

Histopathology Foundation Models Enable Accurate Ovarian Cancer Subtype Classification

Jack Breen^{†1}, Katie Allen^{2,3}, Kieran Zucker³, Lucy Godson⁴, Nicolas M. Orsi^{*2,3}, and Nishant Ravikumar^{*1}

¹Centre for Computational Imaging and Simulation Technologies in Biomedicine (CISTIB), School of Computing, University of Leeds, UK

²Leeds Institute of Medical Research at St James's, School of Medicine, University of Leeds, UK

³Leeds Cancer Centre, St James's University Hospital, Leeds, UK

⁴National Pathology Imaging Cooperative (NPIC), Leeds Teaching Hospitals NHS Trust, Leeds, UK

*Indicates joint senior authors

Large pretrained transformers are increasingly being developed as generalised 'foundation models' which can underpin powerful task-specific artificial intelligence models. Histopathology foundation models show great promise across many tasks, but analyses have typically been limited by arbitrary hyperparameters that were not tuned to the specific task/dataset. We report the most rigorous single-task validation conducted to date of a histopathology foundation model, and the first performed in the context of ovarian cancer morphological subtyping. Attention-based multiple instance learning classifiers were compared using vision transformer and ResNet features generated through varied preprocessing and pretraining procedures. The training set consisted of 1864 whole slide images from 434 ovarian carcinoma cases at Leeds Teaching Hospitals NHS Trust. Five-class classification performance was evaluated using the F1 score, AUROC, and balanced accuracy through five-fold cross-validation, and these cross-validation models were ensembled for evaluation on a hold-out test set and an external set from the Transcanadian study. Reporting followed the TRIPOD+AI checklist. The vision transformer-based histopathology foundation model, UNI, performed best in every evaluation, with five-class balanced accuracies of 88% (95% CI: 82-94%) and 93% (87-98%) in hold-out internal and external testing, compared to the best ResNet model scores of 68% (60-76%) and 81% (70-90%), respectively. Normalisations and augmentations aided the generalisability of ResNet-based models, but these still did not match the performance of UNI, which gave the best external performance in any ovarian cancer subtyping study to date. UNI also gave significant improvements over an ImageNet-pretrained transformer in eight out of nine evaluations (p-values 0.0005-0.0306). Histopathology foundation models offer a clear benefit to subtyping, improving classification performance to a degree where clinical utility is tangible, albeit with an increased computational burden. Such models could provide a second opinion to histopathologists diagnosing challenging cases and may improve the accuracy, objectivity, and efficiency of pathological diagnoses overall.

Key Words

Computer Vision, Digital Pathology, Computational Pathology, Ovarian Carcinoma

[†]Corresponding author - scjib@leeds.ac.uk.

INTRODUCTION

Ovarian cancer is the eighth most common cancer in women worldwide and typically has a poor prognosis, with 324,000 diagnosed cases translating to 207,000 deaths annually [1]. It is represented by an array of histological (morphological) subtypes with distinct prognoses and treatment options [2]. Five carcinoma subtypes account for approximately 90% of all ovarian cancers - high-grade serous (HGSC, 70%), endometrioid (EC, 11%), clear cell (CCC, 10%), low-grade serous (LGSC, 5%), and mucinous carcinomas (MC, 4%) [3–5].

Histological subtyping is an essential component of the diagnostic process, but it can be challenging. From an individual slide, pathologists only exhibit concordance on an ovarian cancer diagnosis around 80% of the time [6]. As a result, they often request ancillary tests (such as P53 immunohistochemistry) or seek a second opinion from a gynaecological subspecialty expert, with associated logistical and financial burdens. With increasing cancer rates [1] and complexity in diagnostic testing, histopathology services are increasingly struggling to meet demand. For example, most histopathology departments in the UK routinely resort to outsourcing work or hiring temporary staff [7], despite the UK being one of the countries with the most pathologists per capita [8]. Any delays resulting from demand outstripping diagnostic resources risk catastrophic impacts on patient outcomes, with a four-week delay in cancer treatment being associated with an approximately 10% increased mortality rate among patients [9].

Conceptually, artificial intelligence (AI) may offer clinical value by giving an efficient second opinion to histopathologists, streamlining the diagnostic process and perhaps offering additional support when subspecialty experts are not readily available [10]. However, AI models for ovarian cancer diagnosis have not yet demonstrated clinical utility, with most research being small-scale prototyping [11] without regulatory approval for clinical use in Europe or the United States [12]. AI for ovarian cancer subtyping has constituted a small field of research where, aside from our work [13, 14], research has almost exclusively been published by a single group [15–21]. While the accuracy of such models has increased over time, the best models still only achieve around 80% accuracy [20], and lack sufficient real-life testing.

One issue limiting AI in histopathology is that whole slide images (WSIs) are orders of magnitude too large for conventional (single instance) models, therefore multiple instance learning (MIL) is often employed [22]. In MIL, individual patches (the ‘instances’) are separately processed and then aggregated to learn information about a WSI. These models are impractical to train end-to-end with such large images, so frozen patch feature extractors are often used. As such, any limitation in the pretrained feature extractor can limit downstream classification performance.

In applying MIL to WSI-level classification, many researchers have used ImageNet-pretrained ResNets [23] for patch feature extraction [13, 19, 24–27]. ImageNet (a set of 1.4 million natural images from 1000 classes) [28] is popular for model pretraining as the quantity and diversity of images enables the creation of a multi-purpose feature set. However, these generic features are likely to be suboptimal and computationally inefficient when applied to histopathology images, which contain a relatively homogeneous and restricted set of shapes and colours, with subtle differences being relevant to diagnostic decisions [5, 29].

Recently, some researchers have attempted to create histopathology ‘foundation models’ [30], using self-supervised learning techniques to generate broad histopathological feature sets which are not specific to a single organ/cancer type. These have grown rapidly, from tens of thousands of WSIs used to train models with tens of millions of parameters in 2022 and early 2023 [31–36], to hundreds of thousands of WSIs and hundreds of millions of parameters more recently [37–40]. One such model was even trained using 1.5 million WSIs [30]. Foundation models have typically been based on vision transformers (ViTs), utilizing the impressive scalability of transformers seen across many fields, most notably with large language models [41, 42]. Histopathology foundation models have exhibited impressive performance across diverse tasks [34, 39], although analyses have been relatively shallow, without thorough hyperparameter tuning and rigorous statistical comparison of downstream models. Consequently, it is unclear whether models were applied optimally (especially those exhibiting sub-optimal performance), and whether differences between them were significant. Furthermore, many analyses have been conducted using single-centre data, limiting the assessment of models’ generalisability.

In this study, we present the most comprehensive validation conducted to date comparing the standard ResNet50 feature extractor with histopathology foundation models, specifically in the context of ovarian cancer subtyping. This includes comparing whether the performance of ResNet-based MIL classifiers can match those of newer transformer-based MIL classifiers through normalisation, augmentation, and improved tissue detection techniques. The analysis includes rigorous hyperparameter tuning and evaluations through five-fold cross-validation, hold-out testing, and external validation, and was conducted with the largest repository of ovarian cancer WSIs used in any AI study to date.

METHODS

Data

A training set of 1864 adnexal tissue WSIs was retrospectively collected from 434 cases of ovarian carcinoma treated at Leeds Teaching Hospitals NHS Trust between 2008 and 2022. Cases were only included if a gynaecological pathologist had diagnosed them as one of the five most common epithelial ovarian cancer subtypes (HGSC, LGSC, CCC, MC, EC). A histopathologist (KA) independently verified all diagnoses, removing any cases with discrepancies. Several representative haematoxylin and eosin (H&E)-stained adnexal tissue glass slides were selected for each case, with only formalin-fixed, paraffin-embedded (FFPE) samples used. These samples were prepared and digitised at 40x magnification using a Leica Aperio AT2 scanner. The population-level class imbalance was reflected in the training set (Table 1), with the least common subtype (LGSC) represented by only 92 WSIs from 21 cases, compared to 1266 WSIs from 308 cases for the most common subtype (HGSC). The training set contained both primary and interval debulking surgery (IDS) specimens, the inclusion of which was previously found to be beneficial to subtype classification [43].

An independent class-balanced hold-out test set was collected through the same protocol, consisting of 100 primary surgery specimen WSIs from 30 patients. An external set of 80 WSIs from 80 patients was accessed from the *Transcanadian Study*

[29]. These had been digitised using an AperioScope scanner and made available at 20x magnification, alongside subtype labels that had been determined by a gynaecological pathologist. This enabled an analysis of generalisability to different slide preparation and scanning procedures.

Carcinoma Subtype	Training WSIs (Patients)	Hold-out WSIs (Patients)	External WSIs (Patients)
High-Grade Serous (HGSC)	1266 (308)	20 (7)	30 (30)
Low-Grade Serous (LGSC)	92 (21)	20 (6)	9 (9)
Clear Cell (CCC)	198 (45)	20 (7)	20 (20)
Endometrioid (EC)	209 (38)	20 (5)	11 (11)
Mucinous (MC)	99 (22)	20 (5)	10 (10)
Total	1864 (434)	100 (30)	80 (80)

Table 1. Dataset breakdown for the training (cross-validation) set, independent internal hold-out test set, and external validation set. Numbers in brackets indicate the number of unique patients.

Feature Extraction

The first computational step in the classification pipeline (Figure 3) was tissue segmentation. A significant proportion of most WSIs is non-tissue background which can be discarded using saturation thresholding, where only the pixels with saturation higher than the threshold are labelled as tissue. Otsu thresholding [44] automatically determines the threshold for each image by minimising the variance within the high-saturation and low-saturation groups. Saturation thresholding is computationally efficient, but risks including artifacts such as bubbles, pen marks, and coverslip edges in the foreground. While more robust (and complex) tissue segmentation techniques exist [45, 46], we focused on simple approaches as the attention mechanism in the classification models should learn to ignore any remaining artifacts. We compared the CLAM [24] default static threshold (8/255) to Otsu thresholding with parameters manually adjusted to qualitatively improve the segmentation. The Otsu procedure was found to remove some unstained tissue and artifacts, while also missing some small areas of stained tissue which may have contained diagnostically relevant information (Figure 1).

Normalisation and augmentation techniques control data variability, which is particularly important for generalisability in histopathology, where varied staining and scanning procedures between labs result in chromatic differences [13]. Normalisation reduces variability, adjusting images into a consistent colour space to allow models to learn general features. We investigated two commonly used [47] stain normalisation techniques - Reinhard normalisation [48] and Macenko normalisation [49]. These approaches work in logarithmic colour spaces, where stains behave linearly, making them

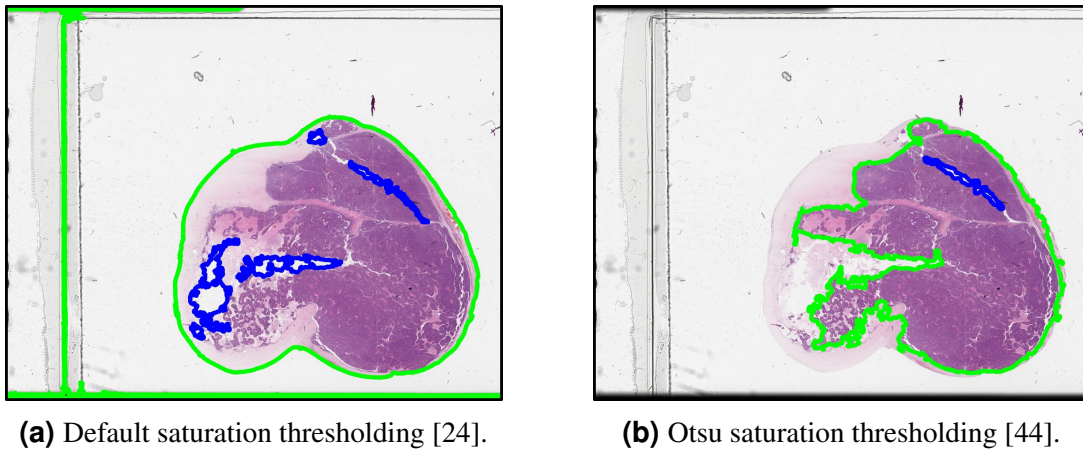


Figure 1. Examples of two saturation thresholding tissue segmentation approaches, with green outlines indicating tissue and blue outlines indicating holes within the tissue regions. This example contains a coverslip edge which is incorrectly identified as foreground by the default approach, and a small amount of stained tissue which is excluded by Otsu thresholding.

easier to separate and manipulate. Reinhard normalisation is a standard normalisation technique applied in $l\alpha\beta$ space (radiance l , blue-yellow α , red-green β). Macenko normalisation uses singular value decomposition to separate stain and saturation values, before scaling stain values in logarithmic RGB space. Basic RGB normalisations were also applied to all images (after any other colour adjustments) to match the ImageNet and histopathology-specific pretraining procedures. While many more sophisticated stain normalisation techniques have been developed, it is unclear whether any such approach is better than Macenko normalisation overall [47].

Augmentation techniques conversely increase the variability of the training data to allow the model to learn a more general domain. For such large images, training end-to-end to allow for online data augmentation (adjustments during training) is extremely computationally intensive [50]. Some researchers have attempted to apply online augmentations in the embedding space using generative models [26, 51], though this adds an extra layer of complexity to an already resource-intensive model pipeline. Instead, offline augmentation creates a finite set of augmented versions of the original data, artificially increasing the diversity of training data to a lesser extent than online augmentation. We investigated colour augmentations which adjusted the brightness, contrast, saturation and hue of each patch using parameters from a previous study [52], which we found to give plausible altered colours (Figure 2).

Preprocessing techniques were compared using a baseline ImageNet-pretrained ResNet50 encoder with the default feature extraction settings from the CLAM repository [24] (static saturation thresholding with no augmentation or normalisation). Comparisons were made using Reinhard normalisation, Macenko normalisation, Otsu thresholding, Otsu thresholding with Macenko normalisation, and colour augmentations to increase the effective training set size by factors of 5x, 10x, and 20x.

Feature extractor architectures were compared using ResNet50 [23], ResNet18 [23], and a large vision transformer (ViT-L) [53], all pretrained using ImageNet [28].

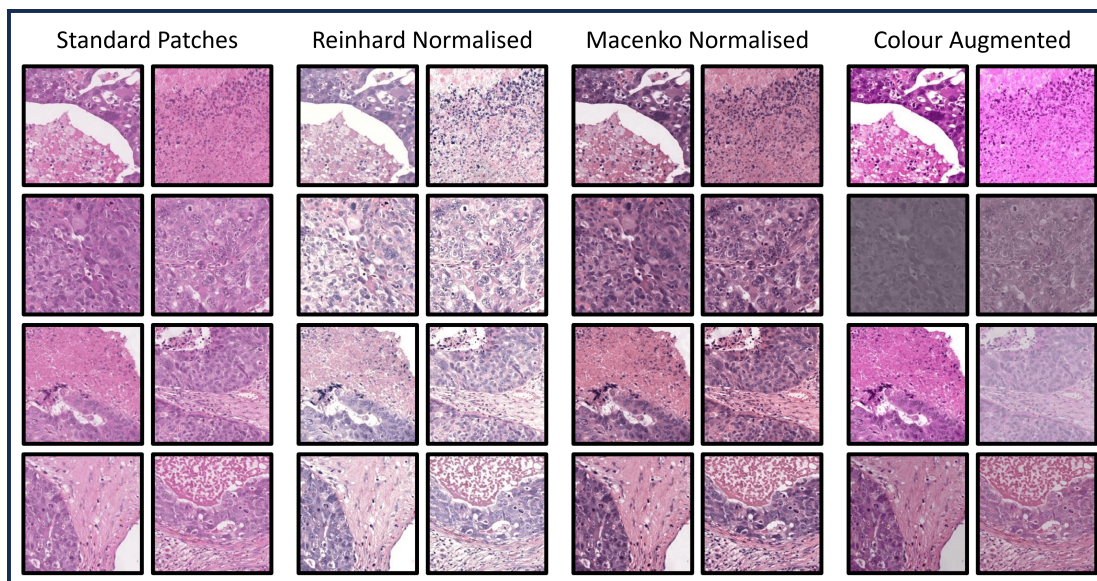


Figure 2. Tissue normalisation and augmentation procedures illustrated using 256x256 pixel patches from a single whole slide image at 10x magnification.

The ResNet50 outputs were taken from the end of the third residual block (as in CLAM [24]) to give 1024 features per input patch. The ResNet18 does not have a layer this large, so 512 features were extracted from the end of the fourth residual block instead. ViT-L was applied without a final fully connected layer to give 1024 features per patch. ViT-L was the largest feature extractor, with 303M model parameters compared to only 9M for ResNet50 and 11M for ResNet18.

Pretraining datasets were compared using the ResNet18 and ViT-L models, each with an ImageNet-pretrained and a histopathology-pretrained version. ImageNet-pretraining for ResNet models had been conducted using the original 1,000 class ImageNet dataset alone, whereas, the ViT-L was first trained on the much larger set of nearly 22,000 classes, and then fine-tuned to the same set of 1,000 classes. The reported ImageNet classification accuracies were 80.9%, 69.8%, and 85.1% for ResNet50, ResNet18, and ViT-L, respectively.

Where the generic feature extractors had been pretrained through supervised classification, the histopathology domain-specific feature extractor pretraining was unsupervised, leveraging large quantities of diverse data without the need for extensive labelling. The ResNet18 was trained with 25,000 WSIs from 57 datasets (including 1376 ovarian cancer WSIs) [31], using the contrastive self-supervised method SimCLR [54]. The domain-specific ViT-L model ‘UNI’ was trained with 100,000 WSIs from 3 data centres (including 5796 gynaecological WSIs) [39], using the newer self-supervised method DINOv2 [55]. While the domain-specific ResNet18 was one of the first histopathology foundation models to be published, UNI was the largest accessible foundation model at the time of this study.

Tuning and Evaluation

The downstream performance of each feature extraction method was evaluated through an attention-based multiple instance learning (ABMIL) [56] five-class ovarian cancer

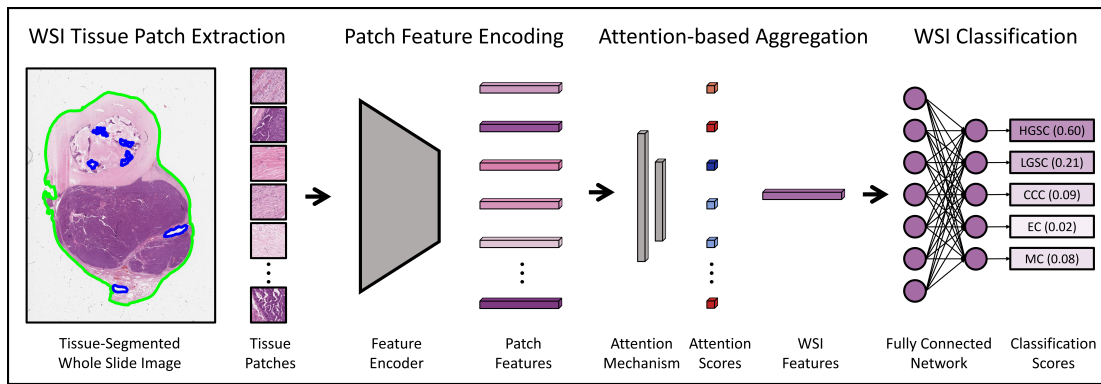


Figure 3. Attention-based multiple instance learning (ABMIL) [56] model pipeline for ovarian cancer subtyping, showing the classification of a high-grade serous carcinoma slide processed without normalisation or augmentation.

subtype classifier (Figure 3). In ABMIL, the patch features were passed through a trainable attention layer which assigned each patch an attention score (between 0 and 1) representing the relative importance of the patch in downstream classification. An attention-weighted average of the patch features generated a WSI feature set, which was classified through a fully connected neural network with one output node per class given the multi-class setting. The outputs were passed through the softmax function to generate the (uncalibrated) classification probabilities for each subtype, with the maximum taken as the predicted class.

For internal data, non-overlapping 1024x1024 pixel tissue patches were extracted at the native 40x tissue magnification, then downsampled to 256x256 pixels at 10x apparent magnification, which was previously found to be optimal using the ResNet50 encoder [14]. For external data, 512x512 pixel tissue patches were extracted and downsampled to achieve the same 256x256 pixels at 10x apparent magnification. Patch features were extracted using these 256x256 patches for all models except the domain-specific ResNet18 and both vision transformers, where they were first resized to 224x224 pixels to match the model pretraining strategies.

Classifiers were tuned using iterative grid search where typically two hyperparameters were adjusted at a time, with the best taken forward to the next iteration. Ten hyperparameters were optimised on the average loss of the five-fold validation sets. Seven hyperparameters directly influenced the Adam optimiser [57], controlling the learning rate, learning rate decay proportion and patience, first and second moment decay, optimisation stability, and L2 regularisation rate. The remaining hyperparameters controlled the model size (the dimension of the attention layer and subsequent fully connected layer), and the proportions of parameter dropout and data dropout during training. Models were trained using a balanced cross-entropy loss and class-weighted sampling to help account for the class imbalance in the training set. Over 150 unique hyperparameter configurations were evaluated during the tuning of each model.

Models were evaluated using the macro F1 score, balanced accuracy, and macro-averaged area under the receiver operating characteristic curve (AUROC). These metrics assessed different aspects of classification performance, with AUROC giving a holistic but imbalanced overview of discriminative power, F1 giving a balanced measure of

predictive performance at a specific threshold, and balanced accuracy representing realistic clinical performance. Stratified five-fold cross-validation (split 60-20-20 train-val-test at the case level to avoid data leakage) was employed during training. In hold-out testing and external validation, the predictions of the five cross-validation models were averaged to generate an ensembled classification. All results were reported using the mean and 95% confidence intervals from 10,000 iterations of bootstrapping.

Paired t-tests were used to test for statistically significant differences between results across the five cross-validation folds. The baseline and non-baseline preprocessing types within each model were compared, along with comparisons between each baseline model architecture, with p-values adjusted for multiple testing using a false discovery rate correction when more than two models were compared [58]. Results were considered *statistically significant* given an adjusted p-value < 0.05 . This manuscript was prepared following the TRIPOD+AI (Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis + Artificial Intelligence) checklist [59] to ensure thorough reporting, with the completed checklist available in Appendix D. The PyTorch-based code used in this study is available at https://github.com/scjbb/Ovarian_Features. Experiments were conducted using an NVIDIA A100 GPU and 32 AMD EPYC7742 CPUs @3.4GHz.

RESULTS

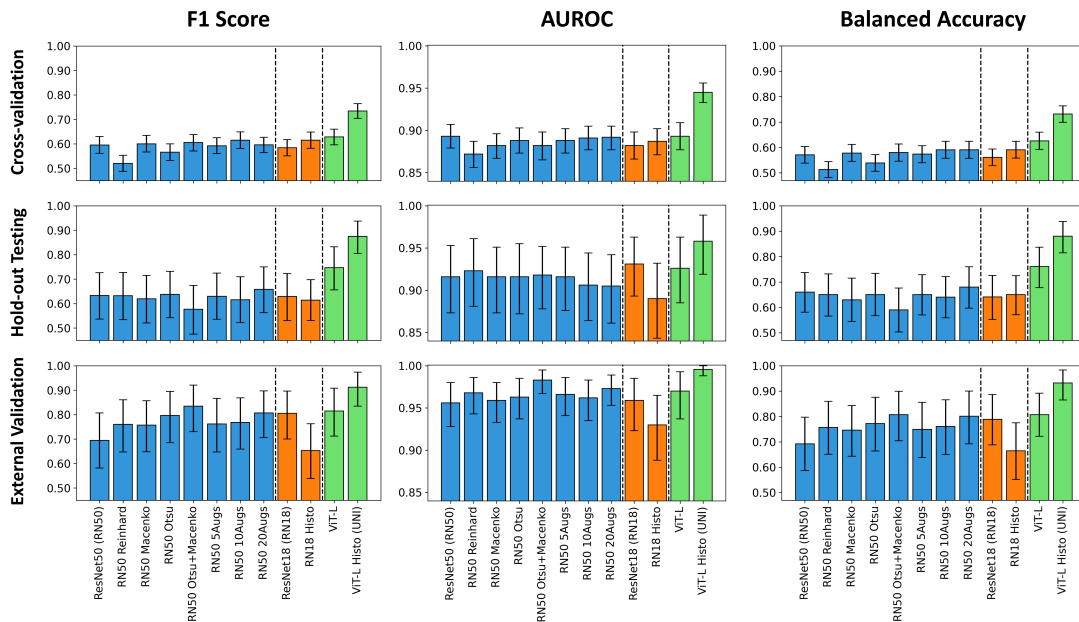


Figure 4. Results of the three validations for each model (mean and 95% confidence interval generated by 10,000 iterations of bootstrapping). Blue indicates ResNet50 (RN50)-based models, orange indicates ResNet18 (RN18)-based models, and green indicates vision transformer (ViT-L)-based models. Hold-out testing and external validation results are based on an ensemble of cross-validation models. Precise values are tabulated in Appendix B.

The highest performance across every metric in each validation was achieved by the histopathology-pretrained vision transformer, UNI. Results are visualised in Figure 4 and tabulated in Appendix 4, and p-values are tabulated in Appendix C. In cross-validation, UNI achieved an F1 score of 0.734 (95% CI: 0.704-0.764), AUROC of 0.945 (0.933-0.956), and a balanced accuracy of 73.2% (69.9-76.4%). In hold-out testing, this improved to an F1 score of 0.875 (0.805-0.937), AUROC of 0.957 (0.919-0.989), and a balanced accuracy of 88.0% (81.5-93.8%). In external testing, UNI achieved an F1 score of 0.912 (0.835-0.974), AUROC of 0.996 (0.988-1.000), and a balanced accuracy of 93.2% (86.5-98.3). Confusion matrices for the UNI model are shown in Tables 2-4. The optimal hyperparameters from tuning each model are reported in Appendix A.

		Predicted Subtype				
		HGSC	LGSC	CCC	EC	MC
Actual Subtype	HGSC	1165	46	28	25	2
	LGSC	39	43	7	3	0
	CCC	29	10	154	3	2
	EC	21	4	2	173	9
	MC	1	0	4	28	66

Table 2. Confusion matrix for predictions of the histopathology-pretrained vision transformer foundation model, UNI, in five-fold cross-validation. Correct classifications are shown in **bold**.

		Predicted Subtype				
		HGSC	LGSC	CCC	EC	MC
Actual Subtype	HGSC	18	0	0	2	0
	LGSC	0	14	2	2	2
	CCC	3	0	17	0	0
	EC	1	0	0	19	0
	MC	0	0	0	0	20

Table 3. Confusion matrix for predictions of the five-fold ensembled histopathology-pretrained vision transformer foundation model, UNI, in the internal test set. Correct classifications are shown in **bold**.

Preprocessing Techniques

In internal validations, varied preprocessing techniques had only a modest effect on the classification performance. In cross-validation, the ResNet50 baseline achieved an F1 score of 0.596, AUROC of 0.893, and a balanced accuracy of 57.1%. No pre-processing method improved the F1 score or balanced accuracy by more than 0.02, and no method improved AUROC at all. Performance was notably decreased by Reinhard normalisation and by Otsu thresholding, with the F1 score falling by 0.076 and 0.030 respectively.

In hold-out testing, the baseline model achieved an F1 score of 0.634, AUROC of 0.916, and a balanced accuracy of 66.0%. Most other models had a slightly lower F1 score and balanced accuracy, and a similar AUROC. Only the 20x colour augmentation

		Predicted Subtype				
		HGSC	LGSC	CCC	EC	MC
Actual Subtype	HGSC	27	0	1	2	0
	LGSC	0	9	0	0	0
	CCC	0	1	19	0	0
	EC	0	0	0	10	1
	MC	0	0	0	1	9

Table 4. Confusion matrix for predictions of the five-fold ensemble histopathology-pretrained vision transformer foundation model, UNI, in the external test set. Correct classifications are shown in **bold**.

appeared to improve overall performance, increasing F1 by 0.023 and balanced accuracy by 0.020, but reducing AUROC by 0.012.

In external validation, every preprocessing method improved performance over the baseline. From a baseline F1 score of 0.696, AUROC of 0.956, and balanced accuracy of 69.2%, all methods increased F1 score and balanced accuracy by over 0.05, and AUROC by over 0.002. The greatest external performances were found by combining Otsu thresholding with Macenko normalisation and by 20x colour augmentations, which each increased baseline performance by over 0.1 F1 score and balanced accuracy, and over 0.016 AUROC. None of the performance differences were statistically significant in any validation.

Model Architectures

The vision transformer gave the best baseline F1 score and balanced accuracy across all validations, often by a wide margin. For example, in the hold-out test set, ViT-L achieved an F1 score of 0.747 and a balanced accuracy of 76.0%, compared to the next best (ResNet50) baseline scores of 0.634 and 66.0%, respectively. AUROC scores were similar between model architectures, with each baseline model performing best for one of the three evaluation datasets. Performance was similar between the ResNet18 and ResNet50 baselines for internal data (within 0.02 by each metric), but ResNet18 generalised much better to the external data (F1 score 0.804 vs 0.696, AUROC 0.959 vs 0.956, balanced accuracy 79.0% vs 69.2%, respectively). The differences in performance were not statistically significant, except between ViT-L and ResNet50 on the hold-out F1 score ($p = 0.001$), and between ViT-L and both ResNet50 ($p = 0.004$) and ResNet18 ($p = 0.034$) on the hold-out balanced accuracy.

Pretraining Datasets

Histopathology-pretraining slightly improved the ResNet18 performance in cross-validation, but made performance worse in hold-out testing and external validation compared to the baseline ResNet18. The only statistically significant difference was found in the hold-out testing AUROC, where the histopathology-pretrained model was worse than the baseline model. The domain-specific ResNet18 did show the same trend of improvement on hold-out and external data compared to cross-validation data, but to a lesser extent than other models.

Histopathology-pretraining drastically improved performance for the vision trans-

former, with UNI giving the best performance by a wide margin in every experiment. The improvement over the baseline ViT-L was statistically significant for all metrics in cross-validation (p-values: 0.010 F1, 0.004 AUROC, 0.028 BalAcc) and hold-out testing (p-values: 0.002 F1, <0.001 AUROC, 0.003 BalAcc), and for two metrics in external validation (p-values: 0.009 AUROC, 0.031 BalAcc).

DISCUSSION

In this study, we thoroughly compared the effects of different feature extraction strategies for the classification of ovarian carcinoma subtypes. The results indicated that the UNI foundation model features were much better than non-domain-specific and ResNet-based features. In contrast, the performance of the ResNet18 foundation model was relatively poor and was often outperformed by the generic (ImageNet-pretrained) ResNet18 feature set. The performance of the UNI model was particularly impressive considering that it was trained using 20x magnification data and applied here at 10x, whereas the ResNet18 foundation model training did include some 10x magnification data.

Different preprocessing techniques often had little impact on internal performance (likely due to the homogeneity of the single-centre dataset), but aided the generalisability to external data. Reinhard and Macenko normalisations gave similar performances, though Reinhard normalisation slightly reduced cross-validation performance. There was a modest positive trend between the number of augmentations used and the resulting model performance which may continue beyond the x20 augmentations used herein, though this may not be worth the considerable computational burden since the normalisation approaches gave similar performance to 20 augmentations. While normalisation and augmentation techniques improved the model generalisability of the standard ImageNet-pretrained ResNet50 model, the UNI foundation model still consistently outperformed such models. These normalisation and augmentation techniques may be able to further improve the UNI foundation model performance.

Performance was generally higher in hold-out testing than in cross-validation, and was higher still in external validation. The five-fold ensemble used in hold-out and external testing may have been beneficial, though given the relatively small size of these test sets, the results are uncertain. Thus, part of the difference in performance between datasets may be attributed to random chance. Cross-validation performance may have also been limited by the inclusion of IDS samples, which exhibit variable diagnostic quality due to chemotherapy-induced visual changes, such as increased cellular damage and reduced tumour size. The improvement in performance on hold-out and external data was relatively modest for the histopathology-pretrained ResNet18, which may indicate that it had focused more on the IDS samples and had overfit to the training domain. It is unclear why the highest performance was on the external dataset across all experiments. One explanation is that this set only contained one representative slide per patient and the slides were largely devoid of artifacts. This may reflect a sampling bias in the external test set, ensuring high data quality and representing a best-case scenario.

The results of this study are similar to the only other studies to use a large ovarian cancer subtyping dataset (both using 948 WSIs) [20, 21]. One presented a multi-scale graph model [21] and reported an optimal cross-validation F1 score of 0.69 and balanced accuracy of 73%, compared to our 0.73 and 73%, respectively. The

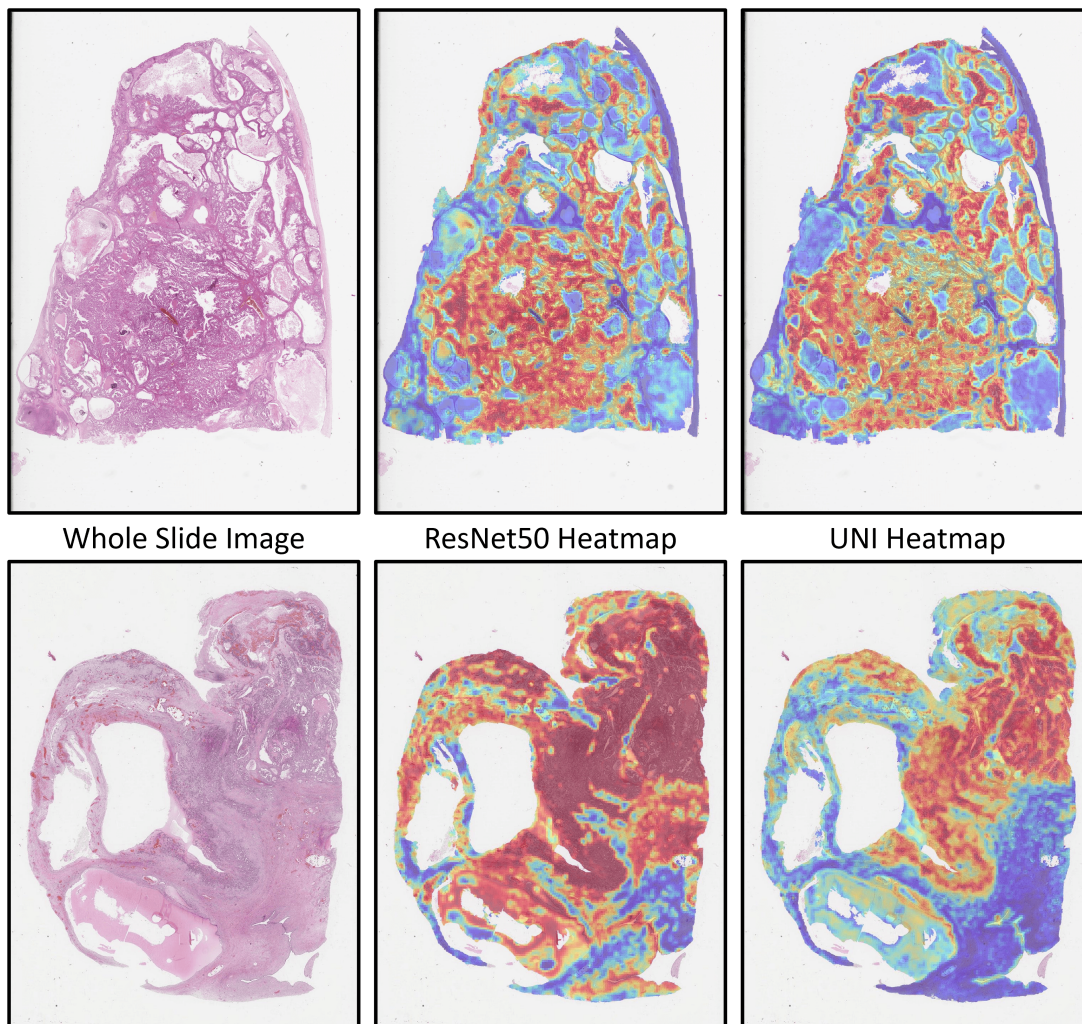


Figure 5. Attention heatmaps from the ABMIL classifier using ImageNet-pretrained ResNet50 features and histopathology-pretrained vision transformer (UNI) features. The upper example shows a typical difference between heatmaps with different diagnoses, and the lower example shows the most extreme qualitative difference found between heatmaps in the internal test set. In both examples, the UNI classification was correct (upper - mucinous, lower - clear cell carcinoma) , and the ResNet50 classification was incorrect (upper - endometrioid, lower - mucinous carcinoma). The heatmaps are based on 256x256 pixel patches with 50% overlap at 10x apparent magnification, with patch sizes only appearing different between slides due to the variable size of resection samples.

other study evaluated four MIL approaches and reported an optimal cross-validation F1 score of 0.79, AUROC of 0.95, and balanced accuracy of 81%, compared to our 0.73, 0.95, and 73%, respectively. This study included external validation using an ensemble of cross-validation models on 60 WSIs, and reported an F1 score of 0.81, AUROC of 0.96, and balanced accuracy of 80%, compared to our 0.91, 1.00, and 93%, respectively. These comparisons are provided for context and should not be considered to be conclusive given the differences in the datasets used. Both previous studies used ImageNet-pretrained encoders and in one [20], tumour segmentation labels were given for most slides, including all external validation slides. A sparsity of publicly available data has limited external validations in previous research [11], though future work may benefit from the 513 ovarian carcinoma WSIs recently released as part of the OCEAN challenge [60]. The external test set performance reported in this study is the highest achieved to date for ovarian carcinoma subtyping [11].

To analyse the differences between UNI and the baseline ResNet50 models, two pathologists (KA and NMO) qualitatively compared the ABMIL attention heatmaps (Figure 5). Most heatmaps were well-focused on tumour and relevant stromal regions for both models, with often only subtle differences between them. The UNI-based heatmaps generally indicated a slightly greater focus on tumour, where the ResNet50 model also paid attention to some stromal regions of variable diagnostic relevance (Appendix E). Attention heatmaps can be useful for identifying potential sources of error but should be interpreted with caution since they cannot provide a complete explanation of classification decisions [61].

The Reinhard normalisation procedure gave the poorest cross-validation performance. As shown in Figure 6, this method was particularly affected by the presence of artifacts in training slides as the normalisation would make these areas appear similar to tissue, making it harder for the attention mechanism to effectively discard these patches. The impact was smaller on Macenko normalisation, which often avoided applying stain colouring to plain background regions. Colour augmentations were unaffected as they could not introduce staining to non-tissue patches (see Figure 6).

The improved subtype classification accuracy from UNI is promising for the potential clinical utility of these models, though more extensive external validation is required to ensure that these models generalise to all relevant sources of variation, especially across different histopathology labs and slide scanners. This should include robustness to lower quality data and artifacts to reduce the burden of quality control, which otherwise risks sacrificing any time savings the models could provide to diagnostic services. Furthermore, it is currently unclear how best to present automatically generated information to pathologists to assist them, rather than to distract, frustrate, or confuse them. This may require improved model interpretability and a measure of model uncertainty, especially considering the existence of rare subtypes which are difficult to collect sufficient data on.

Ideally, algorithms would be made more computationally efficient for use in the clinic, but the vision transformer is much less computationally efficient than the ResNet. This problem is exacerbated by the limited digitisation of histopathology services, with most pathological diagnoses still made under a microscope. AI adoption will be contingent on it being accessible and beneficial given limited computational infrastructure and users who may not be technological experts. While there are various issues inhibiting

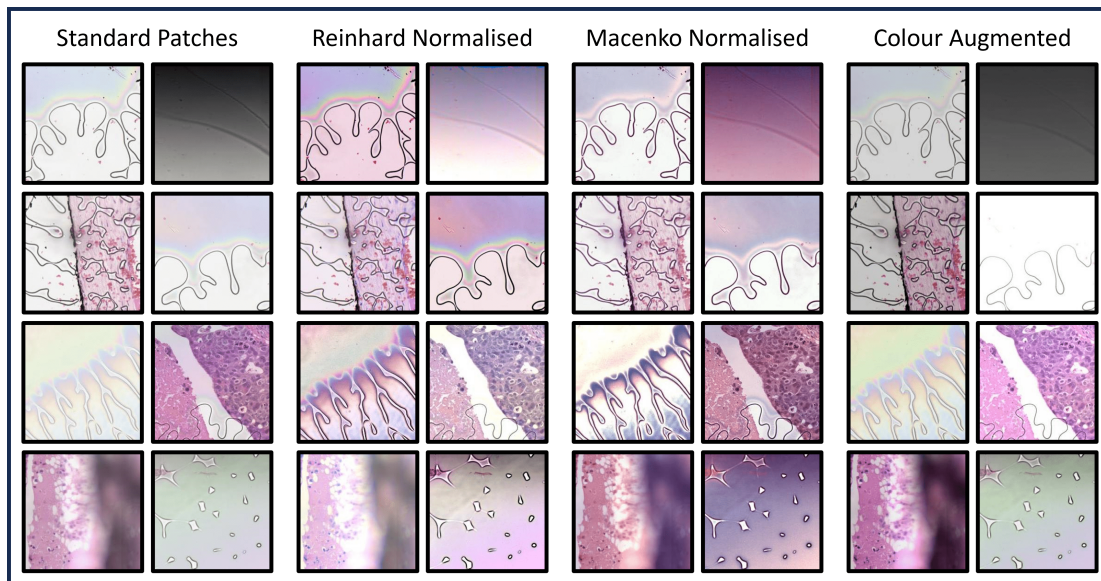


Figure 6. Tissue normalisation and augmentation procedures illustrated using 256x256 pixel patches containing artifacts from the same slide at 10x magnification shown in Figure 2. The normalisation procedures erroneously apply staining to many non-tissue regions, where the colour augmentations are much less affected. This can be seen most clearly in patches where no tissue is present, such as the bottom-right patch in each group.

the clinical translation of ovarian cancer subtyping models, these seem increasingly likely to be overcome in the near future.

CONCLUSION

In this study, we conducted a rigorous validation of feature extraction methods for ovarian cancer subtyping. We found that the features generated by a recently published histopathology-pretrained vision transformer ‘foundation model’, UNI, drastically improved downstream classification performance when compared to a generic transformer or ResNet feature extractor. Several different data preprocessing techniques were evaluated, and while these improved the generalisability of the ResNet models, they were not sufficient to match the performance of the foundation model. Through a five-fold ensemble of ABMIL classifiers, UNI achieved a five-class balanced accuracy of 88% on internal test data and 93% on external test data, compared to 68% and 81% respectively for the best ResNet-based models. While this improved performance may be sufficient for clinical implementation, the need to address logistical hurdles and conduct larger-scale validations remains.

ACKNOWLEDGEMENTS

There was no direct funding for this research. JB is supported by the UKRI Engineering and Physical Sciences Research Council (EPSRC) [EP/S024336/1]. KA is supported by the Tony Bramall Charitable Trust. The funders had no role

in influencing the content of this research. There was no formal study protocol or registration, and no patient or public involvement in this research. This study was conducted retrospectively using human subject data and received approval from the Wales Research Ethics Committee [18/WA/0222] and the Confidentiality Advisory Group [18/CAG/0124]. Approval has not yet been provided for this data to be shared outside of the research group. External data were downloaded from <https://www.medicalimageanalysis.com/data/ovarian-carcinomas-histopathology-dataset> (last accessed 09/04/24). All code used in this research is available at https://github.com/scjbb/Ovarian_Features, along with extended results and details of hyperparameter tuning. For the purpose of open access, the author has applied a Creative Commons Attribution (CC BY) licence to any Author Accepted Manuscript version arising from this submission.

AUTHOR CONTRIBUTIONS

JB created the study protocol with feedback and contributions from all other authors. KA collected data with assistance from NMO. JB conducted all experiments with advice from NR. JB wrote the manuscript, with feedback and contributions from all other authors.

COMPETING INTERESTS

NMO's fellowship is funded by 4D Path. All other authors declare no conflicts of interest.

A HYPERPARAMETERS

The optimal hyperparameters (Table 5) typically did not vary greatly for models using the same feature extraction architecture, with a few notable exceptions. The classifier based on the x5 augmented training data was the smallest by far (and had the smallest stability parameter and learning rate decay factor), with only 0.1M parameters compared to the next smallest at 0.7M (the baseline, Macenko, and x20 augmented ResNet50 models). The larger sets of augmented data gave similar hyperparameters to the other ResNet50 models. The Reinhard and Macenko models used fewer patches per slide (400) than the other models (600-1000), though the combined Otsu Macenko model used a standard quantity (1000). The Otsu model had a much lower dropout (0.1) than other ResNet models (0.3-0.6), though again the combined Otsu Macenko model had a standard amount (0.3). The regularisation (weight decay) parameter was much greater in the baseline ViT-L model (0.1) than in any other model, including the histopathology-pretrained ViT-L model (0.001). Other hyperparameters were relatively stable within a given feature extractor.

Some hyperparameters varied greatly between model architectures, especially the learning rate and dropout. The learning rate was much smaller for ViT-L models (0.00001-0.00005) than ResNet50 models (0.001-0.002), and decayed much faster (every 10 epochs rather than every 15-30 epochs). The optimal dropout was 0 for both ViT-L models, where it was 0.3-0.6 for most ResNet-based models. The Adam optimiser first and second moment decay parameters were also typically higher in ViT-L and ResNet18 models than in ResNet50 models. Other hyperparameters were relatively consistent between model architectures. The size of the classifiers ranged from 0.1M (x5 augmentations) to 1.6M parameters (baseline ResNet18), with the optimal UNI-based classifier using 0.8M. These ABMIL classifiers were orders of magnitude smaller than the feature extraction models.

Model	Learning Rate	Weight Decay	First Moment Decay	Second Moment Decay	Stability Parameter	LR Decay Patience	LR Decay Factor	Model Size	Drop Out	Max Patches
ResNet50 (RN50)	2e-3	1e-3	0.75	0.95	1e-2	20	0.75	[512,128]	0.4	800
RN50 Reinhard	2e-3	1e-3	0.75	0.95	1e-2	25	0.75	[512,256]	0.4	400
RN50 Macenko	2e-3	1e-3	0.85	0.95	1e-2	15	0.75	[512,128]	0.3	400
RN50 Otsu	2e-3	1e-3	0.75	0.95	1e-2	15	0.9	[512,256]	0.1	600
RN50 Otsu+Macenko	2e-3	1e-4	0.75	0.99	1e-3	25	0.9	[512,256]	0.3	1000
RN50 5Augs	1e-3	1e-4	0.8	0.99	1e-4	25	0.6	[128,32]	0.4	700
RN50 10Augs	2e-3	1e-3	0.8	0.99	1e-2	20	0.75	[512,256]	0.4	700
RN50 20Augs	1e-3	1e-4	0.7	0.999	1e-3	20	0.75	[512,128]	0.6	1000
ResNet18 (RN18)	1e-4	1e-5	0.8	0.99	1e-4	20	0.9	[1024,256]	0.5	700
RN18 Histo	2e-4	1e-4	0.9	0.99	1e-4	20	0.9	[512,512]	0.6	1000
ViT-L	5e-5	1e-1	0.85	0.999	1e-3	10	0.35	[512,384]	0.0	800
ViT-L Histo (UNI)	1e-5	1e-3	0.9	0.999	1e-5	10	0.75	[512,256]	0.0	1000

Table 5. The final hyperparameters of each model determined by an iterative grid search tuning procedure using five cross-validation folds. The model size is presented as the number of parameters in the attention layer and subsequent fully connected layer. RN18 - ResNet18. RN50 - ResNet50. ViT-L - Vision Transformer.

B RESULTS TABLES

Model Architecture	Preprocessing Approach	F1 Score	AUROC	Balanced Accuracy
ResNet50	Baseline	0.596 (0.561-0.630)	<u>0.893</u> (0.879-0.907)	57.1% (53.8-60.4%)
	Reinhard Normalisation	0.520 (0.488-0.553)	0.872 (0.856-0.887)	51.3% (48.2-54.4%)
	Macenko Normalisation	0.601 (0.567-0.635)	0.882 (0.867-0.896)	57.8% (54.5-61.2%)
	Otsu Thresholding	0.566 (0.532-0.600)	0.888 (0.873-0.903)	53.9% (50.6-57.2%)
	Otsu + Macenko	0.605 (0.571-0.638)	0.882 (0.865-0.898)	58.0% (54.6-61.4%)
	5x Colour Augmentation	0.592 (0.560-0.625)	0.888 (0.873-0.902)	57.4% (54.0-60.7%)
	10x Colour Augmentation	<u>0.615</u> (0.581-0.649)	0.891 (0.877-0.905)	<u>59.1%</u> (55.7-62.4%)
20x Colour Augmentation	0.596 (0.564-0.627)	0.892 (0.877-0.905)	<u>59.1%</u> (55.7-62.4%)	
ResNet18	Baseline	0.584 (0.551-0.617)	0.882 (0.866-0.898)	56.1% (52.8-59.4%)
	Histo-pretraining	<u>0.615</u> (0.582-0.648)	<u>0.887</u> (0.871-0.902)	<u>59.1%</u> (55.8-62.4%)
ViT-L	Baseline	0.628 (0.596-0.660)	0.893 (0.877-0.909)	62.6% (59.2-66.0%)
	Histo-pretraining (UNI)	0.734 (0.704-0.764)	0.945 (0.933-0.956)	73.2% (69.9-76.4%)

Table 6. Results of five-fold cross-validation. Results are reported as the mean and 95% confidence intervals (in brackets) from 10,000 iterations of bootstrapping. The greatest overall results are shown in **bold** and the greatest results within each model architecture are underlined for each metric.

Model Architecture	Preprocessing Approach	F1 Score	AUROC	Balanced Accuracy
ResNet50	Baseline	0.634 (0.537-0.726)	0.916 (0.873-0.953)	66.0% (58.1-73.7%)
	Reinhard Normalisation	0.632 (0.534-0.727)	<u>0.923</u> (0.881-0.961)	65.0% (56.6-73.2%)
	Macenko Normalisation	0.620 (0.521-0.715)	0.915 (0.873-0.951)	63.0% (54.4-71.5%)
	Otsu Thresholding	0.637 (0.542-0.732)	0.916 (0.872-0.955)	65.0% (56.7-73.4%)
	Otsu + Macenko	0.577 (0.475-0.674)	0.918 (0.878-0.952)	59.0% (50.3-67.6%)
	5x Colour Augmentation	0.630 (0.536-0.725)	0.916 (0.876-0.951)	65.0% (57.0-72.9%)
	10x Colour Augmentation	0.616 (0.522-0.710)	0.906 (0.864-0.944)	64.0% (55.9-72.1%)
20x Colour Augmentation	<u>0.657</u> (0.563-0.750)	0.904 (0.861-0.942)	<u>68.0%</u> (59.7-76.0%)	
ResNet18	Baseline	<u>0.628</u> (0.530-0.723)	<u>0.930</u> (0.893-0.963)	64.0% (55.3-72.6%)
	Histo-pretraining	0.613 (0.531-0.698)	0.890 (0.843-0.932)	<u>65.0%</u> (57.1-72.5%)
ViT-L	Baseline	0.747 (0.656-0.832)	0.926 (0.885-0.963)	76.0% (67.8-83.7%)
	Histo-pretraining (UNI)	0.875 (0.805-0.937)	0.957 (0.919-0.989)	88.0% (81.5-93.8%)

Table 7. Results of hold-out testing, with predictions generated by an ensemble of the five-fold classification models. Results are reported as the mean and 95% confidence intervals (in brackets) from 10,000 iterations of bootstrapping. The greatest overall results are shown in **bold** and the greatest results within each model architecture are underlined for each metric.

Model Architecture	Preprocessing Approach	F1 Score	AUROC	Balanced Accuracy
ResNet50	Baseline	0.696 (0.582-0.807)	0.956 (0.928-0.980)	69.2% (58.7-79.7%)
	Reinhard Normalisation	0.761 (0.647-0.861)	0.968 (0.943-0.986)	75.8% (65.1-86.0%)
	Macenko Normalisation	0.756 (0.648-0.857)	0.959 (0.933-0.980)	74.5% (64.3-84.3%)
	Otsu Thresholding	0.797 (0.685-0.895)	0.963 (0.937-0.985)	77.2% (66.4-87.6%)
	Otsu + Macenko	<u>0.834</u> (0.730-0.921)	<u>0.983</u> (0.967-0.995)	<u>80.5%</u> (70.4-89.9%)
	5x Colour Augmentation	0.762 (0.647-0.866)	0.966 (0.941-0.986)	74.9% (63.8-85.6%)
	10x Colour Augmentation	0.768 (0.659-0.869)	0.962 (0.935-0.983)	76.1% (65.0-86.6%)
20x Colour Augmentation	0.806 (0.706-0.897)	0.973 (0.953-0.989)	80.0% (69.2-90.0%)	
ResNet18	Baseline	<u>0.804</u> (0.700-0.896)	<u>0.959</u> (0.923-0.985)	<u>79.0%</u> (68.8-88.6%)
	Histo-pretraining	0.653 (0.539-0.763)	0.930 (0.888-0.965)	66.5% (55.2-77.5%)
ViT-L	Baseline	0.814 (0.712-0.908)	0.970 (0.937-0.993)	80.7% (72.2-89.2%)
	Histo-pretraining (UNI)	<u>0.912</u> (0.835-0.974)	<u>0.996</u> (0.988-1.000)	<u>93.2%</u> (86.5-98.3%)

Table 8. Results of external validation, with predictions generated by an ensemble of the five-fold classification models. Results are reported as the mean and 95% confidence intervals (in brackets) from 10,000 iterations of bootstrapping. The greatest overall results are shown in **bold** and the greatest results within each model architecture are underlined for each metric.

C RESULTS OF HYPOTHESIS TESTING

Evaluation p-values	Model 1	Model 2	F1 Score	AUROC	Balanced Accuracy
Cross-Validation	Baseline	Reinhard	0.379	0.467	0.282
	Baseline	Macenko	0.959	0.616	0.945
	Baseline	Otsu	0.635	0.775	0.282
	Baseline	Otsu+Macenko	0.739	0.154	0.844
	Baseline	5Augs	0.959	0.649	0.945
	Baseline	10Augs	0.635	0.649	0.660
	Baseline	20Augs	0.959	0.979	0.707
Hold-out Testing	Baseline	Reinhard	0.794	0.553	0.833
	Baseline	Macenko	0.794	0.553	0.833
	Baseline	Otsu	0.794	0.553	0.894
	Baseline	Otsu+Macenko	0.794	0.658	0.833
	Baseline	5Augs	0.794	0.579	0.833
	Baseline	10Augs	0.972	0.553	0.925
	Baseline	20Augs	0.794	0.553	0.833
External Validation	Baseline	Reinhard	0.968	0.637	0.938
	Baseline	Macenko	0.968	0.715	0.970
	Baseline	Otsu	0.968	0.785	0.800
	Baseline	Otsu+Macenko	0.643	0.354	0.689
	Baseline	5Augs	0.968	0.354	0.970
	Baseline	10Augs	0.968	0.637	0.970
	Baseline	20Augs	0.318	0.354	0.553

Table 9. Resulting p-values from paired t-tests comparing the subtype classification results from five cross-validation folds using different preprocessing techniques using a ResNet50 feature extractor. False discovery rate p-value adjustments were applied to account for multiple testing [58].

Evaluation p-values	Model 1	Model 2	F1 Score	AUROC	Balanced Accuracy
Cross-Validation	ResNet50	ResNet18	0.736	0.731	0.736
	ResNet18	ViT-L	0.351	0.731	0.160
	ViT-L	ResNet50	0.181	0.773	0.081
Hold-out Testing	ResNet50	ResNet18	0.297	0.068	0.563
	ResNet18	ViT-L	0.062	0.219	0.034
	ViT-L	ResNet50	0.001	0.071	0.004
External Validation	ResNet50	ResNet18	0.374	0.773	0.418
	ResNet18	ViT-L	0.196	0.128	0.208
	ViT-L	ResNet50	0.052	0.128	0.056

Table 10. Resulting p-values from paired t-tests comparing the subtype classification results from five cross-validation folds using different feature extraction architectures. False discovery rate p-value adjustments were applied to account for multiple testing [58]. p-values from model architecture comparison. Values less than 0.05 are shown in **bold**.

Evaluation p-values	Model 1	Model 2	F1 Score	AUROC	Balanced Accuracy
Cross-Validation	ResNet18	ResNet18 Histo	0.499	0.103	0.322
	ViT-L	ViT-L Histo (UNI)	0.010	0.004	0.028
Hold-out Testing	ResNet18	ResNet18 Histo	0.562	0.015	0.865
	ViT-L	ViT-L Histo (UNI)	0.002	<0.001	0.003
External Validation	ResNet18	ResNet18 Histo	0.101	0.056	0.123
	ViT-L	ViT-L Histo (UNI)	0.101	0.009	0.031

Table 11. Resulting p-values from paired t-tests comparing the subtype classification results from five cross-validation folds using different preprocessing techniques using a ResNet50 feature extractor. Values less than 0.05 are shown in **bold**.

D TRIPOD+AI REPORTING CHECKLIST



Version: 11-January-2024

Section/Topic	Item	Development / evaluation ¹	Checklist item	Reported on page
TITLE				
<i>Title</i>	1	D;E	Identify the study as developing or evaluating the performance of a multivariable prediction model, the target population, and the outcome to be predicted	1
ABSTRACT				
<i>Abstract</i>	2	D;E	See TRIPOD+AI for Abstracts checklist	1
INTRODUCTION				
<i>Background</i>	3a	D;E	Explain the healthcare context (including whether diagnostic or prognostic) and rationale for developing or evaluating the prediction model, including references to existing models	2-3
	3b	D;E	Describe the target population and the intended purpose of the prediction model in the context of the care pathway, including its intended users (e.g., healthcare professionals, patients, public)	2
	3c	D;E	Describe any known health inequalities between sociodemographic groups	2
<i>Objectives</i>	4	D;E	Specify the study objectives, including whether the study describes the development or validation of a prediction model (or both)	3
METHODS				
<i>Data</i>	5a	D;E	Describe the sources of data separately for the development and evaluation datasets (e.g., randomised trial, cohort, routine care or registry data), the rationale for using these data, and representativeness of the data	3-4
	5b	D;E	Specify the dates of the collected participant data, including start and end of participant accrual; and, if applicable, end of follow-up	3-4
<i>Participants</i>	6a	D;E	Specify key elements of the study setting (e.g., primary care, secondary care, general population) including the number and location of centres	3-4
	6b	D;E	Describe the eligibility criteria for study participants	3-4
	6c	D;E	Give details of any treatments received, and how they were handled during model development or evaluation, if relevant	3-4
<i>Data preparation</i>	7	D;E	Describe any data pre-processing and quality checking, including whether this was similar across relevant sociodemographic groups	3-5
<i>Outcome</i>	8a	D;E	Clearly define the outcome that is being predicted and the time horizon, including how and when assessed, the rationale for choosing this outcome, and whether the method of outcome assessment is consistent across sociodemographic groups	3
	8b	D;E	If outcome assessment requires subjective interpretation, describe the qualifications and demographic characteristics of the outcome assessors	3-4
	8c	D;E	Report any actions to blind assessment of the outcome to be predicted	N/A
<i>Predictors</i>	9a	D	Describe the choice of initial predictors (e.g., literature, previous models, all available predictors) and any pre-selection of predictors before model building	3-4
	9b	D;E	Clearly define all predictors, including how and when they were measured (and any actions to blind assessment of predictors for the outcome and other predictors)	3-4
	9c	D;E	If predictor measurement requires subjective interpretation, describe the qualifications and demographic characteristics of the predictor assessors	N/A
<i>Sample size</i>	10	D;E	Explain how the study size was arrived at (separately for development and evaluation), and justify that the study size was sufficient to answer the research question. Include details of any sample size calculation	3-4
<i>Missing data</i>	11	D;E	Describe how missing data were handled. Provide reasons for omitting any data	N/A
<i>Analytical methods</i>	12a	D	Describe how the data were used (e.g., for development and evaluation of model performance) in the analysis, including whether the data were partitioned, considering any sample size requirements	6-8
	12b	D	Depending on the type of model, describe how predictors were handled in the analyses (functional form, rescaling, transformation, or any standardisation).	4-6
	12c	D	Specify the type of model, rationale ² , all model-building steps, including any hyperparameter tuning, and method for internal validation	6-8
	12d	D;E	Describe if and how any heterogeneity in estimates of model parameter values and model performance was handled and quantified across clusters (e.g., hospitals, countries). See TRIPOD-Cluster for additional considerations ³	N/A
	12e	D;E	Specify all measures and plots used (and their rationale) to evaluate model performance (e.g., discrimination, calibration, clinical utility) and, if relevant, to compare multiple models	7-8
	12f	E	Describe any model updating (e.g., recalibration) arising from the model evaluation, either overall or for particular sociodemographic groups or settings	N/A
	12g	E	For model evaluation, describe how the model predictions were calculated (e.g., formula, code, object, application programming interface)	6-8
<i>Class imbalance</i>	13	D;E	If class imbalance methods were used, state why and how this was done, and any subsequent methods to recalibrate the model or the model predictions	7
<i>Fairness</i>	14	D;E	Describe any approaches that were used to address model fairness and their rationale	N/A
<i>Model output</i>	15	D	Specify the output of the prediction model (e.g., probabilities, classification). Provide details and rationale for any classification and how the thresholds were identified	6-7

¹ D=items relevant only to the development of a prediction model; E=items relating solely to the evaluation of a prediction model; D;E=items applicable to both the development and evaluation of a prediction model

² Separately for all model building approaches.

³ TRIPOD-Cluster is a checklist of reporting recommendations for studies developing or validating models that explicitly account for clustering or explore heterogeneity in model performance (eg, at different hospitals or centres). Debray et al, *BMJ* 2023; 380: e071018 [DOI: 10.1136/bmj-2022-071018]

<i>Training versus evaluation</i>	16	D;E	Identify any differences between the development and evaluation data in healthcare setting, eligibility criteria, outcome, and predictors	3-4
<i>Ethical approval</i>	17	D;E	Name the institutional research board or ethics committee that approved the study and describe the participant-informed consent or the ethics committee waiver of informed consent	15
OPEN SCIENCE				
<i>Funding</i>	18a	D;E	Give the source of funding and the role of the funders for the present study	14-15
<i>Conflicts of interest</i>	18b	D;E	Declare any conflicts of interest and financial disclosures for all authors	15
<i>Protocol</i>	18c	D;E	Indicate where the study protocol can be accessed or state that a protocol was not prepared	15
<i>Registration</i>	18d	D;E	Provide registration information for the study, including register name and registration number, or state that the study was not registered	15
<i>Data sharing</i>	18e	D;E	Provide details of the availability of the study data	15
<i>Code sharing</i>	18f	D;E	Provide details of the availability of the analytical code ⁴	15
PATIENT & PUBLIC INVOLVEMENT				
<i>Patient & Public Involvement</i>	19	D;E	Provide details of any patient and public involvement during the design, conduct, reporting, interpretation, or dissemination of the study or state no involvement.	15
RESULTS				
<i>Participants</i>	20a	D;E	Describe the flow of participants through the study, including the number of participants with and without the outcome and, if applicable, a summary of the follow-up time. A diagram may be helpful.	3-4
	20b	D;E	Report the characteristics overall and, where applicable, for each data source or setting, including the key dates, key predictors (including demographics), treatments received, sample size, number of outcome events, follow-up time, and amount of missing data. A table may be helpful. Report any differences across key demographic groups.	3-4
	20c	E	For model evaluation, show a comparison with the development data of the distribution of important predictors (demographics, predictors, and outcome).	3-4
<i>Model development</i>	21	D;E	Specify the number of participants and outcome events in each analysis (e.g., for model development, hyperparameter tuning, model evaluation)	3-4
<i>Model specification</i>	22	D	Provide details of the full prediction model (e.g., formula, code, object, application programming interface) to allow predictions in new individuals and to enable third-party evaluation and implementation, including any restrictions to access or re-use (e.g., freely available, proprietary) ⁵	4-8
<i>Model performance</i>	23a	D;E	Report model performance estimates with confidence intervals, including for any key subgroups (e.g., sociodemographic). Consider plots to aid presentation.	8
	23b	D;E	If examined, report results of any heterogeneity in model performance across clusters. See TRIPOD Cluster for additional details ³ .	N/A
<i>Model updating</i>	24	E	Report the results from any model updating, including the updated model and subsequent performance	N/A
DISCUSSION				
<i>Interpretation</i>	25	D;E	Give an overall interpretation of the main results, including issues of fairness in the context of the objectives and previous studies	11-14
<i>Limitations</i>	26	D;E	Discuss any limitations of the study (such as a non-representative sample, sample size, overfitting, missing data) and their effects on any biases, statistical uncertainty, and generalizability	11-14
<i>Usability of the model in the context of current care</i>	27a	D	Describe how poor quality or unavailable input data (e.g., predictor values) should be assessed and handled when implementing the prediction model	11-14
	27b	D	Specify whether users will be required to interact in the handling of the input data or use of the model, and what level of expertise is required of users	13
	27c	D;E	Discuss any next steps for future research, with a specific view to applicability and generalizability of the model	13

From: Collins GS, Moons KGM, Dhiman P, et al. *BMJ* 2024;385:e078378. doi:10.1136/bmj-2023-078378

⁴ This relates to the analysis code, for example, any data cleaning, feature engineering, model building, evaluation.

⁵ This relates to the code to implement the model to get estimates of risk for a new individual.

E ATTENTION HEATMAP ANALYSIS

To compare the UNI model to the baseline ResNet50 model, two pathologists (KA and NMO) qualitatively assessed the attention heatmaps for ten class-balanced example WSIs from the internal hold-out test set. These WSIs (shown in Figures 5, 7, 8) were selected from those in which a different classification had been determined by each model (specifically using the first-fold model of the five-model ensemble). Out of 39 total disagreements, the UNI-based model gave the correct classification in 26 cases, the ResNet50-based model in 3 cases, and neither was correct in 10 cases. The pathologists were not told which heatmap corresponded to which model, nor what classifications had been predicted by the models.

Overall, the heatmaps were determined to be similar between models, with both giving high attention to tumour regions and low attention to most stroma regions. Where differences occurred, the ResNet50-based model would typically give high attention to a larger tissue area, often including relevant stromal features (e.g. necrosis and psammoma bodies), but sometimes also including irrelevant stroma. The pathologists expressed a preference for the UNI-based heatmap in four cases and the ResNet50-based heatmap in three cases, with no preference expressed for the remaining three cases due to their overwhelming similarity. In eight of the selected cases, the UNI model had correctly determined the classification, including all three cases in which the pathologists had preferred the ResNet50-based heatmap. In these cases, the UNI model did not appear to give sufficient attention to all relevant tissue, though it still determined the correct classification.

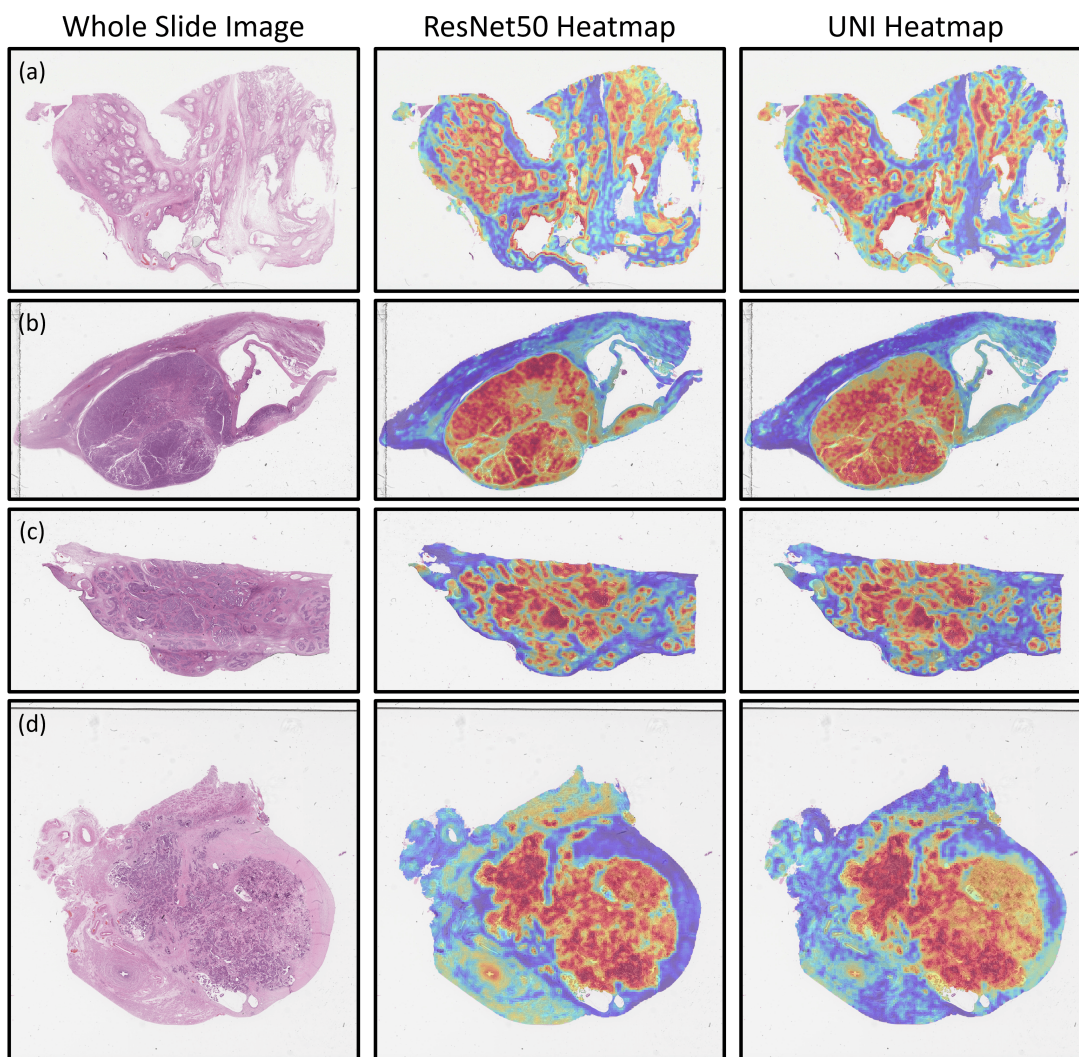


Figure 7. Attention heatmaps from the ABMIL classifier using ImageNet-pretrained ResNet50 features and histopathology-pretrained vision transformer (UNI) features, where the classification differed between the two models. (a) Ground truth: MC, ResNet50: CCC, UNI: MC. (b) Ground truth: CCC, ResNet50: HGSC, UNI: CCC. (c) Ground truth: EC, ResNet50: HGSC, UNI: EC. (d) Ground truth: LGSC, ResNet50: HGSC, UNI: LGSC.

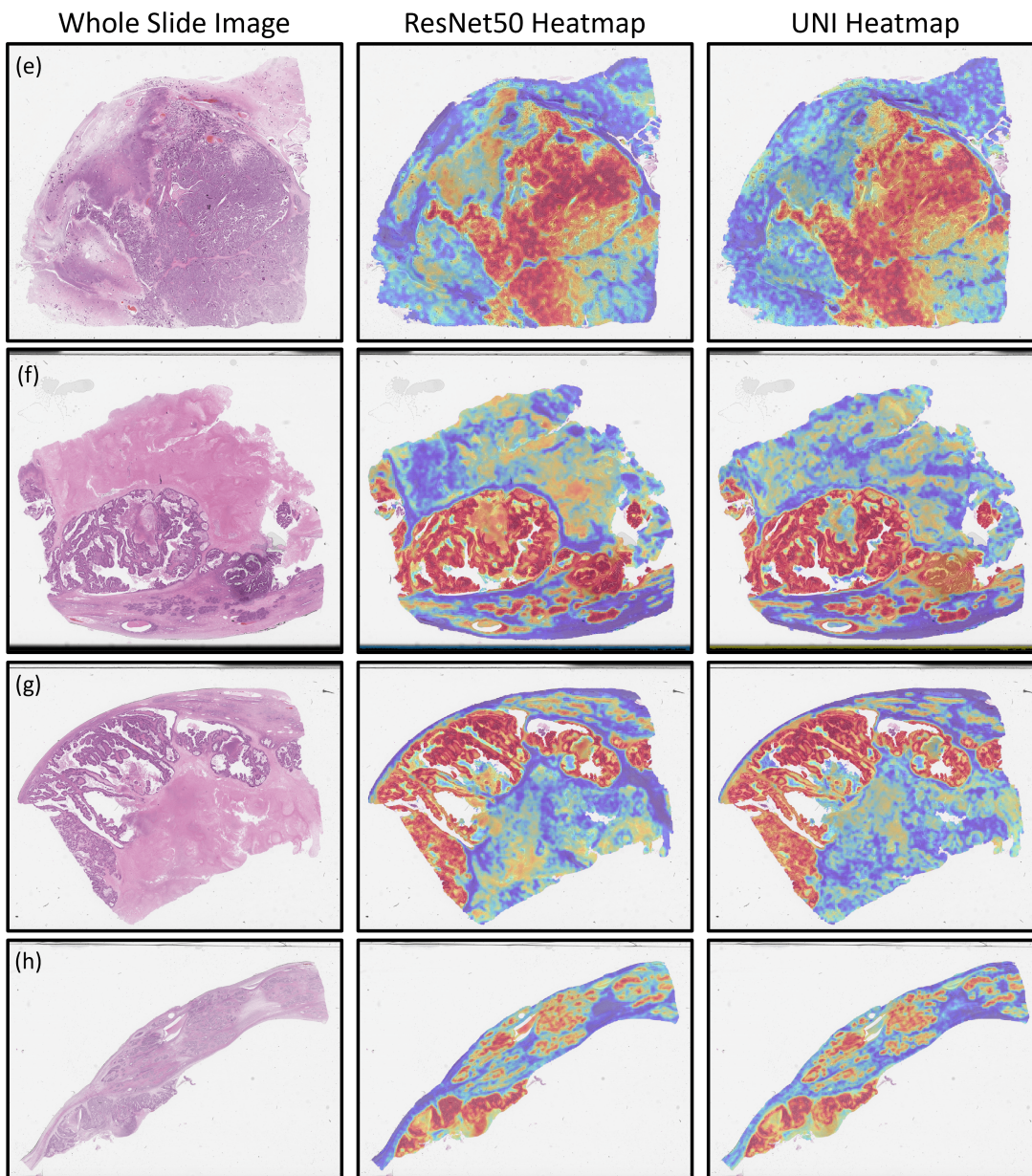


Figure 8. Attention heatmaps from the ABMIL classifier using ImageNet-pretrained ResNet50 features and histopathology-pretrained vision transformer (UNI) features, where the classification differed between the two models. (e) Ground truth: LGSC, ResNet50: HGSC, UNI: LGSC. (f) Ground truth: HGSC, ResNet50: HGSC, UNI: EC. (g) Ground truth: HGSC, ResNet50: HGSC, UNI: EC. (h) Ground truth: EC, ResNet50: HGSC, UNI: EC.

REFERENCES

- [1] Bray F, Laversanne M, Sung H, et al. Global cancer statistics 2022: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA: A Cancer Journal for Clinicians*. 2024;1-35.
- [2] Köbel M, Kalloger SE, Boyd N, et al. Ovarian carcinoma subtypes are different diseases: implications for biomarker studies. *PLoS medicine*. 2008;5(12):e232.
- [3] Peres LC, Cushing-Haugen KL, Köbel M, et al. Invasive epithelial ovarian cancer survival by histotype and disease stage. *JNCI: Journal of the National Cancer Institute*. 2019;111(1):60-8.
- [4] Moch H. Female genital tumours: WHO Classification of Tumours, Volume 4. *WHO Classification of Tumours*. 2020;4.
- [5] Vroobel K. Overview of Ovarian Tumours: Pathogenesis and General Considerations. In: *Pathology of the Ovary, Fallopian Tube and Peritoneum*. Springer; 2024. p. 95-113.
- [6] Köbel M, Bak J, Bertelsen BI, et al. Ovarian carcinoma histotype determination is highly reproducible, and is improved through the use of immunohistochemistry. *Histopathology*. 2014;64(7):1004-13.
- [7] Royal College of Pathologists. Meeting pathology demand: Histopathology workforce census. *RCPATH*; 2018. Available from: <https://www.rcpath.org/static/952a934d-2ec3-48c9-a8e6e00fcdca700f/Meeting-Pathology-Demand-Histopathology-Workforce-Census-2018.pdf>.
- [8] Wilson ML, Fleming KA, Kuti MA, et al. Access to pathology and laboratory medicine services: a crucial gap. *The Lancet*. 2018;391(10133):1927-38.
- [9] Hanna TP, King WD, Thibodeau S, et al. Mortality due to cancer treatment delay: systematic review and meta-analysis. *bmj*. 2020;371.
- [10] Allen KE, Adusumilli P, Breen J, et al. Artificial Intelligence in Ovarian Digital Pathology. In: *Pathology of the Ovary, Fallopian Tube and Peritoneum*. Springer; 2024. p. 731-49.
- [11] Breen J, Allen K, Zucker K, et al. Artificial intelligence in ovarian cancer histopathology: a systematic review. *NPJ Precision Oncology*. 2023;7(1):83.
- [12] Matthews GA, McGenity C, Bansal D, et al. Public evidence on AI products for digital pathology. *medRxiv*. 2024:2024-02.
- [13] Breen J, Allen K, Zucker K, et al. Efficient subtyping of ovarian cancer histopathology whole slide images using active sampling in multiple instance learning. In: Tomaszewski JE, Ward AD, editors. *Medical Imaging 2023: Digital and Computational Pathology*. vol. 12471. International Society for Optics and Photonics. SPIE; 2023. p. 1247110.
- [14] Breen J, Allen K, Zucker K, et al. Reducing Histopathology Slide Magnification Improves the Accuracy and Speed of Ovarian Cancer Subtyping. *arXiv preprint arXiv:231113956*. 2023.
- [15] BenTaieb A, Li-Chang H, Huntsman D, et al. Automatic diagnosis of ovarian carcinomas via sparse multiresolution tissue representation. In: *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part I 18*. Springer;

2015. p. 629-36.
- [16] BenTaieb A, Nosrati MS, Li-Chang H, et al. Clinically-inspired automatic classification of ovarian carcinoma subtypes. *Journal of pathology informatics*. 2016;7(1):28.
 - [17] BenTaieb A, Li-Chang H, Huntsman D, et al. A structured latent model for ovarian carcinoma subtyping from histopathology slides. *Medical image analysis*. 2017;39:194-205.
 - [18] Levine AB, Peng J, Farnell D, et al. Synthesis of diagnostic quality cancer pathology images by generative adversarial networks. *The Journal of pathology*. 2020;252(2):178-88.
 - [19] Boschman J, Farahani H, Darbandsari A, et al. The utility of color normalization for ai-based diagnosis of hematoxylin and eosin-stained pathology images. *The Journal of Pathology*. 2022;256(1):15-24.
 - [20] Farahani H, Boschman J, Farnell D, et al. Deep learning-based histotype diagnosis of ovarian carcinoma whole-slide pathology images. *Modern Pathology*. 2022;35(12):1983-90.
 - [21] Mirabadi AK, Archibald G, Darbandsari A, et al. GRASP: GRAPh-Structured Pyramidal Whole Slide Image Representation. *arXiv preprint arXiv:240203592*. 2024.
 - [22] Gadermayr M, Tschuchnig M. Multiple instance learning for digital pathology: A review of the state-of-the-art, limitations & future potential. *Computerized Medical Imaging and Graphics*. 2024:102337.
 - [23] He K, Zhang X, Ren S, et al. Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*; 2016. p. 770-8.
 - [24] Lu MY, Williamson DF, Chen TY, et al. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature biomedical engineering*. 2021;5(6):555-70.
 - [25] Shao Z, Bian H, Chen Y, et al. Transmil: Transformer based correlated multiple instance learning for whole slide image classification. *Advances in neural information processing systems*. 2021;34:2136-47.
 - [26] Zaffar I, Jaume G, Rajpoot N, et al. Embedding space augmentation for weakly supervised learning in whole-slide images. In: *2023 IEEE 20th International Symposium on Biomedical Imaging (ISBI)*. IEEE; 2023. p. 1-4.
 - [27] Godson L, Alemi N, Nsengimana J, et al. Immune subtyping of melanoma whole slide images using multiple instance learning. *Medical Image Analysis*. 2024;93:103097.
 - [28] Russakovsky O, Deng J, Su H, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*. 2015;115:211-52.
 - [29] Köbel M, Kalloger SE, Baker PM, et al. Diagnosis of ovarian carcinoma cell type is highly reproducible: a transcanadian study. *The American journal of surgical pathology*. 2010;34(7):984-93.
 - [30] Vorontsov E, Bozkurt A, Casson A, et al. Virchow: A Million-Slide Digital Pathology Foundation Model. *arXiv preprint arXiv:230907778*. 2023.
 - [31] Ciga O, Xu T, Martel AL. Self supervised contrastive learning for digital histopathology. *Machine Learning with Applications*. 2022;7:100198.
 - [32] Chen RJ, Chen C, Li Y, et al. Scaling vision transformers to gigapixel images via

- hierarchical self-supervised learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition; 2022. p. 16144-55.
- [33] Wang X, Yang S, Zhang J, et al. Transformer-based unsupervised contrastive learning for histopathological image classification. *Medical image analysis*. 2022;81:102559.
- [34] Kang M, Song H, Park S, et al. Benchmarking self-supervised learning on diverse pathology datasets. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition; 2023. p. 3344-54.
- [35] Filiot A, Ghermi R, Olivier A, et al. Scaling self-supervised learning for histopathology with masked image modeling. *medRxiv*. 2023:2023-07.
- [36] Wang W, Ma S, Xu H, et al. When an image is worth 1,024 x 1,024 words: A case study in computational pathology. *arXiv preprint arXiv:231203558*. 2023.
- [37] Azizi S, Culp L, Freyberg J, et al. Robust and data-efficient generalization of self-supervised machine learning for diagnostic imaging. *Nature Biomedical Engineering*. 2023;7(6):756-79.
- [38] Campanella G, Kwan R, Fluder E, et al. Computational Pathology at Health System Scale—Self-Supervised Foundation Models from Three Billion Images. *arXiv preprint arXiv:231007033*. 2023.
- [39] Chen RJ, Ding T, Lu MY, et al. Towards a general-purpose foundation model for computational pathology. *Nature Medicine*. 2024:1-13.
- [40] Dippel J, Feulner B, Winterhoff T, et al. RudolfV: A Foundation Model by Pathologists for Pathologists. *arXiv preprint arXiv:240104079*. 2024.
- [41] Achiam J, Adler S, Agarwal S, et al. Gpt-4 technical report. *arXiv preprint arXiv:230308774*. 2023.
- [42] Touvron H, Martin L, Stone K, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:230709288*. 2023.
- [43] Allen KE, Breen J, Hall G, et al. #900 Comparative evaluation of ovarian carcinoma subtyping in primary versus interval debulking surgery specimen whole slide images using artificial intelligence. *International Journal of Gynecologic Cancer*. 2023;33(Suppl 3):A429-30.
- [44] Otsu N. A threshold selection method from gray-level histograms. *Automatica*. 1975;11(285-296):23-7.
- [45] Janowczyk A, Zuo R, Gilmore H, et al. HistoQC: an open-source quality control tool for digital pathology slides. *JCO clinical cancer informatics*. 2019;3:1-7.
- [46] Shakhawat H, Hossain S, Kabir A, et al. Review of artifact detection methods for automated analysis and diagnosis in digital pathology. In: *Artificial Intelligence For Disease Diagnosis And Prognosis In Smart Healthcare*. CRC Press; 2023. p. 177-202.
- [47] Breen J, Zucker K, Allen K, et al. Generative Adversarial Networks for Stain Normalisation in Histopathology. In: *Applications of Generative AI*. Cham: Springer International Publishing; 2024. p. 227-47.
- [48] Reinhard E, Adhikhmin M, Gooch B, et al. Color transfer between images. *IEEE Computer graphics and applications*. 2001;21(5):34-41.
- [49] Macenko M, Niethammer M, Marron JS, et al. A method for normalizing histology slides for quantitative analysis. In: *2009 IEEE international symposium on biomedical imaging: from nano to macro*. IEEE; 2009. p. 1107-10.

- [50] Dooper S, Pinckaers H, Aswolinskiy W, et al. Gigapixel end-to-end training using streaming and attention. *Medical Image Analysis*. 2023;88:102881.
- [51] Shao Z, Dai L, Wang Y, et al. Augdiff: Diffusion based feature augmentation for multiple instance learning in whole slide image. *arXiv preprint arXiv:230306371*. 2023.
- [52] Wang Y, Farnell D, Farahani H, et al. Classification of epithelial ovarian carcinoma whole-slide pathology images using deep transfer learning. *arXiv preprint arXiv:200510957*. 2020.
- [53] Dosovitskiy A, Beyer L, Kolesnikov A, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:201011929*. 2020.
- [54] Chen T, Kornblith S, Norouzi M, et al. A simple framework for contrastive learning of visual representations. In: *International conference on machine learning*. PMLR; 2020. p. 1597-607.
- [55] Oquab M, Darcet T, Moutakanni T, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:230407193*. 2023.
- [56] Ilse M, Tomczak J, Welling M. Attention-based deep multiple instance learning. In: *International conference on machine learning*. PMLR; 2018. p. 2127-36.
- [57] Kingma DP, Ba J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:14126980*. 2014.
- [58] Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*. 1995;57(1):289-300.
- [59] Collins GS, Moons KG, Dhiman P, et al. TRIPOD+ AI statement: updated guidance for reporting clinical prediction models that use regression or machine learning methods. *bmj*. 2024;385.
- [60] Asadi-Aghbolaghi M, Farahani H, Zhang A, et al. Machine Learning-driven Histotype Diagnosis of Ovarian Carcinoma: Insights from the OCEAN AI Challenge. *medRxiv*. 2024:2024-04.
- [61] Bibal A, Cardon R, Alfter D, et al. Is attention explanation? an introduction to the debate. In: *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*; 2022. p. 3889-900.