

An attention-based multi-resolution model for prostate whole slide image classification and localization

Jiayun Li

University of California, Los Angeles
jiayunli@g.ucla.edu

Arkadiusz Gertych

Cedars-Sinai Medical Center
Arkadiusz.Gertych@cshs.org

William Speier

University of California, Los Angeles
speier@ucla.edu

Wenyuan Li

University of California, Los Angeles
liwenyuan.zju@gmail.com

Beatrice S. Knudsen

Cedars-Sinai Medical Center
Beatrice.Knudsen@cshs.org

Corey W. Arnold

University of California, Los Angeles
cwarnold@ucla.edu

Abstract

*Histology review is often used as the ‘gold standard’ for disease diagnosis. Computer aided diagnosis tools can potentially help improve current pathology workflows by reducing examination time and interobserver variability. Previous work in cancer grading has focused mainly on classifying pre-defined regions of interest (ROIs), or relied on large amounts of fine-grained labels. In this paper, we propose a two-stage attention-based multiple instance learning model for slide-level cancer grading and weakly-supervised ROI detection and demonstrate its use in prostate cancer. Compared with existing Gleason classification models, our model goes a step further by utilizing visualized saliency maps to select informative tiles for fine-grained grade classification. The model was primarily developed on a large-scale whole slide dataset consisting of 3,521 prostate biopsy slides with only slide-level labels from 718 patients. The model achieved state-of-the-art performance for prostate cancer grading with an accuracy of 85.11% for classifying benign, low-grade (Gleason grade 3+3 or 3+4), and high-grade (Gleason grade 4+3 or higher) slides on an independent test set.*¹

1. Introduction

Prostate cancer is the most common and second deadliest cancer in men in the U.S, accounting for nearly 1 in 5 new cancer diagnoses [38]. Gleason grading of biopsied tissue

is a key component in patient management and treatment selection [6, 41]. The Gleason score (GS) is determined by the two most prevalent Gleason patterns in the tissue section. Gleason patterns range from 1 (G1), representing tissue that is close to normal glands, to 5 (G5), indicating more aggressive cancer. Patients with high risk cancer (*i.e.* $GS > 7$ or $G4 + G3$) are usually treated with radiation, hormonal therapy, or radical prostatectomy, while those with low- to intermediate-risk prostate cancer (*i.e.* $GS < 6$ or $G3 + G4$) are candidates for active surveillance.

Currently, pathologists need to scan through a histology slide, searching for relevant regions on which to ascertain Gleason scores. This process can be time-consuming and prone to observer variability [19, 23, 17]. Therefore, computer aided diagnosis (CAD) tools can benefit clinical practice by identifying relevant regions and serving as a second reader. However, there are many unique challenges in developing CAD tools for whole slide images (WSIs), such as the very large image size, the heterogeneity of slide contents, the insufficiency of fine-grained labels, and possible artifacts caused by pen markers and stain variations.

In this paper, we developed an attention-based multiple instance learning (MIL) model that can not only predict slide-level labels, but also provide visualization of relevant regions using inherent attention maps. Unlike previous work that relied on labor intensive labels, such as manually drawn regions of interest (ROIs) around glands, our model is trained using only slide-level labels, known as weak labels, which can be easily retrieved from pathology reports. In our proposed two-stage model, suspicious regions are detected at a lower resolution (*e.g.* 5x), and further analyzed at a higher resolution (*e.g.* 10x), which is similar to patholo-

¹This paper appears at CVPR 2019 Towards Causal, Explainable and Universal Medical Visual Diagnosis (MVD) Workshop.

gists’ diagnostic process. To the best of our knowledge, this is the first work that utilizes weakly-supervised attention maps and MIL to select ROIs and classify prostate biopsy slides. Our model was trained and validated on a dataset of 2,661 biopsy slides from 491 patients. The model achieved state-of-the-art performance, with a classification accuracy of 85.11% on a held-out testset consisting of 860 slides from 227 patients.

2. Related Work

ROI-level classification. Early work on WSI analysis mainly focused on classifying small ROIs, which usually were selected by pathologists from the large tissue slide [12, 14, 31]. However, this does not accurately reflect the true clinical task as to ensure completeness, pathologists must grade the entire tissue section rather than sub-selected representative ROIs. This makes models based on ROIs unsuitable for automated Gleason grading [11].

Slide-level classification. Instead of relying on ROIs, more recent research has focused on slide-level classification. Nagpal *et al.* [30] developed a two-stage Gleason classification model. In the first-stage, a tile-level classifier was trained with over 112 million annotated tiles from prostatectomy slides. In the second stage, predictions from the first stage were summarized to a K-nearest neighbor classifier for Gleason scoring. They achieved an average accuracy of 70% in four-class Gleason group classification (1, 2, 3, or 4-5). However, these methods [11, 30, 44] required a well-trained tile-level classifier, which can only be developed on a dataset with manually drawn ROIs or slides with homogeneous tissue contents. Moreover, they did not incorporate information embedded in slide-level labels.

To address these challenges, previous work has proposed using an MIL framework for WSI classification [9], where the slide was represented as a bag and tiles within the bag were modeled as instances in the bag [15, 4, 20]. MIL models can be roughly divided into two types [1, 20]: instance-based [33, 29, 34] and bag-based [2, 5, 10]. Bag-based methods project instance features into low-dimensional representations and often demonstrate superior performance for bag-level classification tasks [20, 1]. However, as bag-level methods lack the ability to predict instance-level labels, they are less interpretable and thus sub-optimal for problems where obtaining instance labels is important [25, 18, 37, 32]. Ilse *et al.* [20] proposed an attention-based deep learning model that can achieve comparable performances to bag-level models without losing interpretability. A low-dimensional instance embedding, an attention mechanism for aggregating instance-level features, and a final bag-level classifier were all parameterized with a neural network. They applied the model on two histology datasets consisting of small tiles extracted from WSIs and demonstrated promising performance. However, they did not ap-

ply the model on larger and more heterogeneous WSIs. Also, attention maps were only used for a visualization method. Campanella *et al.* [4] applied a instance-level MIL model for binary prostate biopsy slide classification (*i.e.* cancer versus non-cancer). Their model was developed on a large dataset consisting of 12,160 biopsy slides, and achieved over 95% area under the curve of the receiver operating characteristic (AUROC). Yet, they did not address the more difficult grading problem. Built upon the attention-based MIL model [20], our model further improves the attention mechanism with instance dropout [40]. Instead of only using the attention map for visualization, we utilize it to automatically localize informative areas, which then get analyzed at higher resolution for cancer grading.

3. Methods

Our model is trained with slide-level annotations in an MIL framework. Specifically, k $N \times N$ tiles $x_i, i \in [1, k]$ can be extracted from a given WSI, which usually contains gigabytes of pixels. Different from supervised computer vision models, in which the label for each tile is provided, only the label for the slide (*i.e.* the set of tiles) is available. This problem can be modeled with MIL by considering tiles as instances and the entire slide as a bag. In section 3.1, we introduce the deep attention-based MIL model [20] and the instance dropout method [40]. The attention-based informative tile selection method is discussed in section 3.2. The overview of our two-stage classification model is described in 3.3.

3.1. Attention-based multiple instance learning

The attention-based MIL model uses a convolutional neural network (CNN) as the backbone to extract instance-level features. An attention module $f(\cdot)$ is added before the softmax classifier to learn weight distribution $\alpha = \alpha_1, \alpha_2, \dots, \alpha_k$ for k instances, which indicates importance of k instances for predicting the current bag-level label y . The $f(\cdot)$ can be modeled by a multilayer perceptron (MLP). If we denote a set of d dimensional feature vectors from k instances as $\mathbf{V} \in \mathbb{R}^{k \times d}$, the attention for the i th instance can be defined in Eq(1),

$$\alpha_i = \text{Softmax}[U^T(\tanh(\mathbf{W}_v \mathbf{v}_i^T))] \quad (1)$$

where $U \in \mathbb{R}^{h \times n}$ and $\mathbf{W} \in \mathbb{R}^{h \times d}$ are learnable parameters, n is the number of classes, and h is the dimension of the hidden layer. Then each tile can have a corresponding attention value learned from the module. Bag-level embedding can be obtained by multiplying learned attentions with instance features.

The attention distribution provides a way to localize informative tiles for the current model prediction. However, the attention-based MIL method suffers from the same

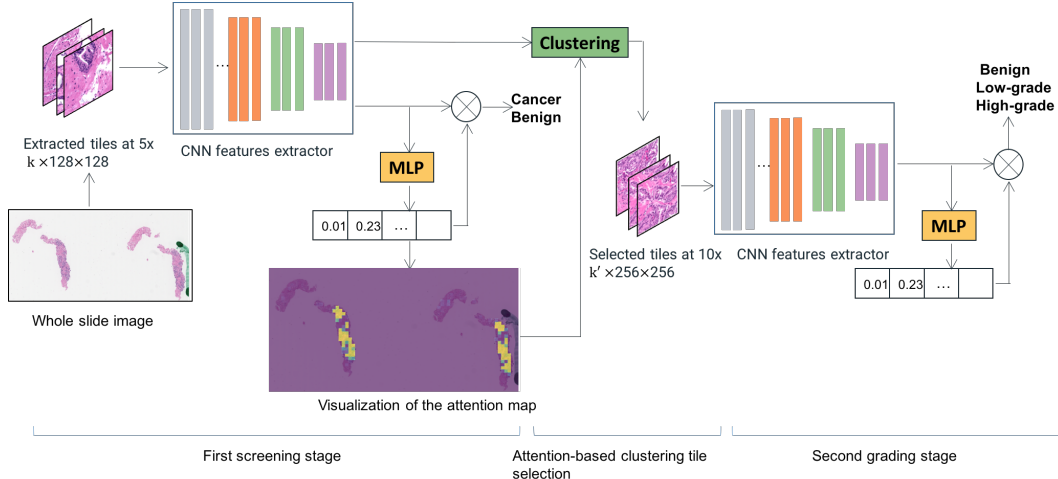


Figure 1. The overview of our two-stage attention-based whole slide image classification model. The first stage is trained with tiles at 5x for cancer versus non-cancer classification. Informative tiles identified using instance features and attention maps from the first stage are selected to be analyzed in the second stage at a higher resolution for cancer grading.

problem as many saliency detection models [43, 40, 16, 27]. Specifically, the model may only focus on the most discriminative input instead of all relevant regions. This problem may not have a large effect on the bag-level classification. Nevertheless, it could affect the integrity of the attention map and therefore affect the performance of our second stage model. To address this challenge, we utilize a similar method to [40]. During training, we randomly drop different instances in the bag by setting their pixel values to the mean RGB value of the training dataset [40]; in testing all instances will be used. This method forces the network to discover more relevant instances instead of only relying on the most discriminative ones.

3.2. Attention-based region selection

In section 3.1, we describe the attention-based MIL model and how to improve the learned attention map using instance dropout. One intuitive way to select informative tiles with attention maps is to rank them by attention values and select the top k percentile. However, this method is highly reliant upon the quality of the learned attention maps, which may not be perfect, especially when there is no explicit supervision.

To address this problem, we incorporate information from instance feature vectors \mathbf{V} . Specifically, instances are clustered into n clusters based on instance features. Principle component analysis (PCA) is applied to reduce the dimension of features before clustering. Thus, instances that share similar semantic features will be grouped together. The average attention value for cluster i with m tiles can be computed $\bar{\alpha}_i = \frac{1}{m} \sum_{i=1}^n \alpha_i$ and normalized so that $\bar{\alpha}$ sums to 1. The intuition is that clusters with higher aver-

age attention are more likely to contain relevant information for slide classification (*e.g.* given a cancerous slide, clusters containing stroma or benign glands should have lower attention values compared with those containing cancerous regions). Based on this, the number of tiles to be selected from each cluster can be determined by the total number of tiles and the average attention of the cluster.

3.3. Two-stage whole slide image classification model

In this section, we discuss how to incorporate the aforementioned methods into a two-stage WSI classification model. WSIs often contain several gigabytes of pixels, which practically impossible to fit into GPU memory. However, most regions on the WSIs are stroma or benign glands, which do not contribute to the final diagnosis. In clinical practice, pathologists usually scan through an entire slide at low magnification (*e.g.* 5x), identify areas that may contain cancer, and closely examine these regions at a higher magnification (*e.g.* 10x or 20x).

Inspired by the workflow of pathologists, we developed a two-stage classification model. In the first screening stage, 128×128 tiles are extracted from each slide at 5x magnification and fed into a binary MIL model for cancer versus non-cancer classification. Informative tiles are identified by using attention maps and instance features from the 5x model as described in 3.2. Then, the second grading stage model uses selected tiles at 10x to classify the slide into benign, low-grade, or high-grade prostate cancer. Selected tiles are at the same location, but at a higher resolution as those in the first screening stage. Figure 1 shows the overview of our two-stage WSI classification model.

4. Experiment

In section 4.1, we introduce the dataset and the preprocessing pipeline used. Details about model implementation and training are discussed in section 4.2.

4.1. Dataset

Cedars Sinai dataset. CNN feature extractors for both stages were pre-trained with a relatively small dataset with manually drawn ROIs from the Department of Pathology at Cedars-Sinai Medical Center (IRB approval numbers: Pro00029960 and Pro00048462) [13, 21, 26, 28]. The dataset contains two parts. 1) 513 tiles of size 1200×1200 extracted from prostatectomies of 40 patients, which contain low-grade pattern (Gleason grade 3), high-grade pattern (Gleason grade 4 and 5), benign (BN), and stromal areas. These tiles were annotated by pathologists at the pixel-level. 2) 30 WSIs from prostatectomies of 30 patients. These slides were annotated by a pathologist who circled and graded the major foci of tumor as either low-grade, high-grade, or BN areas.

The scanning objective for all slides and tiles was set at 20x ($0.5 \mu\text{m}$ per pixel). To use this dataset for tile classification, we randomly sampled 11,595 tiles of size 256×256 at 10x from annotated regions. We will refer this dataset as the tile-level dataset in the following sections.

UCLA dataset. The MIL model is further trained with a large-scale dataset with only slide-level annotations. The dataset contains prostate biopsy slides from the Department of Pathology and Laboratory Medicine at the University of California, Los Angeles (UCLA). We randomly sampled a balanced number of low-grade, high-grade, and benign cases, resulting in 3,521 slides from 718 patients. We randomly divided the dataset based on patients for model training, validation, and testing to ensure the same patient would not be included in both training and testing. Labels for these slides were retrieved from pathology reports. For simplicity, we will refer this dataset as the slide-level dataset in the following sections.

Data preprocessing. Since WSIs may contain a lot of background regions and pen marker artifacts, we converted the slide at the lowest available magnification into HSV color space and thresholded on the hue channel to generate a mask for tissue areas. Morphological operations such as dilation and erosion were applied to fill in small holes and remove isolated points from tissue masks. Then, a set of instances (*i.e.* tiles) for one bag (*i.e.* slide) of size 256×256 at 10x was extracted from the grid with 12.5% overlap. Tiles that contained less than 80% tissue regions were removed from analysis. The number of tiles in the majority of slides ranged from 100 to 300. The same color normalization algorithm [35] was performed on tiles from both UCLA and Cedars Sinai datasets. Tiles at 10x were downsampled to 5x for the first stage of model training.

4.2. Implementation Details

Blue ratio selection. Most previous work on WSI classification utilizes the blue ratio image to select relevant regions [7, 3, 24]. The blue ratio image as defined in Eq(2) reflects the concentration of the blue color, so it can detect regions with the most nuclei.

$$\text{BR} = \frac{100 \times B}{1 + R + G} \times \frac{256}{1 + R + G + B} \quad (2)$$

where R , G , B are the red, green and blue channels in the RGB image. The top k percentile of tiles with highest blue ratio are selected. We used this method, br-two-stage, as the baseline for ROI detection.

CNN feature extractor. As suggested by the previous study [4], we adopted the Vgg11 model with batch normalization (Vgg11bn) as the backbone for the feature extractor in both 5x and 10x models [39]. The Vgg11bn was initialized with weights pretrained on ImageNet [8]. The feature extractor was first trained on the tile-level dataset for tile classification. After that, the fully connected layers were replaced by a 1×1 convolutional layer to reduce the feature map dimension, outputs of which were flattened and used as instance feature vectors \mathbf{V} in the MIL model for slide classification. The batch size of the tile-level model was set to 50, the initial learning rate was set to $1e^{-5}$. Adam [22] was used for model optimization.

Two-stage classification model. The first stage model was developed for cancer versus non-cancer classification. We transferred the knowledge from the tile-level dataset by initializing the feature extractor with learned weights. The feature extractor was initially fixed, while the attention module and classification layer were trained with a learning rate at $1e^{-4}$ for 10 epochs. Then we fine-tuned the last two convolutional blocks for the Vgg11bn model with a learning rate of $1e^{-5}$ for the feature extractor, and a learning rate of $1e^{-4}$ for the classifier for 90 epochs. Learning rates were reduced by 0.1 if the validation loss did not decrease for the last 10 epochs. The instance dropout rate was set to 0.5. Feature maps of size $512 \times 4 \times 4$ were reduced to $64 \times 4 \times 4$ after the 1×1 convolution, and then flattened to form a 1024×1 vector. A fully connected layer embedded it into a 1024×1 instance feature vector. The size of the hidden layer in the attention module h was set to 512. The model with the highest accuracy on the validation set was utilized to generate attention maps. PCA was used to reduce the dimension of the instance feature vector to 32. K-means clustering was then performed to group similar tiles. The number of clusters was set to 4. Hyper-parameters were tuned on the validation set. Selected tiles at 10x were fed into the second-stage grading model. Similarly, we initialized the feature extractor with weights learned from the tile-level classification. The model was trained for five epochs with the feature extractor fixed. Other hyperparameters were the same as the

Table 1. Model performances on whole slide image classification for prostate cancer

Models	Accuracy (%)	Dataset	Classification task
Zhou <i>et al.</i> [44]	75.00	368 slides	G3 + G4 and G4 + G3 slides
Xu <i>et al.</i> [42]	79.00	312 slides	GS 6, GS 7, and GS 8 slides
Nagpal <i>et al.</i> [30]	70.00	112 million patches and 1490 slides	4 Gleason groups
Ours	85.11	3521 slides	benign, low-grade, high-grade slides

first-stage model. Both tile- and slide-classification models were implemented in PyTorch 0.4, and trained using one NVIDIA Titan X GPU.

5. Results

We have summarized the performance of most state-of-the-art models for prostate WSIs classification in Table 1. The confusion matrix for our best model is shown in Figure 2. As shown in Table 1, the task of Zhou *et al.*'s work [44] is the closest to the presented study, with the main difference being that we included a benign class. The work by Xu *et al.* can be considered relatively easy compared with our task, since differentiating G3 + G4 versus G3 + G4 is non-trivial [30, 44] and often has the largest inter-observer variability. The model developed by Nagpal *et al.* [30] achieved a lower accuracy compared with our model. However, their model predicted more classes, but relied on tile-level labels, which may not be directly comparable.

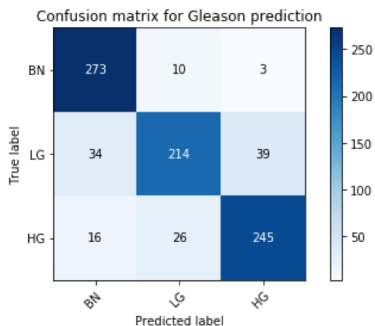


Figure 2. Confusion matrix for Gleason grade classification on the test set

We performed several experiments to evaluate the effects of different components on model performance. Specifically, in experiment *att-two-stage*, we selected informative tiles based only on attention maps generated from the first stage model, while in the *att-cluster-two-stage* model, both instance features and attention maps were used as discussed in section 3.2. Since blue ratio-based tile selection is the most commonly used method, we implemented the *br-two-stage* model to evaluate the effectiveness of the attention-based ROI detection. To investigate the instance dropout,

we trained another model without instance dropout, *att-no-dropout*. To evaluate the contribution of knowledge transferred from the Cedars dataset, we trained a model without transfer learning. For simplicity, we denoted this model as *no-transfer*. The *one-stage* model was trained with tiles only from 5x.

Table 2. Test performances for different multiple instance learning models for whole slide image classification

Models	Accuracy (%)
one-stage	77.80
br-two-stage	80.11
att-two-stage	81.86
att-no-dropout	79.65
no-transfer	84.30
att-cluster-two-stage	85.11

From Table 2, we can see that the model with clustering-based attention achieved the best performance with the average accuracy over 7% higher than the *one-stage* model, over 5% higher than the vanilla attention model (*i.e.* *att-no-dropout*). All two-stage models outperformed the *one-stage*, which utilized all tiles at 5x to predict cancer grading. This is likely due to the fact that important visual features, such as those from nuclei, may only be available at higher resolution. As discussed in section 3.1, attention maps learned in the weakly-supervised model are likely to be only focused on the most discriminative regions instead of the whole part, which could potentially harm model performance.

As shown in Figure 3, clustering with instance features reduced false positive tiles. Pen markers, which may indicate potential suspicious areas, were drawn by pathologists during the diagnosis. We did not use this information for model training, since it was not always available. In Figure 4, we demonstrated the effect of instance dropout. The attention map trained without instance dropout failed to identify the entire region of interest.

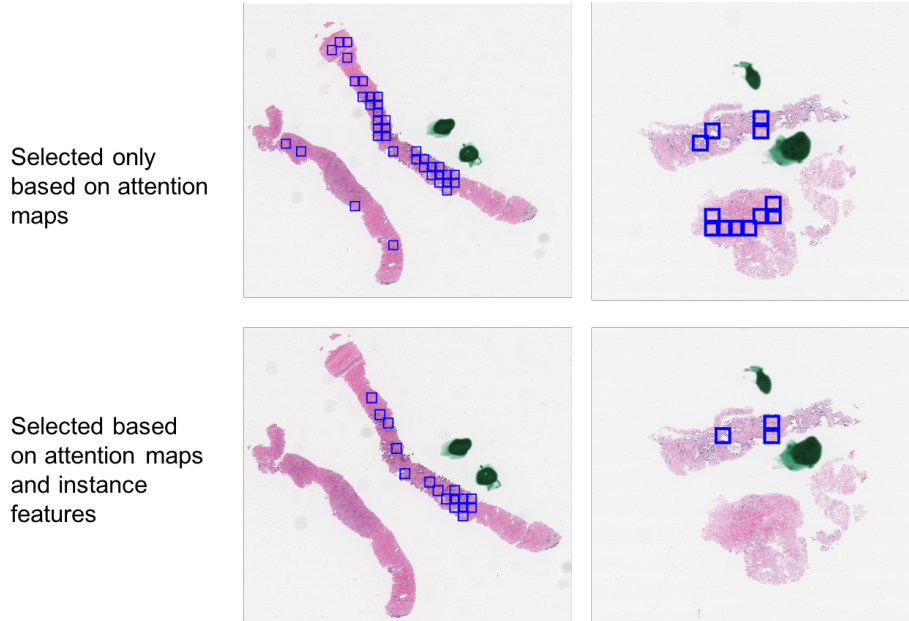


Figure 3. Visualization of selected tiles based on different methods. Each blue box indicates one selected tile.

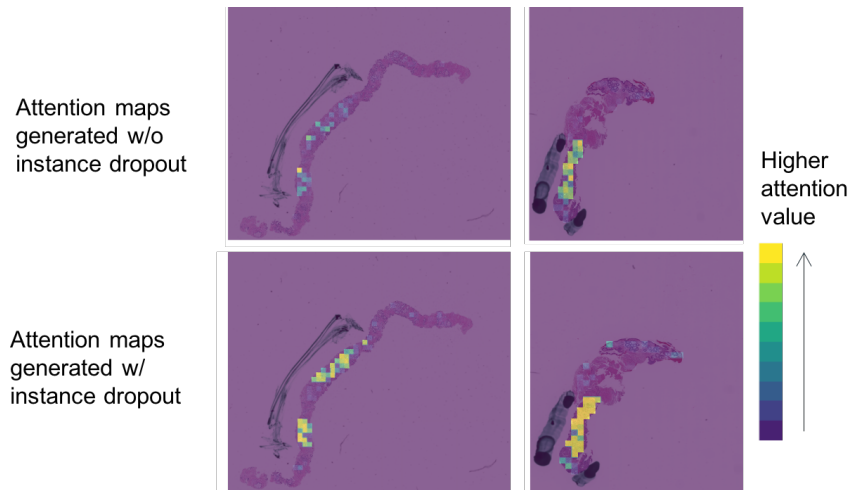


Figure 4. Visualization of the model trained with or without instance dropout.

6. Discussion and Future Work

In this paper, we developed an attention-based two-stage model for WSI classification on a large dataset with thousands of slides from hundreds of patients. Our model was trained to classify low-grade, high-grade, and benign slides. The model achieved an average accuracy of 85.11%, which is over 5% higher compared with the vanilla attention mechanism in [20], and we believe is state-of-the-art performance in prostate biopsy slide classification. In addition, the inherent attention mechanism enhances the interpretability of the classification results.

There are some limitations of this work. Attention maps were only implicitly evaluated using the performance from the second stage model. Annotations or assessment from pathologists are needed for a better evaluation. Moreover, we only included two resolutions (*i.e.* 5x and 10x), which may not be sufficient to capture nucleoli-related features. In future work, higher resolutions will be used. As shown in the results, using transfer learning only slightly improved model performance. The reason could be that we only initialized the feature extractor with learned weights. However, more powerful transfer learning techniques such as [36] will be investigated in the future work.

7. Acknowledgements

The authors would like to acknowledge support from the UCLA Radiology Department Exploratory Research Grant Program (16-0003) and NIH/NCI awards R21CA220352 and P50CA092131. This research was also enabled in part by GPUs donated by NVIDIA Corporation.

References

- [1] Jaume Amores. Multiple instance classification: Review, taxonomy and comparative study. *Artificial intelligence*, 201:81–105, 2013.
- [2] Stuart Andrews, Ioannis Tsochantaridis, and Thomas Hofmann. Support vector machines for multiple-instance learning. In *Advances in neural information processing systems*, pages 577–584, 2003.
- [3] Eirini Arvaniti and Manfred Claassen. Coupling weak and strong supervision for classification of prostate cancer histopathology images. *arXiv preprint arXiv:1811.07013*, 2018.
- [4] Gabriele Campanella, Vitor Werneck Krauss Silva, and Thomas J Fuchs. Terabyte-scale deep multiple instance learning for classification and localization in pathology. *arXiv preprint arXiv:1805.06983*, 2018.
- [5] Yixin Chen, Jinbo Bi, and James Ze Wang. Miles: Multiple-instance learning via embedded instance selection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(12):1931–1947, 2006.
- [6] Marc A DallEra, Peter C Albertsen, Christopher Bangma, Peter R Carroll, H Ballentine Carter, Matthew R Cooperberg, Stephen J Freedland, Laurence H Klotz, Christopher Parker, and Mark S Soloway. Active surveillance for prostate cancer: a systematic review of the literature. *European urology*, 62(6):976–983, 2012.
- [7] Oscar Jiménez del Toro, Manfredo Atzori, Sebastian Otálora, Mats Andersson, Kristian Eurén, Martin Hedlund, Peter Rönquist, and Henning Müller. Convolutional neural networks for an automatic classification of prostate tissue slides with high-grade gleason score. In *Medical Imaging 2017: Digital Pathology*, volume 10140, page 101400O. International Society for Optics and Photonics, 2017.
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [9] Thomas G Dietterich, Richard H Lathrop, and Tomás Lozano-Pérez. Solving the multiple instance problem with axis-parallel rectangles. *Artificial intelligence*, 89(1-2):31–71, 1997.
- [10] Lin Dong. *A comparison of multi-instance learning algorithms*. PhD thesis, The University of Waikato, 2006.
- [11] Scott Doyle, Michael Feldman, John Tomaszewski, and Anant Madabhushi. A boosted bayesian multiresolution classifier for prostate cancer detection from digitized needle biopsies. *IEEE transactions on biomedical engineering*, 59(5):1205–1218, 2012.
- [12] Reza Farjam, Hamid Soltanian-Zadeh, Kourosh Jafari-Khouzani, and Reza A Zoroofi. An image analysis approach for automatic malignancy determination of prostate pathological images. *Cytometry Part B: Clinical Cytometry: The Journal of the International Society for Analytical Cytology*, 72(4):227–240, 2007.
- [13] Arkadiusz Gertych, Nathan Ing, Zhaoxuan Ma, Thomas J Fuchs, Sadri Salman, Sambit Mohanty, Sanica Bhele, Adriana Velásquez-Vacca, Mahul B Amin, and Beatrice S Knudsen. Machine learning approaches to analyze histological images of tissues from radical prostatectomies. *Computerized Medical Imaging and Graphics*, 46:197–208, 2015.
- [14] Lena Gorelick, Olga Veksler, Mena Gaed, José A Gómez, Madeleine Moussa, Glenn Bauman, Aaron Fenster, and Aaron D Ward. Prostate histopathology: Learning tissue component histograms for cancer detection and classification. *IEEE transactions on medical imaging*, 32(10):1804–1818, 2013.
- [15] Le Hou, Dimitris Samaras, Tahsin M Kurc, Yi Gao, James E Davis, and Joel H Saltz. Patch-based convolutional neural network for whole slide tissue image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2424–2433, 2016.
- [16] Qibin Hou, PengTao Jiang, Yunchao Wei, and Ming-Ming Cheng. Self-erasing network for integral object attention. In *Advances in Neural Information Processing Systems*, pages 549–559, 2018.
- [17] Cheng Cheng Huang, Max Xiangtian Kong, Ming Zhou, Andrew B Rosenkrantz, Samir S Taneja, Jonathan Melamed, and Fang-Ming Deng. Gleason score 3+ 4= 7 prostate cancer with minimal quantity of gleason pattern 4 on needle biopsy is associated with low-risk tumor in radical prostatectomy specimen. *The American journal of surgical pathology*, 38(8):1096–1101, 2014.
- [18] Xiao Huang, Jingyuan Zhang, Dingcheng Li, and Ping Li. Knowledge graph embedding based question answering. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, pages 105–113. ACM, 2019.
- [19] Peter A Humphrey. Gleason grading and prognostic factors in carcinoma of the prostate. *Modern pathology*, 17(3):292, 2004.
- [20] Maximilian Ilse, Jakub M Tomczak, and Max Welling. Attention-based deep multiple instance learning. *arXiv preprint arXiv:1802.04712*, 2018.
- [21] Nathan Ing, Zhaoxuan Ma, Jiayun Li, Hootan Salemi, Corey Arnold, Beatrice S Knudsen, and Arkadiusz Gertych. Semantic segmentation for prostate cancer grading by convolutional neural networks. In *Medical Imaging 2018: Digital Pathology*, volume 10581, page 105811B. International Society for Optics and Photonics, 2018.
- [22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [23] Hugh J Lavery and Michael J Droller. Do gleason patterns 3 and 4 prostate cancer represent separate disease states? *The Journal of urology*, 188(5):1667–1675, 2012.

- [24] Peter Lawson, Jordan Schupbach, Brittany Terese Fasy, and John W Sheppard. Persistent homology for the automatic classification of prostate cancer aggressiveness in histopathology images. In *Medical Imaging 2019: Digital Pathology*, volume 10956, page 109560G. International Society for Optics and Photonics, 2019.
- [25] Christy Y Li, Xiaodan Liang, Zhiting Hu, and Eric P Xing. Knowledge-driven encode, retrieve, paraphrase for medical image report generation. *arXiv preprint arXiv:1903.10122*, 2019.
- [26] Jiayun Li, Karthik V Sarma, King Chung Ho, Arkadiusz Gertych, Beatrice S Knudsen, and Corey W Arnold. A multi-scale u-net for semantic segmentation of histological images from radical prostatectomies. In *AMIA Annual Symposium Proceedings*, volume 2017, page 1140. American Medical Informatics Association, 2017.
- [27] Kunpeng Li, Ziyang Wu, Kuan-Chuan Peng, Jan Ernst, and Yun Fu. Tell me where to look: Guided attention inference network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9215–9223, 2018.
- [28] Wenyuan Li, Jiayun Li, Karthik V Sarma, King Chung Ho, Shiwen Shen, Beatrice S Knudsen, Arkadiusz Gertych, and Corey W Arnold. Path r-cnn for prostate cancer diagnosis and gleason grading of histological images. *IEEE transactions on medical imaging*, 2018.
- [29] Oded Maron and Tomás Lozano-Pérez. A framework for multiple-instance learning. In *Advances in neural information processing systems*, pages 570–576, 1998.
- [30] Kunal Nagpal, Davis Foote, Yun Liu, Ellery Wulczyn, Fraser Tan, Niels Olson, Jenny L Smith, Arash Mohtashami, James H Wren, Greg S Corrado, et al. Development and validation of a deep learning algorithm for improving gleason scoring of prostate cancer. *arXiv preprint arXiv:1811.06497*, 2018.
- [31] Kien Nguyen, Anil K Jain, and Bikash Sabata. Prostate cancer detection: Fusion of cytological and textural features. *Journal of pathology informatics*, 2, 2011.
- [32] Emanuele Pesce, Samuel Joseph Withey, Petros-Pavlos Ypsilantis, Robert Bakewell, Vicky Goh, and Giovanni Montana. Learning to detect chest radiographs containing pulmonary lesions using visual attention networks. *Medical image analysis*, 53:26–38, 2019.
- [33] Jan Ramon and Luc De Raedt. Multi instance neural networks. 2000.
- [34] Vikas C Raykar, Balaji Krishnapuram, Jinbo Bi, Murat Dundar, and R Bharat Rao. Bayesian multiple instance learning: automatic feature selection and inductive transfer. In *ICML*, volume 8, pages 808–815, 2008.
- [35] Erik Reinhard, Michael Adhikhmin, Bruce Gooch, and Peter Shirley. Color transfer between images. *IEEE Computer graphics and applications*, 21(5):34–41, 2001.
- [36] Jian Ren, Ilker Hacihaliloglu, Eric A Singer, David J Foran, and Xin Qi. Adversarial domain adaptation for classification of prostate histopathology whole-slide images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 201–209. Springer, 2018.
- [37] Jo Schlemper, Ozan Oktay, Michiel Schaap, Mattias Heinrich, Bernhard Kainz, Ben Glocker, and Daniel Rueckert. Attention gated networks: Learning to leverage salient regions in medical images. *Medical image analysis*, 53:197–207, 2019.
- [38] Rebecca L Siegel, Kimberly D Miller, and Ahmedin Jemal. Cancer statistics, 2019. *CA: a cancer journal for clinicians*, 69(1):7–34, 2019.
- [39] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [40] Krishna Kumar Singh and Yong Jae Lee. Hide-and-seek: Forcing a network to be meticulous for weakly-supervised object and action localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 3544–3553. IEEE, 2017.
- [41] Jeffrey J Tosoian, H Ballentine Carter, Abbey Lepor, and Stacy Loeb. Active surveillance for prostate cancer: current evidence and contemporary state of practice. *Nature Reviews Urology*, 13(4):205, 2016.
- [42] Hongming Xu, Sunhoh Park, and Tae Hyun Hwang. Automatic classification of prostate cancer gleason scores from digitized whole slide tissue biopsies. *bioRxiv*, page 315648, 2018.
- [43] Xiaolin Zhang, Yunchao Wei, Jiashi Feng, Yi Yang, and Thomas S Huang. Adversarial complementary learning for weakly supervised object localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1325–1334, 2018.
- [44] Naiyun Zhou, Andrey Fedorov, Fiona Fennessy, Ron Kikinis, and Yi Gao. Large scale digital prostate pathology image analysis combining feature extraction and deep neural network. *arXiv preprint arXiv:1705.02678*, 2017.