# VISUAL ATTENTION ANALYSIS OF PATHOLOGISTS EXAMINING WHOLE SLIDE IMAGES OF PROSTATE CANCER

*Souradeep Chakraborty[1], Ke Ma[1], Rajarsi Gupta[2], Beatrice Knudsen[3], Gregory J. Zelinsky[1,4],*
*Joel H. Saltz[2], Dimitris Samaras[1]*

[1] Department of Computer Science, Stony Brook University, NY, USA
[2] Department of Biomedical Informatics, Stony Brook University, NY, USA
[3] Department of Pathology, University of Utah School of Medicine, Utah, USA
[4] Department of Psychology, Stony Brook University, NY, USA

## ABSTRACT

We study the attention of pathologists as they examine whole-slide images (WSIs) of prostate cancer tissue using a digital microscope. To the best of our knowledge, our study is the first to report in detail how pathologists navigate WSIs of prostate cancer as they accumulate information for their diagnoses. We collected slide navigation data (i.e., viewport location, magnification level, and time) from 13 pathologists in 2 groups (5 genitourinary (GU) specialists and 8 general pathologists) and generated visual attention heatmaps and scanpaths. Each pathologist examined five WSIs from the TCGA PRAD dataset, which were selected by a GU pathology specialist. We examined and analyzed the distributions of visual attention for each group of pathologists after each WSI was examined. To quantify the relationship between a pathologist's attention and evidence for cancer in the WSI, we obtained tumor annotations from a genitourinary specialist. We used these annotations to compute the overlap between the distribution of visual attention and annotated tumor region to identify strong correlations. Motivated by this analysis, we trained a deep learning model to predict visual attention on unseen WSIs. We find that the attention heatmaps predicted by our model correlate quite well with the ground truth attention heatmap and tumor annotations on a test set of 17 WSIs by using various spatial and temporal evaluation metrics.

*Index Terms*— Prostate cancer, visual attention, tumor segmentation, digital histopathology

## 1. INTRODUCTION

Research on attention tracking in digital histopathology images has been evolving [1, 2, 3] in the field of medical imaging. Being able to analyze and predict the visual attention behavior of a pathologist during the examination of WSIs is useful in developing computer-assisted training and clinical decision support systems. For example, pathologists in training can benefit from visualizing and comparing their attention behavior to experienced pathologists with specialty expertise with the goal of improving interobserver variability in tasks such as Gleason grading in prostate cancer.

Earliest works to interpret attention behavior of pathologists as they view WSIs of cancer tissue samples, include [4] that conducted eye tracking studies to determine the effect of tumor architecture on the grading of prostate cancer images and [5] that examined the spatial coupling between eye-gaze and mouse cursor positions during WSI viewing, and showed that mouse movement tracking data could be a reliable indicator of a physician's attention and diagnostic behavior. [3] explored the complex interactions between pathologists and WSIs that guide diagnostic decision-making using eye-tracking studies. The study in [6] used a web-based viewer to record pathologist behavior in order to characterize diagnostic search patterns (scanning and drilling) using viewport attention data from a digital microscope. Similarly, we also used the viewport location, time stamp, and zoom level data in our study. More recently, the works of [1] and [2] have highlighted the importance of eye tracking and revealed expertise-related differences in medical image inspection and diagnosis. However, significantly less focus has been allocated on visualizing the spatiotemporal distribution of visual attention in relation with the tumor regions. Also, previous studies have not presented a model that can predict visual attention of pathologists during whole slide image viewing.

In this paper, we provide a detailed analysis of how WSIs of prostate cancer were examined by pathologists, which is further correlated with tumor annotations to analyze the differences in the viewing patterns of GU specialists and general pathologists. To the best of our knowledge, we are the first to study the visual attention of pathologists and predict the distribution of attention during the histopathologic evaluation of prostate cancer in WSIs. We measured viewing behavior by collecting slide navigation data while pathologists examined prostate cancer WSIs as a reasonable surrogate to capture visual attention. We constructed visual attention heatmaps from capturing viewports during examination that we then compared to tumor regions annotated by a GU expert in order to ascertain where the pathologists focused their attention.
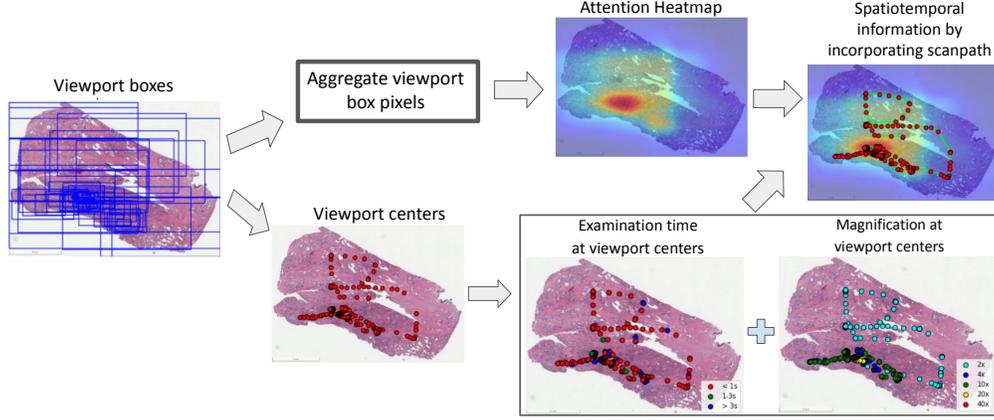
**Fig. 1**. Processing of the captured attention data to obtain visual attention heatmaps and scanpaths.

We also analyzed how visual attention varies across different magnification levels. We used various evaluation metrics to compare the spatial and spatiotemporal attention distribution to the annotated tumor regions. Our study is proof of concept for collecting slide navigation data to generate surrogate visual trajectories and associated attention heatmaps to understand how pathologists identify cancer presence in WSIs. Future studies will explore intra- and inter-observer variability during the evaluation of more complex cancer features.

## 2. METHODS

### 2.1. Data collection

We used QuIP caMicroscope, a web-based toolset for digital pathology data management and visualization [7] to record the attention data of 13 pathologists as they viewed Prostate cancer whole slide images. We asked each participant to provide their expertise level as: (1) boarded and practising anatomic pathologist, (2) resident or fellow. Our data came from 5 Genitourinary (GU) specialists and 8 general pathologists. Each subject was instructed to view five whole side images on our web based viewing interface. A GU specialist selected the five WSIs we used for our study among 342 WSIs of the TCGA-PRAD dataset [8]. Average viewing time per slide per pathologist was 103 seconds. No compensation was provided to the pathologists for participation.

Also, in order to compare our attention data with tumor evidence, we collected tumor region annotations from a GU specialist on the five WSIs in this study. Please see these tumor annotations in Sec. 3 (second column in Fig. 3(a)).

### 2.2. Data processing

Our attention data provides us information about where and how long pathologists looked at, and the gaze shifts they made from one region to another while viewing the WSIs. Fig. 1, shows our pipeline to process the collected attention data (using viewport boxes). Our attention data comes in two forms: (1) attention heatmaps, (2) attention scanpaths.

**Attention Heatmap**: The attention heatmap shows the aggregate spatial distribution of pathologist attention during WSI viewing. We assign a value of 1 at every image pixel within a viewport box and sum up the values over all viewport boxes to construct our attention heatmap. Next, we normalize this heatmap to obtain the final attention heatmap, which is:

$$HM'^I_{Attn.}(x,y) = G^\sigma * \sum_{v=1}^{V} (\sum_{v_x^s}^{v_x^e} \sum_{v_y^s}^{v_y^e} 1)$$

$$HM^I_{Attn.} = \frac{HM'^I_{Attn.} - min(HM'^I_{Attn.})}{max(HM'^I_{Attn.}) - min(HM'^I_{Attn.})}$$

(1)

where, $HM'^I_{Attn.}$ is the intermediate attention heatmap, $HM^I_{Attn.}$ is the final normalized attention heatmap, $V$ is the number of viewport boxes on a WSI $I$, and $v_x^s$, $v_x^e$, $v_y^s$, $v_y^e$ denote the starting and the ending $x$ and $y$ coordinates of the viewport box $v$ respectively. $G^\sigma$ is a 2D gaussian with $\sigma = 16$ pixels that smooths the intermediate heatmap $HM'^I_{Attn.}$.

**Attention Scanpath**: The attention heatmap only captures aggregate spatial attention distribution. In order to capture the viewing trajectories, we produce attention scanpaths from the collected viewport data. We stack the viewport centers of every viewport box, $v$ in WSI $I$ in order to construct our scanpath, $SP^I_{Attention}$ as $SP^I_{Attention} = [v^1_{center}, v^2_{center}, ..., v^{(V-1)}_{center}, v^V_{center}]$, where the viewport center of a viewport box $i$ is $v^i_{center} = (\frac{v_x^s+v_x^e}{2}, \frac{v_y^s+v_y^e}{2})$.

In Fig. 1, we also show the viewing examination time and the magnification levels at the viewport centers. In the depicted WSI instance, the pathologist spent less than 1 second at most of the viewports (indicated by the red viewport centers) and mostly viewed the slide at 2X and 4X magnification.

### 2.3. Predicting attention heatmaps using deep learning

Here we discuss the deep learning model we trained for predicting attention heatmap over a WSI. We formulate attention
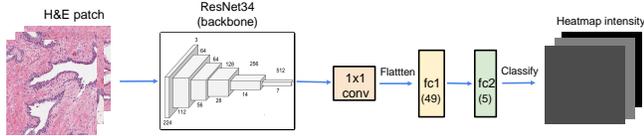
**Fig. 2**. Our model, ProstAttNet for predicting visual attention on whole slide images of Prostate cancer.

prediction as a classification task where the goal is to classify a WSI patch into one of the $N$ attention intensity bins ($N = 5$ in our study). During training, we discretize the average pixel intensity of every heatmap patch into an intensity bin. At inference time, we assign the average pixel intensity of a bin to the image patch according to the predicted class in order to construct the patch-wise heatmaps. We reconstruct the final attention heatmap over the WSI by assembling the predicted patch-wise heatmaps followed by gaussian smoothing and map normalization.
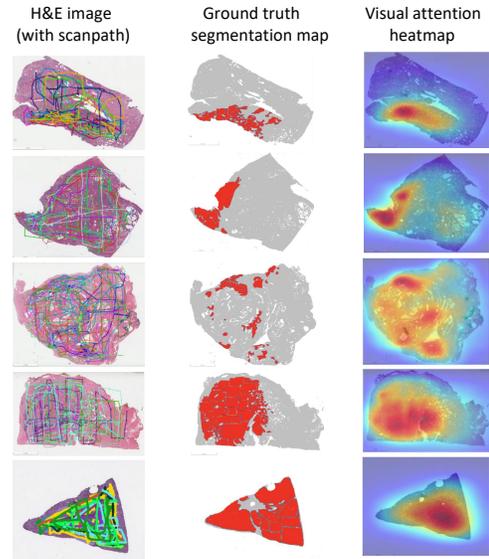
We trained a CNN model, ProstAttNet (Prostate AttentionNet), to classify an image patch into one of the 5 attention intensity levels using the pre-trained ResNet34 model [9] as the backbone followed by a $1 \times 1$ convolutional layer and two fully connected (fc) layers. We depict our attention prediction model, ProstAttNet in Fig. 2. See suppl. for training details.
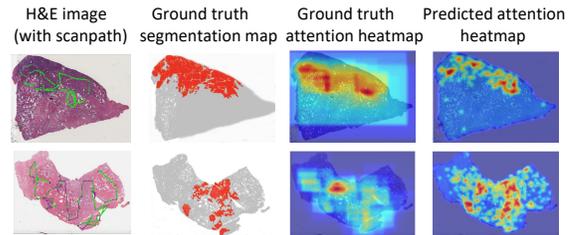
## 3. RESULTS

### 3.1. Qualitative Evaluation

Fig. 3(a) shows ground truth tumor segmentation maps and attention heatmaps for the five H&E WSIs of this study. Attention heatmaps have a high spatial correlation with tumor locations in the ground truth tumor segmentation map. We also predicted attention heatmaps using our ProstAttNet model on a test dataset of 17 whole slide images (from the TCGA-PRAD dataset). Fig. 3(b) compares prediction results on two WSI instances from our test dataset. The ground truth attention heatmap was constructed from slide navigation data collected from th GU specialist only for the purpose of validating model predictions. We see that the predicted attention heatmap correlates well with the ground truth attention heatmap and the ground truth tumor segmentations.

We also compared the viewing behavior of the GU specialists and the general pathologists in our study. In Fig. 4, we show how pathologists examined WSIs of Prostate cancer by overlaying navigation scanpaths on H&E images and generating attention heatmaps to compare overall behavior of all pathologists, which is subdivided to evaluate potential differences in behavior between genitourinary (GU) specialists and general pathologists for one of the cases in our study. We also generate attention heatmaps with respect to the viewport magnification level to evaluate visual behavior at 4X, 10X, 20X, and 40X. We see a high concurrence in the viewed regions at every magnification level among both groups of pathologists. In addition, we also show the average viewing time of the



(a)



(b)

**Fig. 3**. (a) Comparison of the visual attention heatmaps with the ground truth tumor segmentation maps. We see a strong correlation between the two maps. (b) Comparison of the predicted attention heatmap using our ProstAttNet model with the ground truth segmentation map and the ground truth attention heatmap constructed from attention data collected from a Genitourinary specialist on two test WSI instances.

two groups of pathologists across four different magnification levels. While we see good concurrence within the groups in terms of the attended regions, we also observe some differences in their averaged attention heatmaps. While the general pathologists did not seem to have allocated sufficient attention to the bottom right tumor regions, the GU specialists looked at more of the tumor region (left and right high intensity areas) as seen in the overall averaged heatmap in the second row.

### 3.2. Quantitative Evaluation

To evaluate how well attention data predicts tumor regions (from pathologists' annotations), we used the Cross correlation (CC) metric to compare our attention heatmaps and the Semantic Sequence Score (SSS) metric to evaluate the scanpaths constructed from the viewport centers. We compared the attention heatmaps with tumor proability maps constructed from the tumor annotations. All compared maps were downsampled by a factor of 1/16 (compared to origi-
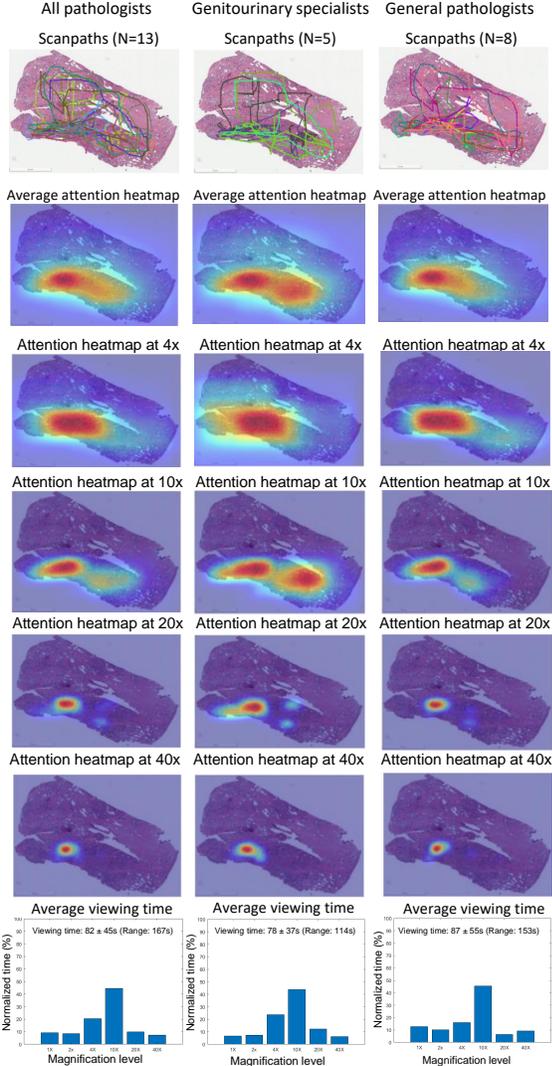
**Fig. 4**. Comparison of attention heatmaps (average and magnification level-wise), scanpaths and average viewing time of the Genitourinary specialists and the general pathologists on the TCGA-EJ-7328 WSI from the TCGA-PRAD dataset.

| Case | All pathologists | | GU specialists | | General pathologists | |
|------|------|------|------|------|------|------|
| | CC | SSS | CC | SSS | CC | SSS |
| TCGA-EJ-7328 | 0.729 | 0.421 | 0.765 | 0.390 | 0.706 | 0.441 |
| TCGA-HC-A8D1 | 0.877 | 0.412 | 0.881 | 0.464 | 0.874 | 0.380 |
| TCGA-G9-6384 | 0.712 | 0.481 | 0.725 | 0.562 | 0.704 | 0.430 |
| TCGA-EJ-7315 | 0.780 | 0.390 | 0.787 | 0.445 | 0.776 | 0.355 |
| TCGA-G9-6494 | 0.473 | 0.398 | 0.437 | 0.496 | 0.495 | 0.336 |
| Average | 0.714 | 0.420 | 0.719 | 0.472 | 0.711 | 0.388 |

**Table 1**. Quantitative evaluation of attention data and tumor annotations (by pathologist), CC is the Cross-Correlation score for comparing attention heatmap with ground truth segmentation map and SSS is Semantic Sequence Score comparing inter-observer similarity on attention scanpaths.

during whole slide image viewing. To obtain SSS we modified the Sequence Score (SS) metric [12] (generally used to compare scanpaths on natural images) by substituting clusters based on eye fixations with clusters based on Gleason grades. Specifically, we first convert a scanpath into a string representing the sequence of Gleason grades corresponding to the viewport box centers (e.g. $G_3$-$G_5$-$G_4$-$G_3$, $G_4$-$G_4$-$G_3$-$G_5$-$G_5$, etc. where $G_n$ represents Gleason grade $n$). We then use a string matching algorithm [13] to measure string similarity.

Table. 1 lists the different evaluation metrics discussed above for the GU specialists and the general pathologists on the five WSIs. A high average cross correlation score of 0.714 across all pathologists suggests a strong spatial correlation between the distribution of attention and the ground truth tumor locations. Despite slight differences in the viewing behaviors of the two groups, independent t-test analysis indicates that differences in the degree of correlation between tumor region and attention heatmaps are not significant ($p > 0.52$ averaged over the 5 images). However, the consistently high SSS for GU specialists compared to general pathologists indicates more consistent viewing behavior within this group compared to general pathologists. Using ProstAttNet, we obtained average CC score of 0.453 between predicted attention and ground truth attention heatmaps on our test set of 17 WSIs. The predicted attention heatmap overlapped with the ground truth segmentation map with a CC score of 0.532.

## 4. CONCLUSION

We have shown how pathologists allocate attention while viewing prostate cancer WSIs and presented a deep learning model that predicts visual attention. Our data visualization schema allow us to understand how pathologists examine WSIs to identify cancer. We find a strong correlation between the locations of tumor, annotated by a GU specialist, and attention heatmaps of 13 pathologists. In the future, we will collect attention data in a larger study with more WSIs in order to improve our attention prediction model. Also, we plan to leverage attention data to train a deep learning model for tumor segmentation and Gleason grading to test whether we can circumvent the need for extensive annotation.

nal image size) for computational reasons. We convolved the tumor segmentation map with a 2D gaussian ($\sigma = 16$ pixels) and normalized this map to obtain the tumor probability map for comparison with the attention heatmap. In order to ensure that the distributions of the attention heatmap and the tumor probability map are similar, we perform histogram matching [10] of the two maps as a pre-processing step [11].

We use the cross-correlation score (CC) by interpreting the attention heatmap, $HM_{Atten.}$ and the tumor probability map, $PM_{Tumor}$ as random variables and measure the linear relationship between them. High positive CC values occur where the compared maps have values of similar magnitude. The Semantic Sequence Score (SSS) metric we use captures the inter-observer scanpath similarity across the two groups of pathologists in terms of the grade of tumor regions traversed

## 5. COMPLIANCE WITH ETHICAL STANDARDS

Data collection was conducted within the ethical guidelines established and overseen by the Stony Brook University Institutional Review Board.

## 6. ACKNOWLEDGEMENTS

## 7. SUPPLEMENTARY: TRAINING DETAILS OF PROSTATTNET

During training of our ProstAttNet model, we froze the pre-trained ResNet34 weights and allowed gradient flow over the fully connected layers of our model. For training the model, we used image patches (each of size $500 \times 500$) extracted from four whole slide images at 10x magnification (the most frequent magnification level used by pathologists during WSI viewing per our analysis). The remaining single slide was used for validation. We generated a total of 58.5K image patches for training with labels corresponding to 13 pathologists per patch. We also performed data augmentation [14] by introducing color jitter, random horizontal and vertical image flips during training. We trained this model using the Cross-Entropy loss between the predicted and the ground truth heatmap intensity bin. We used the Adam optimizer [15] with initial learning rate of 0.005. Training converged within 20 epochs with a total training time of 8.5 hours on a Nvidia Titan-Xp GPU. We used the PyTorch deep learning library for the model implementation.

## 8. REFERENCES

[1] Tad T Brunyé, Trafton Drew, Kathleen F Kerr, Hannah Shucard, Donald L Weaver, and Joann G Elmore, "Eye tracking reveals expertise-related differences in the time-course of medical image inspection and diagnosis," *Journal of Medical Imaging*, vol. 7, no. 5, pp. 051203, 2020.

[2] Ellhia Sudin, Dorina Roy, Nourdin Kadi, Panagiotis Triantafyllakis, Guprit Atwal, Alastair Gale, Ian Ellis, David Snead, and Yan Chen, "Eye tracking in digital pathology: identifying expert and novice patterns in visual search behaviour," in *Medical Imaging 2021: Digital Pathology*. International Society for Optics and Photonics, 2021, vol. 11603, p. 116030Z.

[3] Tad T Brunyé, Ezgi Mercan, Donald L Weaver, and Joann G Elmore, "Accuracy is in the eyes of the pathologist: the visual interpretive process and diagnostic accuracy with digital whole slide images," *Journal of biomedical informatics*, vol. 66, pp. 171–179, 2017.

[4] Dario Bombari, Braulio Mora, Stephan C Schaefer, Fred W Mast, and Hans-Anton Lehr, "What was i thinking? eye-tracking experiments underscore the bias that architecture exerts on nuclear grading in prostate cancer," *PLoS One*, vol. 7, no. 5, pp. e38023, 2012.

[5] Vignesh Raghunath, Melissa O Braxton, Stephanie A Gagnon, Tad T Brunyé, Kimberly H Allison, Lisa M Reisch, Donald L Weaver, Joann G Elmore, and Linda G Shapiro, "Mouse cursor movement and eye tracking data as an indicator of pathologists' attention when viewing digital whole slide images," *Journal of pathology informatics*, vol. 3, 2012.

[6] Ezgi Mercan, Linda G Shapiro, Tad T Brunyé, Donald L Weaver, and Joann G Elmore, "Characterizing diagnostic search patterns in digital breast pathology: Scanners and drillers," *Journal of digital imaging*, vol. 31, no. 1, pp. 32–41, 2018.

[7] Joel Saltz, Ashish Sharma, Ganesh Iyer, Erich Bremer, Feiqiao Wang, Alina Jasniewski, Tammy DiPrima, Jonas S Almeida, Yi Gao, Tianhao Zhao, et al., "A containerized software system for generation, management, and exploration of features from whole slide tissue images," *Cancer research*, vol. 77, no. 21, pp. e79–e82, 2017.

[8] M Zuley, R Jarosz, B Drake, D Rancilio, A Klim, K Rieger-Christ, and J Lemmerman, "Radiology data from the cancer genome atlas prostate adenocarcinoma [tcga-prad] collection," *Cancer Imaging Arch*, vol. 9, 2016.

[9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[10] Rafael C Gonzales and BA Fittes, "Gray-level transformations for interactive image enhancement," *Mechanism and Machine Theory*, vol. 12, no. 1, pp. 111–122, 1977.

[11] Candace E Peacock, Taylor R Hayes, and John M Henderson, "Center bias does not account for the advantage of meaning over salience in attentional guidance during scene viewing," *Frontiers in Psychology*, vol. 11, 2020.

[12] Ali Borji, Hamed R Tavakoli, Dicky N Sihite, and Laurent Itti, "Analysis of scores, datasets, and models in

visual saliency prediction," in *Proceedings of the IEEE international conference on computer vision*, 2013, pp. 921–928.

[13] Saul B Needleman and Christian D Wunsch, "A general method applicable to the search for similarities in the amino acid sequence of two proteins," *Journal of molecular biology*, vol. 48, no. 3, pp. 443–453, 1970.

[14] David Tellez, Maschenka Balkenhol, Irene Otte-Höller, Rob van de Loo, Rob Vogels, Peter Bult, Carla Wauters, Willem Vreuls, Suzanne Mol, Nico Karssemeijer, et al., "Whole-slide mitosis detection in h&e breast histology using phh3 as a reference to train distilled stain-invariant convolutional networks," *IEEE transactions on medical imaging*, vol. 37, no. 9, pp. 2126–2136, 2018.

[15] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.