, named HistoArtifacts, is publicly available and can be downloaded from Zenodo.

Feature extractors and classifiers

The feature extractor and classifier are two significant components of most DL models for classification tasks. Feature extractors are crucial in DL algorithms as they help identify critical features in the data. In short, it reduces the dimensionality of the image and facilitates classification from a vector. Based on artifact detection works in the literature [16, 78], we have selected two popular DL

Table 1 Breakdown of the number of patches obtained in each class of the dataset \mathcal{D} , obtained form EMC_{dev} after preprocessing

(Label) class	(35 WSIs) training	(10 WSIs) validation	(10 WSIs) test	Total
(0) Artifact-free	5,249	1,591	965	7,805
(1) Blood	16,743	4,186	5996	26,655
(2) Blur	5,661	754	1,137	7,552
(3) Air bubbles	2,499	1,175	846	4520
(4) Damaged Tissue	2,577	332	1,013	3,922
(5) Folded Tissue	998	114	131	1,243



Fig. 2 Examples of artifact-free and artifact-classes patches in our prepared patch-based dataset \mathcal{D} from EMC_{dev}, and extracted at 40x magnification

architectures as feature extractors due to their smaller parametric size and faster inference: i) DCNN-based MobileNet-v3 [32] architecture, and ii) Vision transformerbased ViT-Tiny [33] architecture.

MobileNetv3: MobileNet-v3 is a SOTA DCNN architecture proposed by Howard et al. [32] and is part of the family of computationally efficient models for small devices by Google. The basic building blocks of MobileNetv3 include depth-wise separable convolutions and inverted residual blocks designed to reduce computational complexity and improve accuracy. MobileNet-v3 is optimized through a combination of hardware-aware network architecture search and novel architecture advances, including the use of hard-swish activation and squeeze-and-excitation modules [32]. This architecture is released in different variants. The large architecture variant (used in this work) has a 5.4M parameter and is lightweight and efficient, making it suitable for computationally efficient image classification pipelines.

Vision Transformer: Vision Transformers (ViTs) have gained attention as a new SOTA for image recognition tasks [26, 27]. ViT architecture breaks down an input image into a series of smaller patches, linearly embeds each patch, adds position embeddings, and then feeds the resulting sequence of vectors to a standard Transformer encoder [82]. This Transformer encoder consists of a stack of identical layers. It uses a self-attention mechanism to focus on different parts of the input by computing a weighted sum of the input features based on their similarity. We use a lightweight and efficient variant of the ViT architecture, ViT-Tiny [33], with 6M parameters for faster inference.

We apply transfer learning to train DL models and update model parameters at each epoch. Assume ϕ represents our feature extractor with θ_f parameters. Then, for the input patch (\mathbf{x}_n) with ground truth (y_n), we get a flattened feature embedding (a_n) using;

$$\phi_{\theta_f}(\mathbf{x}_n) = a_n \quad \text{where} \quad a_n = \{a_1, a_1, ..., a_z\}$$
 (1)

For patch-wise classification, we train classifiers in a binary and multiclass fashion. We appended a three-layer fully connected (FC) classifier (C_{θ_c}) at the end of the feature extractor. Let us denote our DL models with notation ψ_{θ} , where $\theta = \theta_f \cup \theta_c$, denotes the parameter set of both the feature extractor and the classifier. To obtain the output probability vector (P_{y_n}) for the input patch, we apply softmax (σ) to the output logits of the classifier as shown in Eq. (2). For instance, binary models predict (artifact vs. artifact-free), and multiclass models predict (five artifact classes vs. artifact-free), as shown in Eq. (3).

$$P_{y_n}(\mathbf{x}_n) = \psi_{\theta}(\mathbf{x}_n) = \sigma(C_{\theta_c}(\phi_{\theta_f}(\mathbf{x}_n))) = \sigma(C_{\theta_c}(a_n))$$
(2)

$$P_{y_n} = \begin{cases} \left[p_{y_0}, p_{y_1} \right]^T & \text{if binary} \\ \left[p_{y_0}, p_{y_1}, p_{y_2}, p_{y_3}, p_{y_4}, p_{y_5} \right]^T & \text{if multiclass} \end{cases}$$
(3)

Here, y_{p_0} is the probability of being an artifact-free class. In the binary model, y_{p_1} corresponds to artifact class and in the multiclass model, $[y_{p_1}, y_{p_2}, y_{p_3}, y_{p_4}, y_{p_5}]$ are predicted probabilities for blood, blur, air bubbles, damaged tissue, and folded tissue classes respectively. Finally, we calculate cross-entropy loss between the ground truth and the prediction, back-propagate this loss, and update model parameters, θ , at each epoch based on the experimental setup explained in "Implementation details" section.

$$L_{CE}(y_n, P_{y_n}) = \begin{cases} -y_n \cdot log(p_{y_0}) + (1 - y_n) \cdot log(1 - p_{y_0}) & \text{for binary} \\ -\sum_{i=0}^k y_n \cdot log(p_{y_i}) & \text{for multiclass} \end{cases}$$
(4)

To obtain final predictions (\hat{P}_{y_n}) for classes, we apply argmax to P_{y_n} .

$$\hat{P}_{y_n} = \operatorname{argmax}(P_{y_n}) \tag{5}$$

At the inference stage, we establish four DL pipelines using combinations of trained models, i.e., multiclass models (with MobileNet-v3 and ViT-Tiny) and MoEs (combining binary MobileNet-v3 and ViT-Tiny), as explained further in the following sections.

Mixture of experts

The "mixture of experts (MoE)" DL approach is often confused with deep ensembles. A deep ensemble combines DL models trained on the same data using different seed initializations or hyperparameters to learn different aspects of the data [81]. Unlike deep ensemble, in MoE, each DL model is trained for a specific task (blur, fold, blood, folded tissue, and damaged tissue detection) to become a specialist in particular task of data. Instead of applying simple majority voting like deep ensembles, a gating mechanism forms the final prediction, incorporating output from diverse experts and improving robustness.

Our proposed DL scheme is a kind of MoE where we integrate five identical DL architectures (also called base learners or experts) after training on the parts of the data (similar to bagging). Bagging offers the advantage of reducing variance, thus eliminating overfitting by training models on subsets of data. This parallel and data-independent training strategy avoids affecting the results of other experts. We form two MoE-based DL pipelines, ViTs-based MoE and DCNNs-based MoE, by choosing five base learners (DCNN or ViT architectures as explained in "Feature extractors and classifiers" section). All these experts are trained on five overlapping subsets, $\{D_{blood}, D_{blur}, D_{airbubble}, D_{damaged}, D_{folds}\} \in D$. Each sub-dataset

contains a distinct artifact class and the same artifact-free class as shown in Fig. 3. For simplicity, we transform ground truth labels as a positive class with the label '1' for artifact-free and a negative class with the label '0' for the artifact class.

The contingent MoE model, Ω , forms a single prediction using the aggregation function (*G*). *G* is similar to gating, which combines the output probabilities of the experts using a fusion approach. In short, the proposed approach formulates MoE trained on individual artifact morphology detection tasks. For artifact models $\psi_i \in \{\psi_{blood}, \psi_{blur}, \psi_{airbubble}, \psi_{damage}, \psi_{fold}\}$, we only utilize the prediction for negative class (P_{ψ_0}) (a.k.a probability of being an artifact), and fuse binary outputs for Ω as shown in Eq. (7).

$$P_{\Omega_y} = G(\psi_{blood}, \psi_{blur}, \psi_{airbubble}, \psi_{damage}, \psi_{fold}) \quad (6)$$

$$P_{\Omega_{y}} = \begin{cases} 1 - max(P_{\psi_{i_{0}}}) & \text{for artifact-free (positive) class} \\ max(P_{\psi_{i_{0}}}) & \text{for artifact (negative) class} \end{cases}$$
(7)

To evaluate the final prediction (\hat{P}_{Ω_y}) , we adopt a form of meta-learning by placing a constraint on maximizing the sensitivity of the model for the positive (artifact-free) class. Therefore, we introduce a probability threshold, t_s , to handle previously unseen tissue morphology and avoid misclassifying artifact-free patches with potential diagnostic relevance. In other words, if the probability of being a positive class in P_{Ω_y} is higher than t_s , then we assign *artifact-free* label to the patch as shown in Eq. (8). Here, t_s would help to efficiently minimize false negatives without re-training models with a new cohort of WSIs with different tissue types or staining. We determine the best value of t_s by maximizing the true positive rate (sensitivity) in the receiver operating characteristic (ROC) curve over the validation data.

$$\hat{P}_{\Omega_{y}} = \begin{cases} Artifact - free & \text{if } P_{\Omega_{y_{0}}} \ge t_{s} \\ Artifact_{k} & \text{Otherwise } k \in \{1, 2, 3, 4, 5\} \end{cases}$$
(8)

Multiclass models

In case of multiclass models (ψ_{multi}) with predicted probability distribution $P_{\psi_{y_i}} \forall i \in \{0, 1, 2, 3, 4, 5\}$. We find the probability threshold (t_s) by maximizing sensitivity similar to MoE (see "Mixture of experts" section). In other words, if the predicted probability for the artifact-free class is higher than t_s , then the patch is assigned *artifact-free* label. Otherwise, the artifact label with the highest probability value is assigned (see Eq. (9)).

$$\hat{\boldsymbol{p}}_{\boldsymbol{\psi}_{multiy}} = \begin{cases} Artifact - free & \text{with } p_{\psi_{y_0}} & \text{if } p_{\psi_{y_0}} \ge t_s \\ Artifact_k & \text{with } p_{\psi_{y_k}} & max(p_{\psi_{y_1}}, p_{\psi_{y_2}}, \dots, p_{\psi_{y_k}}) & \text{Otherwise} \end{cases}$$
(9)

Post-processing

At the inference stage, we utilize predictions for both artifact detection and QC applications, as illustrated in the post-processing part of Fig. 1. Since the predictions of DL models are patch-based, we need to stitch patches back to see the overall view of the tissue in the WSI structure. However, stitching smaller patches introduces boundary artifacts (blockish appearance) [4]. To avoid this problem, we turn to the matrixfilling approach.

For patch x_i with coordinates (x_0, y_0) , the next consecutive patch $x_{(i+1)}$ holds the difference of sampling stride (s) with coordinates $(x_1, y_1) = (x_{0+s}, y_{0+s})$. Here, s equals the patch size owing to a uniform, non-overlapping grid. For the segmentation map, we use a matrix (M), a downscale version of the original resolution, to assign predicted class k.

$$M[x_0: x_0 + s, y_0: y_0 + s] = k \quad \text{where} \quad s = 224 \quad (\text{patch-size})$$
$$M[x_1: x_1 + s, y_1: y_1 + s] = k \quad \text{where} \quad k = \{0, 1, ...5\}$$
(10)

Since M is down-scaled to sampling stride size, every filled box can be seen as a pixel in the final segmentation map (see 1 in Fig. 4). We use filled-in M for the artifact report to calculate the percentage of predicted patches with artifact class k over the total number of patches (N_{tot}) in the foreground. See 2 in Fig. 4 for an example artifact report for QC.

$$Per_k = \frac{N_k}{N_{tot}} * 100\%$$
 where N_k = Number of patches predicted with class k (11)

We denote the artifact-free post-processed region as ρ . It measures the usefulness of the WSI and can be compared against a predefined threshold τ for assessing its suitability (accepting or discarding) for developing DL algorithms.

$$\rho = \frac{\text{Number of artifact-free pixels} (N_{k_0})}{\text{Total number of pixels in the foreground} (N_{tot})}$$
(12)

To highlight the histologically relevant region, we binarize M to M_{ρ} and treat all artifact classes as a single class, as shown in Eq. (13). The binary mask (M_{ρ}) indicates the potentially histologically relevant RoI (see 3 in Fig. 4). Later, we apply a morphological closing operation to remove small holes in the final mask.

$$M_{\rho_{(i,j)}} = \begin{cases} 1, & \text{if } M_{(i,j)} = k_0 \quad \text{(artifact-free)} \\ 0, & \text{Otherwise} \end{cases}$$
(13)

Finally, obtain artifact-free WSI by performing the Hadamard product between M_{ρ} and the original WSI $(I \in \mathbb{R}^{m \times n})$ with the dimensions of $m \times n$ (see Eq. (14)). Using the nearest interpolation, we resize the M_{ρ} mask



Histological Artifacts

Fig. 3 An overview of the mixture of experts (MoE) formation for artifact detection. Five base learners (either MobileNet-v3 or ViT-Tiny deep learning architectures) are trained on overlapping sub-datasets to learn the distinct morphology of each artifact. Labels are transformed to take the artifact class as a negative class. A fusion function integrates output from all experts to form a predictive probability distribution for the final prediction. A meta-learned probability threshold is applied to maximize the sensitivity of the MoE



Fig. 4 Overview of deep learning pipeline emphasizing the post-processing stage during the inference. *Pre-processing*: The whole slide image (WSI) is split, and every patch is stored with its corresponding coordinate. *Inference*: Every patch is assigned a label using a mixture of experts or multiclass DL models. *Post-processing*: The matrix-based filling method assigns a color to every pixel (in the downscaled version of WSI) at the corresponding coordinate location. Post-processing provides: 1) Segmentation map; 2) Artifact report for quality control; 3) Artifact-free region of interest map, and 4) Artifact-refined WSI for computational analysis

to $m \times n$. Let us denote the element at the *i*-th row and *j*-th column of M_{ρ} as $M_{\rho}(i,j)$, and the corresponding element in *I* as I(i, j). This element-wise operation between M_{ρ} and *I* removes any regions or areas with the presence of artifacts (see 4 in Fig. 4) and $I_{artifact-free}$ can be written as:

$$(I \odot M_{\rho})_{ij} = \begin{bmatrix} M_{\rho(1,1)} \cdot I_{(1,1)} & M_{\rho(1,2)} \cdot I_{(1,2)} & \dots & M_{\rho(1,n)} \cdot I_{(1,n)} \\ M_{\rho(2,1)} \cdot I_{(2,1)} & M_{\rho(2,2)} \cdot I_{(2,2)} & \dots & M_{\rho(2,n)} \cdot I_{(2,n)} \\ \vdots & \vdots & \ddots & \vdots \\ M_{\rho(m,1)} \cdot I_{(m,1)} & M_{\rho(m,2)} \cdot I_{(m,2)} & \dots & M_{\rho(m,n)} \cdot I_{(m,n)} \end{bmatrix}$$

$$(14)$$

Evaluation metrics

For performance comparison, we report accuracy, sensitivity, and the F1-score. Let TP, FN, FP, and TN denote true positive and false negative, false positive, and true negative predictions, respectively. Here, a positive class refers to a patch without artifacts (artifact-free patch). Then, the confusion matrix (CM) is a tabular representation of the model's predictions using TP, FN, FP, and TN. Accuracy is the proportion of correct predictions to the total number of predictions and is defined as Acc. = (TP + TN)/(TP + FN + FP + TN). Sensitivity, also known as recall, measures the proportion of actual positives correctly identified by the model and is termed Sens. = TP/(TP + FN). High sensitivity is essential to retaining potentially relevant (artifact-free) RoIs for

in bold, and the second best results are directimed in each column						
DL architectu	re .	Acc.(%)	<i>F</i> 1	Acc.afree (%)	F1 _{afree}	Sens. _{afree}
DCNNs	MoE	92.08	91.87	97.82	88.66	90.12
	Multiclass	93.48	93.43	94.96	78.64	96.89
	Binary	<u>95.92</u>	<u>95.26</u>	-	-	<u>94.68</u>
ViTs	MoE	94.81	94.53	97.84	89.06	91.92
	Multiclass	94.29	94.48	96.79	83.80	86.84
	Binary	97.45	97.46	-	-	87.25

Table 2 Performance of artifact processing pipelines on the validation set of *EMC_{dev}* cohort "Training and development data". Various DL pipelines, including the mixture of experts (MoE) and multiclass models using SOTA DCNN and ViT architectures, are deployed. A simple binary formulation is used for a fair comparison, and accuracy for the artifact-free class is reported. The best results are marked in bold, and the second-best results are underlined in each column

the diagnostic algorithm. On the other hand, specificity Specs. = TP/(TP + FP), quantifies the performance of a model in distinguishing negative instances from those falsely labeled as positive. In our application, high specificity filters out irrelevant information (artifacts) appearing in relevant RoIs. The F1 score is the harmonic mean of precision and recall and is calculated as $F1 = 2 \cdot (\text{precision} \cdot \text{recall})(\text{precision} + \text{recall}),$ where precision = TP/(TP + FP). For overall segmentation, dice co-efficient is reported. Dice scores the overlap between the predicted segmentation and the ground truth and ranges from 0 to 1, where 1 indicates perfect overlap between the predicted and ground truth segmentation. We use model weights with the lowest validation loss during the training to report these evaluation metrics.

We have considered FLOPS, parameters, and inference time for computational complexity evaluation. FLOPS measures the number of floating-point operations required by a specific algorithm. The number of parameters refers to the learnable parameters in the model that are used to perform operations, where high parameters result in more FLOPS. Finally, inference time is the time the DL model consumes to make predictions over a patch. These metrics, combined, provide a comprehensive understanding of the DL model's performance and computational efficiency, which are crucial for assessing the practical applicability in real-world scenarios.

Implementation details

The code was implemented using Python. The patch extraction was accomplished using the Pyvips¹ library. During the patching, we used torch multiprocessing² to carry out process pooling for faster pre-processing. The extracted patches were standardized to the mean and standard deviation of ImageNet [83] due to transfer

learning over ImageNet weights. To compensate for the scarcity of labeled data, augmentation is applied at each epoch during the training [34, 84]. We used random geometric transformations, including rotations and flips, both horizontally and vertically. Our DL models consist of a feature extractor and a classifier with three fully connected (FC) layers. We used state-of-the-art architectures MobileNet-v3 [32] and ViT-Tiny [33] as backbones for feature extractors. MobileNe-tv3 was borrowed from the Pytorch³ DL framework, and ViT-Tiny was taken from the Timm⁴ library. Both of these backbones were initialized with ImageNet weights. We referred to the best hyperparameter settings from works [16, 58, 78] and fixed final parameters to cross-entropy loss, SGD optimizer, ReduceLRonPlateau scheduler initialized with 0.01, batch size of 128, early stopping of 20 epoch over the validation loss to avoid overfitting, dropout of 0.2, and fixed random seed for reproducibility. All training and inference experiments were done on the Nvidia A100 40GB GPU. The source code is available at Github.

Experimental results and discussion

This section presents experimental results for training and validating DL pipelines of the EMC cohort and discusses their performance on validation, testing, and external data.

Validation on the EMC_{dev} cohort

This experiment aims to evaluate the performance of the proposed MoE and multiclass models for artifact detection task. These pipelines consist of four DL approaches using MoE and multiclass models based on DCNNs (MobileNet-v3 [32]) and ViTs (ViT-Tiny [33]). For simplicity, we will refer to DCNNs or ViTs in the discussion. For a baseline comparison, we also trained binary classification models (DCNN and ViT) using the entire EMC_{dev}

¹ https://libvips.github.io/pyvips/

² https://pytorch.org/docs/stable/multiprocessing.html

³ https://github.com/pytorch/pytorch

⁴ https://timm.fast.ai/

dataset in a binary fashion. In other words, we wanted to compare the benefits and drawbacks of the simpler classification model against a MoE and their computational and performance trade-offs for efficient DL pipelines.

We will first focus on discussing the performance aspect. Table 2 presents classification results over the EMC_{dev} validation subset. We have reported metrics for artifact-free classes to compare them fairly against baseline (binary) models. For better classification performance, we desire high sensitivity to avoid misclassifying artifact-free patches as artifacts and retain potential histologically relevant tissue for automated diagnostics. This is because the artifact detection application is not affected by one artifact class being classified as another. In the end, patches with the presence of any artifacts will be excluded from downstream (diagnostic) applications. Though the baseline models yield the best overall accuracy, they relatively underperform and exhibit lower sensitivity in classifying the artifact-free class. The MoEs outperform multiclass models and baseline models in detecting artifact-free class. Overall, both MoE pipelines give superior results for the positive class and avoid false negatives. However, the DCNN-based multiclass model yields the best sensitivity score. Table 3 shows validation results from other relevant works from the literature as a reference. The reported results can not be directly compared as the methods were trained using different data and varying experimental setup. To present an unbiased view, we test MoEs and multiclass models on unseen data from the same *EMC*_{dev} cohort.

We present generalization results in Table 4. The table reports mixed results when probabilistic thresholding is not applied. To improve the sensitivity over new data, we learn a probability threshold (t_s) using ROC curves of the validation set (see "Mixture of experts" section), as displayed in Fig. 5. We target a 98% sensitivity and obtain different t_s values for each DL pipeline, as reported in Table 4. Interestingly, the DCNN-based pipelines assign higher probability scores to the artifact-free **Table 3** Comparison of the proposed mixture of experts (MoE) against the literature on identical classification tasks. Note that the reported methods were developed using different data under different experimental setups. Thus, the results are provided for reference, not as a direct comparison. The best results are marked in bold, and the second-best results are underlined

Task and method from literature	Accuracy (%)
Folded tissue detection by [85]	92.17
Folded tissue detection by [76]	96.7
Blur detection by [86]	93.2
Blur detection by [71]	94.4
Air bubble detection by [85]	87.33
Air bubble detection by [87]	91.5
Blood detection by [69]	85.0
Damaged tissue detection by [88]	90.0
Five artifacts - MoE-DCNNs (Ours)	97.82
Five artifacts - MoE-ViTs (Ours)	97.84

class, indicating better confidence and stronger learning of histologically relevant morphology than the ViT-based models. Figure 6 reflects similar insight that ViT-based pipelines carry weak differentiation between artifacts and artifact-free patches (see black dotted line). It is fascinating to see that probabilistic thresholding significantly improves the ability to detect artifact-free class, hinting that the proposed MoEs would be the best choice with the fewest false negatives.

To evaluate the computational aspect, Table 5 indicates the computational complexity of all four DL pipelines. Undoubtedly, MoEs have nearly five times more parameters than multiclass models. This is because each MoE combines five binary experts. Comparatively, DCNNbased pipelines can be efficient at the inference stage due to very little patch processing time per second. Similarly, MoEs have lower throughput than simpler multiclass models. As WSIs can be of different sizes, we considered reporting throughput as a more informative metric, as

Table 4 Generalization results on the test set of *EMC*_{dev} cohort "Training and development data". The table presents results over unseen data, with and without probabilistic thresholding. All metrics are calculated for the classification performance over artifact-free class. The best results in each column are marked in bold, and the second-best results are underlined

DL architecture		Without probabilistic threshold				With probabilistic threshold	
		Acc. (%)	F1	Sens.	t _s	F1	Sens.
DCNNs	MoE	97.82	88.66	89.12	0.326	86.15	97.93
	Multiclass	93.58	85.21	94.72	0.341	83.53	95.47
ViTs	MoE	<u>95.61</u>	88.91	<u>90.45</u>	0.052	<u>84.90</u>	<u>97.83</u>
	Multiclass	92.55	82.51	89.94	0.015	70.15	96.54

Table 5 A comparative analysis of computational complexity.Lower values of parameters and flops indicate computationallyefficient models, and higher throughput is desired for fasterinference

DL pipelines	Parameters (M)↓	Flops (B)↓	Throughput (p/sec.)☆
MoE (DCNNs)	17.65	1.13	178
MoE (ViTs)	27.62	5.38	128
Multiclass (DCNN)	3.53	0.22	832
Multiclass (ViT)	5.53	1.08	419

it scales linearly with the number of patches in a WSI. For instance, in skin cancer WSIs, MoEs roughly take 3-4 minutes per WSIs, and for urinary bladder cancer WSIs, which are 3-4 times larger, the end-to-end processing takes 10-12 minutes. We have to make a trade-off in selection, either choosing multiclass DCNN with better computational efficiency but relatively lower performance or based on the best performance. We prioritize classification performance and opt for the two best-performing DL pipelines from Table 4; therefore, we will use MoEs for the following experiments.



Fig. 5 ROC curves for deep learning pipelines over the validation subset. All plots highlight the area under the curves (AUC) score and best probability thresholds for maximizing F1 and sensitivity metrics



Fig. 6 Classification plots for deep learning pipelines over the validation subset. All subplots highlight the delineation (black dotted line) with the estimated value of t_s for probabilistic thresholding

Quantitative evaluation

We perform this experiment to assess the robustness of DL pipelines over external (OoD) data. For this purpose, we chose six WSIs (s_1 - s_6) from external validation data (see "External validation data" section). Note that all these WSIs were prepared and scanned by different laboratories and scanning hardware. Thus, they exhibit vast differences in staining, tissue types, and image acquisition protocols, as displayed in Fig. 7. We did not incorporate color normalization in the artifact processing pipeline due to their additive computational cost and latency [78].

Quantitative assessment is crucial to objectively evaluate the numerical performance, enabling us to compare both the proposed MoEs of DCNNs and ViTs. We require histological correctness that only an expert can provide in the form of ground truths. Therefore, all WSIs were roughly annotated by FK, UK, and AM for different artifacts. At the inference level, every patch from the WSI tissue foreground is used to perform predictions. Therefore, the error is calculated using evaluation metrics reported in "Evaluation metrics" section, where the ground truth label of a patch is only an artifact if it overlaps 70% with artifact annotation as described in "Feature extractors and classifiers" section. Table 6 presents the results for classification and segmentation performance. Since certain artifacts, such as folded tissue, have blurry areas surrounded [10]; one artifact class is likely to be predicted as another. Thus, for simplicity purposes, we report metrics for artifact-free (positive) classes only.

Both MoE pipelines experience a drop in sensitivity over breast cancer (SUH_{inf}) and skin cancer $(INCLIVA_{inf})$ WSIs. This behavior could be due to misclassifying ambiguous regions or susceptibility to specific tissue types. Since SUH_{inf} and $INCLIVA_{inf}$ WSIs are OoD data for our DL pipelines, it is interesting to see that we get high specificity scores. In short, both pipelines ensure that most of the actual artifacts present in the data are

DL pipeline	Cohort	WSIs	F1 _{afree}	Sens. _{afree}	Spec. _{afree}	Dice
MoE of DCNNs	EMC _{inf}	\$1	92.86	93.48	53.76	0.909
		s ₂	89.11	89.61	52.71	0.784
	SUH _{inf}	\$ ₃	70.91	55.07	99.09	0.487
		S4	85.51	79.78	44.57	0.572
	INCLIVA inf	\$5	60.05	43.99	80.53	0.532
		s ₆	37.39	23.55	98.97	0.506
MoE of ViTs	EMCinf	<i>s</i> ₁	93.17	93.01	60.92	0.939
		s ₂	89.34	87.97	63.18	0.795
	SUH _{inf}	\$3	68.79	54.51	79.56	0.367
		S4	87.97	85.63	26.38	0.482
	INCLIVA inf	\$5	78.92	66.02	79.71	0.559
		s ₆	45.49	42.49	42.91	0.412

Table 6 The results for quantitative evaluation for assessing the robustness of the proposed mixture of experts (MoE) approach. Qualitative evaluation is performed on external (out-of-distribution) data. The table reports classification performance corresponding to patch-wise classification and dice scores for overall segmentation maps obtained through artifact processing pipelines



Fig. 7 Hue-Saturation plot shows massive variation in the external (out-of-distribution) data. Random patches from all six WSIs (s_{TS6}) are used to calculate hue and saturation values to observe the depth of H&E staining. WSI acquisition procedures from different laboratories and scanning hardware affect the final appearance of histological images (as shown on the right)

accurately flagged. Dice score in Table 6 shows good segmentation results on the EMC_{inf} cohort. Nevertheless, EMC_{inf} is bladder cancer tissue and may carry more similarity in structural appearance.

Quantitative metrics can miss subtle nuances masked by overall performance scores. Therefore, we observe false predictions of both DCNNs-based MoE and ViTbased MoE over the better-performing cases (s_1 , s_4 , and s_5) and the worse-performing cases (s_2 , s_3 , and s_6) in OoD data. Figures 8 and 9 show ground truths and predictions masks for the better results in each cohort, and Figs. 10 and 11 show the same for the worst results in each cohort. Both MoEs densely predict artifacts in all three examples. Here, false negative instances pertain to regions identified as artifacts but were artefact-free. Conversely, false positives are cases classified as artifact-free



Fig. 8 Visualization of DCNNs-based mixture of experts' predictions with better performance over out-of-distribution data. The image shows the original WSIs (*s*₁, *s*₄, and *s*₅) along with ground truth for artifacts (combined), artifact segmentation map, and a few examples of false predictions. False negative refers to patches detected as artifacts but were artifact-free, and false positive refers to patches detected as artifact-free but belonged to any artifact class

but were labeled as any artifact class. Figure 10 highlights that DCNN-based MoE might be overdoing their job predicting certain artifacts like air bubbles. For instance, in s_6 , the entire WSI has a hazy appearance, with air trapped under most of the tissue. The false predictions for s_6 show that those examples lack cellular features. Likewise, for false positives in the case of s_2 , those specific examples were the boundary of another artifact region and contained some presence of blood. In cases s_2 and s_3 , annotations had some noise, and with the chosen mask overlap, the obtained ground truth was not accurate enough. On the other hand, the ViT-based MoE (in Fig. 11) appears

to be slightly overdoing damage detection. In most false predictions here, we might be dealing with potentially noisy and imprecise ground truth annotations. Therefore, relying on only quantitative analysis is not concrete and conclusive. We require a thorough qualitative analysis by field experts to scrutinize further the strengths and weaknesses of both MoEs in detecting artifacts.

Qualitative evaluation

In this experiment, we perform qualitative evaluations by three field experts to delve deeper into the DL pipelines' behavior and see the holistic view after the artifacts



Fig. 9 Visualization of ViTs-based mixture of experts' predictions with better performance over out-of-distribution (OoD) data. Image shows original WSIs (*s*₁, *s*₄, and *s*₅) along with ground truth for artifacts (combined), artifact segmentation map, and a few examples of false predictions. False negative refers to patches detected as artifacts but were artifact-free, and false positive refers to patches detected as artifact-free but belonged to any artifact class

refinement. While quantitative metrics provide valuable numerical insights into a model's performance, they often fall short of capturing the intricacies of segmentation results. Therefore, assessing whether the model was misclassified due to genuine limitations or imperfections in the ground truth is vital.

Three field experts (P1, P2, and P3) assessed segmentation maps for six WSIs (s_1 - s_6) from three cohorts used in the above experiment. They scored them based on visual interpretation, including how well artifacts were detected, how artifact-free regions were preserved, and the overall diagnostic usability of WSIs after the artifact processing, where field experts scored them from 1 (worst) to 10 (best). Each expert who rated these WSIs was a domain specialist on a specific cancer type (see box plot in Fig. 12). Figure 12 represents the score variability for each task across the six WSIs. The central line in each box represents the median, while the box's upper and lower edges correspond to the interquartile range.

Cohen's Kappa coefficient measures the agreement between experts, where '1' indicates perfect agreement between experts and '0' indicates agreement no better than chance. Figure 13 reveals levels of agreement



Fig. 10 Visualization of DCNNs-based mixture of experts' predictions with worst performance over out-of-distribution (OoD) data. Image shows original WSIs (*s*₂, *s*₃, and *s*₆) along with ground truth for artifacts (combined), artifact segmentation map, and a few examples of false predictions. False negative refers to patches detected as artifacts but were artifact-free, and false positive refers to patches detected as artifact-free but belonged to any artifact class

for each assessment category among the different pairs of experts for DCNNs-based MoE and ViT-based MoE. Vertical dotted lines present the average consensus across three assessment categories for each pair (in corresponding color). Both subplots highlight substantial agreement for overall usability and high average agreement between P1 and P2 (red dashed line) in Fig. 13. In contrast, artifact-free preservation has relatively lower agreement, echoing similar findings across all pairs. Based on the remarks obtained from field experts (see



Fig. 11 Visualization of ViTs-based mixture of experts' predictions with worst performance over out-of-distribution data. Image shows original WSIs (*s*₂, *s*₃, and *s*₆) along with ground truth for artifacts (combined), artifact segmentation map, and a few examples of false predictions. False negative refers to patches detected as artifacts but were artifact-free, and false positive refers to patches detected as artifact-free but belonged to any artifact class

Fig. 12), generally, better results were obtained for bladder cancer WSIs (s_1 , s_2). Although MoEs were too sensitive for detecting blurry areas, their folded and damaged regions were well segmented. In breast cancer WSIs (s_3 , s_4), adipose tissue was predicted as air bubbles (with DCNNs-based MoE) or damaged (with ViTs-based MoE). Note that the training data did not include adipose tissue, primarily fat cells. This situation can be more evident in breast samples because there is more adipose tissue in them than in other cancer types. While adipose tissue



Fig. 12 Scores for qualitative evaluation by field experts (P1, P2 and P3) for different Tasks. The boxplot provides a visual representation of the experts' assessments for predictions of OoD WSIs. The scores were provided on a scale of 1 to 10, with higher scores indicating better performance



Fig. 13 Qualitative evaluation of artifact detection by the mixture of experts (MoE) models over OoD data. Plot (a) represent Cohen's kappa score (on the x-axis) for DCNNs-based MoE and a pair of field experts on the y-axis, and (b) show scores for ViTs-based MoE. Both subplots show agreement by chance for each task. Each pair's average agreement of all three tasks is plotted as a vertical dashed line

can provide valuable contextual information and aid in certain aspects of diagnosis, its absence does not necessarily preclude an accurate assessment of breast cancer. In practical scenarios, there can also be other artifacts, like pen markings, due to manual annotation of RoIs. Since the models are trained for artifact-free class, we may expect MoE to distinguish between artifact-free and other regions. It is due to the fact that pen marking gives a similar visual of folded or blurry class patches as the cellular feature or tissue texture diminishes and becomes hard to observe. Nevertheless, encountering patches with pen markings might disappear as DP workflow becomes standard in all labs, and annotation over glass slides will be done manually rather than with markers.

The particular examples of skin cancer WSIs (s_5 , s_6) had significant air bubbles, leaving a hazy and unclear appearance over the foreground tissue. At the same time, both artifact processing pipelines were overdoing air bubble prediction, and the epidermis was predicted as blood. The performance of both MoEs is worst in these cases; one of the reasons could be the severity of artifacts and significant variation in staining in the WSI. While there is generally substantial agreement among field experts for overall diagnostic usability, there are areas, such as artifact-free preservation, where discrepancies emerge and may be more challenging to achieve. Moreover, considering inter-rater variability, DCNNs-based MoE indicates potential effectiveness for artifact detection and overall diagnostic usability.

By triangulating quantitative and qualitative analysis findings, we conclude that DCNNs-based MoE provides better generalizability and robustness with the trade-off of higher computational cost.

Conclusion

In this work, we established end-to-end deep learning (DL) pipelines, taking whole slide images (WSIs) as input and providing artifact-refined WSIs to enable computational pathology (CPATH) systems to make reliable predictions. For the development of DL pipelines, we propose the mixture of experts (MoE) scheme and multiclass models. The MoE scheme uses five base learners (experts) with underlying state-ofthe-art DL architectures (MobileNet-v3 or ViT-Tiny). The MoE captures the intricacies of different artifact morphologies and dynamically combines predictions using a fusion mechanism to generate predictive probability distribution. Later, a meta-learned probabilistic threshold is applied to improve sensitivity for histologically relevant regions. In rigorous experiments, we performed generalizability and robustness tests over DL pipelines by testing on external cohorts of different tissue types. During the investigation, we found that the MoE scheme with underlying DCNNs attains the best classification and segmentation performance with some computational trade-offs compared to multiclass models. However, if high inference speed is the desired requirement, then multiclass models can be a better alternative with some degree of performance trade-off. Furthermore, during the qualitative evaluation, field experts rated the outcomes and agreed substantially on the overall usability of DCNNs-based MoE.

Our artifact-processing DL pipelines can provide various outcomes, such as a segmentation map, artifact report, artifact-free mask with potential regions of interest with histological relevance, and an artifact-refined WSI for further computational analysis. Overall, the proposed DL solution is efficient and has a significant advantage in equipping the CPATH system with the necessary tools to isolate anomalies (or noise) from affecting automated clinical applications.

Limitations and future work

The proposed work has a limitation in that the DL models were trained on a dataset prepared from a single cohort of data. In future work, these limitations can be overcome by pooling datasets from different cohorts in training and adopting an active learning strategy to adjust meta-learned thresholding parameters for improved sensitivity. Also, by formulating tailored fusion mechanisms for different cancer types. Furthermore, artifact-refined WSIs can be tested with the corresponding diagnostic or prognostic algorithms to assess the usefulness of artifact processing pipelines for clinical practice.

Abbreviations

71001011	100115
WSI	Whole slide image
DL	Deep learning
CPATH	Computational pathology
MoE	Mixture of experts
SOTA	State-of-the-art
DCNN	Deep convolutional neural networks
ViT	Vision transformer
OoD	Out-of-distribution
DP	Digital pathology
QC	Quality control
RGB	Red, Green, Blue
HSI	Hue, Saturation, Intensity
SVM	Support vector machine
H&E	Hematoxylin and Eosin
EMC	Erasmus medical centre
SUH	Stavanger University Hospital
~ .	

Supplementary Information

The online version contains supplementary material available at https://doi.org/10.1186/s12911-024-02676-z.

Supplementary Material 1: Peer review

Acknowledgements

The European Union's Horizon 2020 research and innovation program (CLARIFY) financially supports this research work under the Marie Skłodowska-Curie grant agreement 860627.

Authors' contributions

N.K: Conceptualization, Methodology, Software, Validation, Investigation, Data curation, Writing - Original Draft, Writing - Review & Editing, Visualization; F.K: Formal analysis; Investigation; Data Resources; Writing - Review & Editing; U.K: Formal analysis; Investigation; Data Resources; Writing-Review & Editing; A.M: Formal analysis; Investigation; Data Resources; Writing - Review & Editing; C.M: Formal analysis; Data Resources; Writing - Review & Editing; Funding acquisition; E.J: Formal analysis; Data Resources; Writing - Review & Editing; Funding acquisition; T.Z: Formal analysis; Data Resources; Writing - Review & Editing; Funding acquisition; C.R: Writing - Review & Editing; K.E: Conceptualization; Methodology; Investigation; Supervision; Writing - Review & Editing; Project administration; Funding acquisition.

Funding

Open access funding provided by University of Stavanger & Stavanger University Hospital The European Union's Horizon 2020 research and innovation program (CLARIFY) financially supports this research work under the Marie Skłodowska-Curie grant agreement 860627.

Availability of data and materials

The source code is available at Github. The training and development dataset (named HistoArtifacts) can be downloaded from Zenodo.

Declarations

Ethics approval and consent to participate

This study was performed in line with the principles of the Declaration of Helsinki. The patients/participants provided written informed consent to use their data for secondary purposes. The Erasmus MC Medical Research Committee granted approval from the Institutional Review Board under the reference MEC-2018-1097. The Stavanger University Hospital's data is approved by the Norwegian Regional Committee for Medical and Health Research Ethics under REC, 2010/1241. INCLIVA Biomedical Research Institute granted approval from the Research Ethics Committee (CEIm) of the Hospital Clínico Universitario of Valencia, Spain, under the reference CIm-2020/114.

Consent for publication

Competing interests

All authors consent to the open-access publication of this research work.

The authors declare no competing interests.

Peer review

The peer review reports can be found at https://doi.org/10.1186/ s12911-024-02676-z.

Author details

¹Department of Electrical Engineering and Computer Science, University of Stavanger, 4021 Stavanger, Norway. ²Department of Urology, University Medical Center Rotterdam, Erasmus MC Cancer Institute, 1035 GD Rotterdam, The Netherlands. ³Department of Pathology, Stavanger University Hospital, 4011 Stavanger, Norway. ⁴Department of Chemistry, Bioscience and Environmental Engineering, University of Stavanger, 4021 Stavanger, Norway. ⁵Department of Pathology, INCLIVA Biomedical Research Institute, and University of Valencia, 46010 Valencia, Spain.

Received: 24 July 2024 Accepted: 9 September 2024 Published online: 07 October 2024

References

 National Cancer Institute. Environmental carcinogens and cancer risk. 2015. https://www.cancer.gov/about-cancer/causes-prevention/risk/ substances/carcinogens. Accessed 31 Aug 2023.

- World Cancer Research Fund International. Differences in cancer incidence and mortality across the globe. 2023. https://www.wcrf.org/differences-in-cancer-incidence-and-mortality-across-the-globe/. Accessed 31 Aug 2023.
- Pulumati A, Pulumati A, Dwarakanath BS, Verma A, Papineni RV. Technological advancements in cancer diagnostics: Improvements and limitations. Cancer Rep. 2023;6(2):e1764.
- Khened M, Kori A, Rajkumar H, Krishnamurthi G, Srinivasan B. A generalized deep learning framework for whole-slide image segmentation and analysis. Sci Rep. 2021;11(1):11579.
- Zhu C, Song F, Wang Y, Dong H, Guo Y, Liu J. Breast cancer histopathology image classification through assembling multiple compact CNNs. BMC Med Inform Decis Mak. 2019;19(1):1–17.
- Kanwal N, Amundsen R, Hardardottir H, Janssen EA, Detection Engan K, localization of melanoma skin cancer in histopathological whole slide images. In: 2023 31st European Signal Processing Conference (EUSIPCO). IEEE; 2023. pp. 1128–35.
- Car LT, Papachristou N, Bull A, Majeed A, Gallagher J, El-Khatib M, et al. Clinician-identified problems and solutions for delayed diagnosis in primary care: a PRIORITIZE study. BMC Fam Pract. 2016;17:1–9.
- Pallua J, Brunner A, Zelger B, Schirmer M, Haybaeck J. The future of pathology is digital. Pathol-Res Pract. 2020;216(9):153040.
- Inc DSRS. Digital Science and Research Solutions Inc. Query: "CPATH" OR "Computational Pathology" OR "Digital Pathology". https://app.dimen sions.ai/analytics/publication/overview/timeline?search_mode=conte nt&or_facet_year=2018 &or_facet_year=2019 &or_facet_year=2020 &or_facet_year=2021 &or_facet_year=2022 &or_facet_year=2023 & search_text=Digital%20Pathology &search_type=kws &search_field= full_search. Accessed Aug 2023.
- Kanwal N, Pérez-Bueno F, Schmidt A, Molina R, Engan K. The devil is in the details: Whole Slide Image acquisition and processing for artifacts detection, color variation, and data augmentation: a review. IEEE Access. 2022.
- Campanella G, Rajanna AR, Corsale L, Schüffler PJ, Yagi Y, Fuchs TJ. Towards machine learned quality control: A benchmark for sharpness quantification in digital pathology. Comput Med Imaging Graph. 2018;65:142–51.
- Hosseini MS, Bejnordi BE, Trinh VQH, Chan L, Hasan D, Li X, et al. Computational pathology: a survey review and the way forward. J Pathol Inform. 2024:100357.
- Louis DN, Gerber GK, Baron JM, Bry L, Dighe AS, Getz G, et al. Computational pathology: an emerging definition. Arch Pathol Lab Med. 2014;138(9):1133–8.
- 14. Taqi SA, Sami SA, Sami LB, Zaki SA. A review of artifacts in histopathology. J Oral Maxillofac Pathol. 2018;22(2):279.
- Bindhu P, Krishnapillai R, Thomas P, Jayanthi P. Facts in artifacts. J Oral Maxillofac Pathol. 2013;17(3):397.
- Kanwal N, Eftestøl T, Khoraminia F, Zuiverloon TC, Engan K. Vision Transformers for Small Histological Datasets Learned Through Knowledge Distillation. In: Pacific-Asia Conference on Knowledge Discovery and Data Mining. Springer; 2023. pp. 167–179.
- 17. Wright AI, Dunn CM, Hale M, Hutchins GG, Treanor DE. The effect of quality control on accuracy of digital pathology image analysis. IEEE J Biomed Health Inform. 2020;25(2):307–14.
- Tabatabaei Z, Colomer A, Engan K, Oliver J, Naranjo V, Residual block Convolutional Auto Encoder in Content-Based Medical Image Retrieval. In: 2022 IEEE 14th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP). IEEE; 2022. pp. 1–5.
- Chen C, Chen C, Ma M, Ma X, Lv X, Dong X, et al. Classification of multidifferentiated liver cancer pathological images based on deep learning attention mechanism. BMC Med Inform Dec Making. 2022;22(1):1–13.
- Fuster S, Khoraminia F, Kiraz U, Kanwal N, Kvikstad V, Eftestøl T, et al. Invasive Cancerous Area Detection in Non-Muscle Invasive Bladder Cancer Whole Slide Images. In: 2022 IEEE 14th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP). Nafplio: IEEE; 2022. p. 1–5.
- 21. Litjens G, Kooi T, Bejnordi BE, Setio AAA, Ciompi F, Ghafoorian M, et al. A survey on deep learning in medical image analysis. Med Image Anal. 2017;42:60–88.
- Chen L, Li S, Bai Q, Yang J, Jiang S, Miao Y. Review of image classification algorithms based on convolutional neural networks. Remote Sens. 2021;13(22):4712.

- Lu Z, Xie H, Liu C, Zhang Y. Bridging the gap between vision transformers ers and convolutional neural networks on small datasets. Adv Neural Inf Process Syst. 2022;35:14663–77.
- Zhu H, Chen B, Yang C. Understanding Why VIT Trains Badly on Small Datasets: An Intuitive Perspective. 2023. arXiv preprint arXiv:2302.03751.
- Atabansi CC, Nie J, Liu H, Song Q, Yan L, Zhou X. A survey of Transformer applications for histopathological image analysis: New developments and future directions. Biomed Eng Online. 2023;22(1):96.
- Naseer MM, Ranasinghe K, Khan SH, Hayat M, Shahbaz Khan F, Yang MH. Intriguing properties of vision transformers. Adv Neural Inf Process Syst. 2021;34:23296–308.
- Bhojanapalli S, Chakrabarti A, Glasner D, Li D, Unterthiner T, Veit A. Understanding robustness of transformers for image classification. In: Proceedings of the IEEE/CVF international conference on computer vision. IEEE; 2021. p. 10231–10241.
- Hsu ST, Su YJ, Hung CH, Chen MJ, Lu CH, Kuo CE. Automatic ovarian tumors recognition system based on ensemble convolutional neural network with ultrasound imaging. BMC Med Inform Dec Making. 2022;22(1):298.
- Meng Z, Zhao Z, Li B, Su F, Guo L. A Cervical Histopathology Dataset for Computer Aided Diagnosis of Precancerous Lesions. IEEE Trans Med Imaging. 2021;40(6):1531–41.
- Abe T, Buchanan EK, Pleiss G, Zemel R, Cunningham JP. Deep ensembles work, but are they necessary? Adv Neural Inf Process Syst. 2022;35:33646–60.
- Mohammed A, Kora R. A comprehensive review on ensemble deep learning: Opportunities and challenges. J King Saud Univ-Comput Inform Sci. 2023;35(2):757–74.
- Howard A, Sandler M, Chu G, Chen LC, Chen B, Tan M, et al. Searching for mobilenetv3. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. IEEE; 2019. p. 1314–1324.
- Touvron H, Cord M, Douze M, Massa F, Sablayrolles A, Jégou H. Training data-efficient image transformers & distillation through attention. In: International conference on machine learning. PMLR; 2021. pp. 10347–10357.
- Morales S, Engan K, Naranjo V. Artificial intelligence in computational pathology-challenges and future directions. Digit Signal Process. 2021;119:103196.
- Bulten W, Kartasalo K, Chen PHC, Ström P, Pinckaers H, Nagpal K, et al. Artificial intelligence for diagnosis and Gleason grading of prostate cancer: the PANDA challenge. Nat Med. 2022;28(1):154–63.
- Khoraminia F, Fuster S, Kanwal N, Olislagers M, Engan K, van Leenders GJ, et al. Artificial Intelligence in Digital Pathology for Bladder Cancer: Hype or Hope? A Systematic Review. Cancers. 2023;15(18):4518.
- Gay J, Harlin H, Wetzer E, Lindblad J, Sladoje N. Texture-based oral cancer detection: A performance analysis of deep learning approaches. In: 3rd NEUBIAS Conference. Digitala Vetenskapliga Arkivet (DiVA); 2019.
- Gandomkar Z, Brennan PC, Mello-Thoms C. MuDeRN: Multi-category classification of breast histopathological image using deep residual networks. Artif Intell Med. 2018;88:14–24. https://doi.org/10.1016/j.artmed.2018. 04.005. https://www.sciencedirect.com/science/article/pii/S093336571 7305031
- Wessels F, Schmitt M, Krieghoff-Henning E, Nientiedt M, Waldbillig F, Neuberger M, et al. A self-supervised vision transformer to predict survival from histopathology in renal cell carcinoma. World J Urol. 2023;41(8):2233–41.
- Stegmüller T, Bozorgtabar B, Spahr A, Thiran JP. Scorenet: Learning non-uniform attention and augmentation for transformer-based histopathological image classification. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. IEEE; 2023. p. 6170–6179.
- Perincheri S, Levi AW, Celli R, Gershkovich P, Rimm D, Morrow JS, et al. An independent assessment of an artificial intelligence system for prostate cancer detection shows strong diagnostic accuracy. Mod Pathol. 2021;34(8):1588–95.
- 42. Huang G, Liu Z, Van Der Maaten L, Weinberger KQ. Densely connected convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. IEEE; 2017. p. 4700–4708.
- He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. IEEE; 2016. p. 770–778.

- Szegedy C, Liu W, Jia Y, Sermanet P, Reed S, Anguelov D, et al. Going deeper with convolutions. In: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). 2015. pp. 1–9.
- Caron M, Touvron H, Misra I, Jégou H, Mairal J, Bojanowski P, et al. Emerging properties in self-supervised vision transformers. In: Proceedings of the IEEE/CVF international conference on computer vision. IEEE; 2021. p. 9650–9660.
- Zidan U, Gaber MM, Abdelsamea MM. SwinCup: Cascaded swin transformer for histopathological structures segmentation in colorectal cancer. Expert Syst Appl. 2023;216:119452.
- Srinidhi CL, Ciga O. Martel AL. Deep neural network models for computational histopathology: a survey. Med Image Anal. 2021;67:101813.
- Riasatian A, Babaie M, Maleki D, Kalra S, Valipour M, Hemati S, et al. Fine-Tuning and training of densenet for histopathology image representation using TCGA diagnostic slides. Med Image Anal. 2021;70:102032. https:// www.sciencedirect.com/science/article/pii/S1361841521000785.
- Talo M. Automated classification of histopathology images using transfer learning. Artif Intell Med. 2019;101:101743. https://www.sciencedirect. com/science/article/pii/S0933365719307110.
- Wang Y, Peng T, Duan J, Zhu C, Liu J, Ye J, et al. Pathological Image Classification Based on Hard Example Guided CNN. IEEE Access. 2020;8:114249–58.
- Wang C, Gong W, Cheng J, Qian Y. DBLCNN: Dependency-based lightweight convolutional neural network for multi-classification of breast histopathology images. Biomed Signal Process Control. 2022;73:103451. https://www.sciencedirect.com/science/article/pii/S174680942101048X.
- Gao Z, Hong B, Zhang X, Li Y, Jia C, Wu J, et al. Instance-based vision transformer for subtyping of papillary renal cell carcinoma in histopathological image. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part VIII 24. Springer; 2021. pp. 299–308.
- Schömig-Markiefka B, Pryalukhin A, Hulla W, Bychkov A, Fukuoka J, Madabhushi A, et al. Quality control stress test for deep learning-based diagnostic model in digital pathology. Mod Pathol. 2021;34(12):2098–108.
- Linmans J, Raya G, van der Laak J, Litjens G. Diffusion models for out-of-distribution detection in digital pathology. Med Image Anal. 2024;93:103088.
- Ghaffari Laleh N, Truhn D, Veldhuizen GP, Han T, van Treeck M, Buelow RD, et al. Adversarial attacks and adversarial robustness in computational pathology. Nat Commun. 2022;13(1):5711.
- 56. Kanwal N, Engan K. Extract, detect, eliminate: Enhancing reliability and performance of computational pathology through artifact processing pipelines. Sci Talks. 2024;9.
- 57. Kothari S, Phan JH, Wang MD. Eliminating tissue-fold artifacts in histopathological whole-slide images for improved image-based prediction of cancer grade. J Pathol Inform. 2013;4(1):22.
- Kanwal N, López-Pérez M, Kiraz U, Zuiverloon TC, Molina R, Engan K. Are you sure it's an artifact? Artifact detection and uncertainty quantification in histological images. Comput Med Imaging Graph. 2024;112:102321.
- Salvi M, Acharya UR, Molinari F, Meiburger KM. The impact of pre-and post-image processing techniques on deep learning frameworks: A comprehensive review for digital pathology image analysis. Comput Biol Med. 2021;128:104129.
- Pérez-Bueno F, Vega M, Naranjo V, Molina R, Katsaggelos AK. Super Gaussian priors for blind color deconvolution of histological images. In: 2020 IEEE International Conference on Image Processing (ICIP). IEEE; 2020. pp. 3010–3014.
- Ameisen D, Deroulers C, Perrier V, Bouhidel F, Battistella M, Legrès L, et al. Towards better digital pathology workflows: Programming libraries for high-speed sharpness assessment of Whole Slide Images. Diagn Pathol. 2014;9(1):1–7.
- 62. Shrestha P, Kneepkens R, Vrijnsen J, Vossen D, Abels E, Hulsken B. A quantitative approach to evaluate image quality of whole slide imaging scanners. J Pathol Inform. 2016;7(1):56.
- 63. Bahlmann C, Patel A, Johnson J, Ni J, Chekkoury A, Khurd P, et al. Automated detection of diagnostically relevant regions in H &E stained digital pathology slides. In: Medical Imaging 2012: Computer-Aided Diagnosis, vol. 8315. International Society for Optics and Photonics; 2012. p. 831504.
- 64. Avanaki ARN, Espig KS, Xthona A, Lanciault C, Kimpe TRL. Automatic Image Quality Assessment for Digital Pathology. In: Tingberg A, Lång K,

Timberg P, editors. Breast Imaging. Cham: Springer International Publishing; 2016. p. 431–8.

- Gao D, Padfield D, Rittscher J, McKay R. Automated training data generation for microscopy focus classification. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. Springer; 2010. pp. 446–453.
- Hashimoto N, Bautista PA, Yamaguchi M, Ohyama N, Yagi Y. Referenceless image quality evaluation for whole slide imaging. J Pathol Inform. 2012;3(1):9.
- 67. Palokangas S, Selinummi J, Yli-Harja O, Segmentation of folds in tissue section images. In: 2007 29th annual international Conference of the IEEE Engineering in Medicine and biology society. IEEE; 2007. pp. 5641–4.
- Bautista PA, Yagi Y. Detection of tissue folds in whole slide images. Proceedings of the 31st Annual International Conference of the IEEE Engineering in Medicine and Biology Society: Engineering the Future of Biomedicine, EMBC 2009. IEEE; 2009. p. 3669–3672.
- Swiderska-Chadaj Z, Markiewicz T, Cierniak S, Koktysz R, Automatic quantification of vessels in hemorrhoids whole slide images. In: 2016 17th International Conference Computational Problems of Electrical Engineering (CPEE). IEEE; 2016. pp. 1–4.
- Mercan E, Aksoy S, Shapiro LG, Weaver DL, Brunye T, Elmore JG. Localization of diagnostically relevant regions of interest in whole slide images. In: 2014 22nd International Conference on Pattern Recognition. IEEE; 2014. pp. 1179–1184.
- Albuquerque T, Moreira A, Cardoso JS. Deep Ordinal Focus Assessment for Whole Slide Images. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. IEEE; 2021. p. 657–663.
- Kohlberger T, Liu Y, Moran M, Chen PHC, Brown T, Hipp JD, et al. Wholeslide image focus quality: Automatic assessment and impact on ai cancer detection. J Pathol Inform. 2019;10(1):39.
- Wetteland R, Engan K, Eftestøl T, Kvikstad V, Janssen EAM. Multiclass tissue classification of whole-slide histological images using convolutional neural networks. ICPRAM. 2019;1:320–7.
- Wetteland R, Engan K, Eftestøl T, Kvikstad V, Janssen EA. A Multiscale Approach for Whole-Slide Image Segmentation of five Tissue Classes in Urothelial Carcinoma Slides. Technol Cancer Res Treat. 2020;19:1533033820946787.
- Clymer D, Kostadinov S, Catov J, Skvarca L, Pantanowitz L, Cagan J, et al. Decidual vasculopathy identification in whole slide images using multiresolution hierarchical convolutional neural networks. Am J Pathol. 2020;190(10):2111–22.
- Babaie M, Tizhoosh HR. Deep features for tissue-fold detection in histopathology images. In: European Congress on Digital Pathology. Springer; 2019. pp. 125–132.
- Janowczyk A, Zuo R, Gilmore H, Feldman M, Madabhushi A. HistoQC: an open-source quality control tool for digital pathology slides. JCO Clin Cancer Inform. 2019;3:1–7.
- Kanwal N, Fuster S, Khoraminia F, Zuiverloon TC, Rong C, Engan K, Quantifying the effect of color processing on blood and damaged tissue detection in whole slide images. In: 2022 IEEE 14th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP). IEEE; 2022. pp. 1–5.
- Guo C, Pleiss G, Sun Y, Weinberger KQ. On calibration of modern neural networks. In: International conference on machine learning. PMLR; 2017. pp. 1321–1330.
- Linmans J, Elfwing S, van der Laak J, Litjens G. Predictive uncertainty estimation for out-of-distribution detection in digital pathology. Med Image Anal. 2023;83:102655.
- Zhang C, Bengio S, Hardt M, Recht B, Vinyals O. Understanding deep learning (still) requires rethinking generalization. Commun ACM. 2021;64(3):107–15.
- 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. 2020. arXiv preprint arXiv:2010.11929.
- Deng J, Dong W, Socher R, Li LJ, Li K, Fei-Fei L, Imagenet: A large-scale hierarchical image database. In: 2009 IEEE ICCV. IEEE; 2009. pp. 248–55.
- 84. Wetzer E. Representation learning and information fusion: applications in biomedical image processing. Acta Universitatis Upsaliensis; 2023.
- Shakhawat HM, Nakamura T, Kimura F, Yagi Y, Yamaguchi M. Automatic quality evaluation of whole slide images for the practical use of whole slide imaging scanner. ITE Trans Media Technol Appl. 2020;8(4):252–68.

- Senaras C, Niazi MKK, Lozanski G, Gurcan MN. DeepFocus: detection of out-of-focus regions in whole slide digital images using deep learning. PLoS ONE. 2018;13(10):e0205387.
- 87. Raipuria G, Singhal N. Stress testing vision transformers using common histopathological artifacts. In: Medical Imaging with Deep Learning. 2022.
- Swiderska-Chadaj Z, Markiewicz T, Gallego J, Bueno G, Grala B, Lorent M. Deep learning for damaged tissue detection and segmentation in Ki-67 brain tumor specimens based on the U-net model. Bull Pol Acad Sci Tech Sci. 2018:849–56.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.