

Multi-Scale Attention-based Multiple Instance Learning for Classification of Multi-Gigapixel Histology Images

Made Satria Wibawa¹, Kwok-Wai Lo², Lawrence S. Young³, and Nasir Rajpoot^{1,4}

¹ Tissue Image Analytics Centre, Department of Computer Science, University of Warwick

{made-satria.wibawa,n.m.rajpoot}@warwick.ac.uk

² Department of Anatomical and Cellular Pathology, The Chinese University of Hong Kong

kwlo@cuhk.edu.hk

³ Warwick Medical School, University of Warwick

l.s.young@warwick.ac.uk

⁴ The Alan Turing Institute, London

Abstract. Histology images with multi-gigapixel of resolution yield rich information for cancer diagnosis and prognosis. Most of the time, only slide-level label is available because pixel-wise annotation is labour intensive task. In this paper, we propose a deep learning pipeline for classification in histology images. Using multiple instance learning, we attempt to predict the latent membrane protein 1 (LMP1) status of nasopharyngeal carcinoma (NPC) based on haematoxylin and eosin-stain (H&E) histology images. We utilised attention mechanism with residual connection for our aggregation layers. In our 3-fold cross-validation experiment, we achieved average accuracy, AUC and F1-score 0.936, 0.995 and 0.862, respectively. This method also allows us to examine the model interpretability by visualising attention scores. To the best of our knowledge, this is the first attempt to predict LMP1 status on NPC using deep learning.

Keywords: deep learning, multiple instance learning, attention, H&E, LMP1, NPC

1 Introduction

Analysing histology images for cancer diagnosis and the prognosis is not a task without challenges. A Whole Slide Image (WSI) in $40\times$ magnification with an average resolution of $200,000 \times 150,000$ contains 30 billion pixels and a size of ~ 90 GB in the uncompressed state. Moreover, histology image mainly has a noisy/ambiguous label. Most of the time label of the WSI is assigned on the slide level, for example, the label for cancer stage. This may create ambiguous interpretations by the machine learning/deep learning model, because in such

a large region, not all regions correspond to the slide label. Some regions may contain tumours, and the rest is non-tumour regions e.g. background, stroma, and lymphocyte.

The standard approach to handling multi-gigapixel WSIs is by extracting its region into smaller patches that the machine can process. The noisy label problem also occurs in this approach since we do not have a label for each patch. The multiple instance learning (MIL) paradigm is generally used to overcome this problem [3]. In MIL, each patch is represented as an instance in a bag. Since WSIs have more than one patch, the bag contains multiple instances, hence the name 'multiple' instances learning. During training, only global (slide-level) image labels are required for supervision. Aggregation mechanism is then utilised to summarise all information in instances to make a final prediction.

The selection of aggregation mechanism is one of the vital parameter in MIL. Several aggregation mechanisms have been introduced in MIL for computational histopathology, such as the max and the median score [10] or the average of all patches scores [4,16]. The limitation of such mechanisms is they are not trainable. Therefore, attention mechanism [1] was utilised in the recent MIL model on computational pathology [15,7]. Attention mechanism is trainable, it computes the weights of the instances representation. Strong weights/scores indicate instances are more important to final prediction than instances with low scores.

Several studies have already utilised MIL for cancer diagnosis and prognosis. For example, Lu and colleagues [13] proposed attention-based multiple instance learning for subtyping renal cell carcinoma and non-small-cell lung carcinoma. Campanella and colleagues [2] used multiple instance learning with RNN-based aggregation to discriminate tumour regions in breast, skin, and prostate cancer. No current study examines the use of multiple instance learning in nasopharyngeal carcinoma.

Nasopharyngeal carcinoma (NPC) is a malignancy that develops from epithelial cells within lymphocyte-rich nasopharyngeal mucosa. The high incidence rate of NPC mainly occurs in southern China, southeast Asia and north Africa, which are 50-100 times greater than the rates in other regions of the world. NPC has unique pathogenesis involving genetic, lifestyle and viral (Epstein-Barr virus or EBV) cofactors [19]. EBV infection on NPC encodes several viral oncoprotein, one of the most important ones is latent membrane protein (LMP1). LMP1 is the predominant oncogenic driver of NPC [9]. Overexpression of LMP1 promotes NPC progression by inducing invasive growth of human epithelial and nasopharyngeal cells. Due its nature, LMP1 is an excellent therapeutic target for EBV-associated NPC [11]. However, detection of LMP1 on NPC with standard testing such as immunohistochemistry may cause additional costs and delays the diagnosis. On the other hand, histopathology images are widely available and provide rich information of the cancer on the tissue level. Deep learning-based algorithms have been applied for several tasks and deliver promising results in EBV-related cancer in histopathology. Deep learning has been used to predict EBV status in gastric cancer [21], predict EBV integration sites [12] and predict

microsatellite instability status in gastric cancer [14]. However, despite the advance in computational histopathology and the significant importance of LMP1 in NPC progression, this topic is still an understudy area.

In this paper, we proposed a deep learning pipeline based on the MIL paradigm to predict LMP1 status in NPC. Our main contributions in this paper are: (1) To the best of our knowledge, this is the first study that attempts to utilise deep learning on LMP1 status prediction based on histology images of NPC. (2) We propose multi-scale attention-based with residual connection¹ for aggregation layer in MIL. (3) Our proposed MIL pipeline outperformed other known MIL methods for predicting LMP1 status in NPC.

2 Materials and Methods

2.1 Dataset

This study used two kinds of datasets: the first for the tissue classification task and the second is the LMP1 dataset.

- **Tissue classification dataset.** For this task, we combined two datasets from Kather100K [8] dataset and an internal dataset. Kather100K contains 100,000 non-overlapping image patches from H&E stained histological images of colorectal cancer. The number of tissue classes in this dataset is nine, we exclude the normal colon mucosa class (NORM) from the dataset. We argued that the NORM class is irrelevant in head and neck tissue. Thus, only eight classes from the Kather100K dataset were utilised in this study. The internal dataset consists of three cohorts. The number of images patches in this dataset is 180,344 with seven classes which taken at 20x magnification. We combined images with the same classes from Kather100K and internal datasets in the final dataset.
- **LMP1 dataset.** LMP1 data was collected from an University Hospital. The dataset comprises 101 NPC cases with a H&E stained whole slide images (WSIs) for each case. All WSIs were taken by an Aperio scanner at 40x magnification and 0.250 microns per pixel resolution. Expression of LMP1 was determined by immunohistochemical staining. The proportion score was according to the percentage of tumour cells with positive membrane and cytoplasmic staining (0–100). The intensity score was assigned for the average intensity of positive tumour cells (0,none; 1,weak; 2,intermediate; 3,strong). The LMP1 staining score was the product of proportion and intensity scores, ranging from 0 to 300. The LMP1 expression was categorized into absence/low/negative (score 0–100) and high/positive (score 101–300). Details of the class distribution of the dataset in this study can be seen in Table 1.

¹ Model code: <https://github.com/mdsatria/MultiAttentionMIL>

Table 1. Details of the dataset

Label	Number of Cases
LMP1 positive	25
LMP1 negative	76
Total cases	101

2.2 Tissue Classification

We utilised a pretrained ResNet-50 [6] with ImageNet [5] to classify images into tissue and non-tissue classes. Dataset for tissue classification was split into three parts, i.e. training, validation and testing with the data distribution of 80%, 10% and 10%, respectively. Training images were augmented in several methods, including colour augmentation, kernel filter (sharpening and blurring) and geometric transformation. Before images used in training, we normalised the value of the pixel within the range of -1 and 1.

We trained all layers in the ResNet50 model, using an Adam optimizer with starting learning rate of 1×10^{-3} and weight decay of 5×10^{-4} . Training was conducted on 20 epochs and every 10 epochs, we decreased learning rate by ten times. To regularize the model, we apply a dropout layer with the probability of 0.5 before the fully connected layer. We use cross entropy to calculate the loss of the network which is defined as:

$$CE = \sum_{i=1}^C t_i \log(f(s)_i) \quad (1)$$

where

$$f(s)_i = \frac{e_i^s}{\sum_{j=1}^C e_j^s} \quad (2)$$

t_i is the ground truth label and s_i is the predicted label and $f(s)_i$ is a softmax activation function. In our case, there were eight classes, thus $C \in [0..7]$.

2.3 LMP1 Prediction

The main pipeline of LMP1 prediction is illustrated in Figure 1. After tumour patches were identified in the tissue classification stage, we encoded all the tumour patches into feature vectors. We experimented with two backbone networks for encoding patches into feature vectors, namely ResNet-18 and EfficientNet-B0 [17].

Let B denote the bag, x is the feature vector extracted from backbone network and K is the number of patches in the slide, bag of instances will be $B = \{x_1, \dots, x_K\}$. Within the ResNet-18 and EfficientNet-B0 as the backbone networks, each patch is represented by a 512 and 1280-dimensional features vector, respectively. The number of K will vary depending from the number of

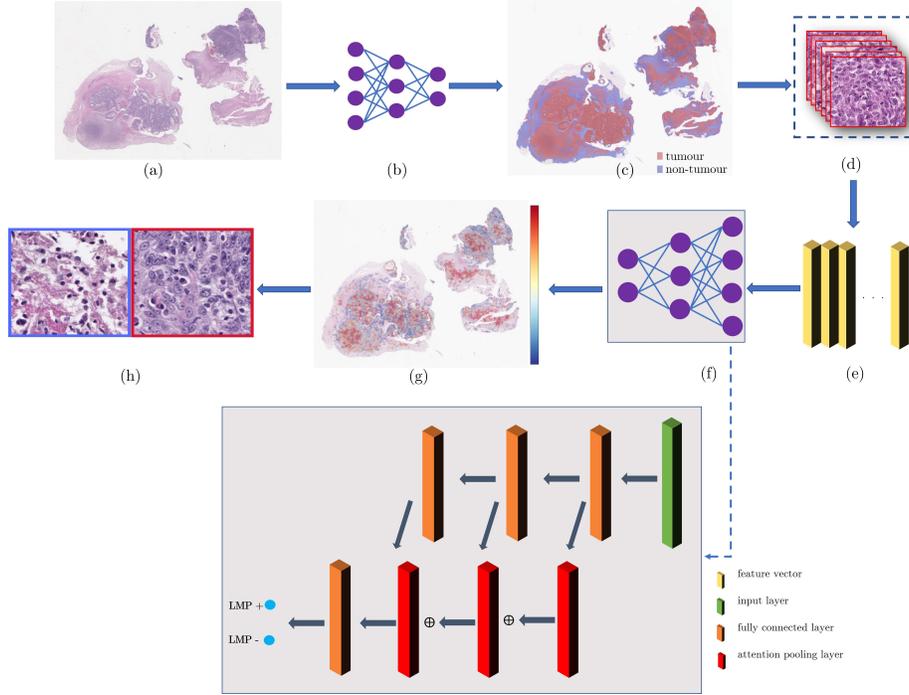


Fig. 1. The workflow of LMP1 Prediction. Patches from WSIs (a) were classified into the tumour and non-tumour regions by the tumour classification model based on ResNet50 (b). The tumour regions of WSIs can be seen as the red area of the heatmap (c). All tumour patches were then used as instances of the bag (d) and encoded into feature vectors (e) with backbone networks. Our MIL model which consists of four fully connected layers, and three attention-based aggregation layers with the residual connection (f) then trained with these feature vectors. The model will generate attention scores for each patch. The red area in the heatmap shows a high attention score (g&h).

tumour patches in their corresponding slide. Slide labels will be inherited into their corresponding patches. Thus, for all bags of instances $\{B_1, B_2, \dots, B_N\}$ with N as the number of slide, every x in B will have same label as their corresponding slide $Y_n \in \{0, 1\}, n = 1 \dots N$.

We used attention mechanism as an aggregation layer as widely used in [7,15,20]. An attention layer in aggregation layer is defined as a weighted sum:

$$z = \sum_{k=1}^K a_k x_k, \quad (3)$$

where

$$a_k = \frac{\exp\{w^\top \tanh(Vx_k^\top)\}}{\sum_{j=1}^K \exp\{w^\top \tanh(Vx_j^\top)\}} \quad (4)$$

and $w \in \mathbb{R}^{L \times 1}$ is weight vector of each instance, $V \in \mathbb{R}^{L \times M}$. \tanh is hyperbolic tangent function which is used in the first hidden layer. The second layer employs a softmax non-linearity function to ensure that the attention weights sum is equal to one. The attention weights can be interpreted as the relative importance of the instance. The higher the attention weight, the more important an instance in the final prediction score. Furthermore, this mechanism also creates a fully trainable aggregation layer.

We applied the ReLU non-linearity function for each fully connected layer except for the last of fully connected layers. Each output from the fully connected layer will go through the next fully connected layer and the aggregation layers. This model has three weighted sums z_1, z_2, z_3 that are then accumulated in the last aggregation via the residual connection. The accumulation of weighted sum from each aggregation layer is then used in the last fully connected layers to predict label of the bag.

We trained our model with learning rate of 1×10^{-4} and with Adam optimization algorithm. All the experiments were implemented on Python language and Pytorch framework. We conducted our experiment on workstation with Intel(R) i5-10500 CPU (3.1GHz), 64 GB RAM and single GPU NVIDIA RTX 3090.

2.4 Metrics for Evaluation

We use the accuracy score to measure our model performance. Due to the imbalanced class distribution in our dataset, we also use F1-score and Area Under the Curve and Receiver Operating Characteristic (AUROC) curve. F1-score is defined as follow:

$$F1 = \frac{TP}{TP + \frac{1}{2}(FP + FN)} \quad (5)$$

where TP, TN, FN, and FP are true positive, true negative, false negative and false positive, respectively. All the experiments were conducted on three-folds cross validation with stratified label. Due to the usage of cross-validation in our experiment, the average and standard deviation of all metrics were also reported.

3 Results

3.1 Tissue Classification

Tissue classifier model achieved accuracy of 0.990 and F1-score of 0.970 on the test set. This model then used on the LMP1 cohort to stratify tumour patches. We extracted non-overlapping patches with the size of 224×224 at a 20x magnification level from LMP1 cohort. We conducted a segmentation on lower-resolution image to find foreground and background region in image. We use the segmentation result as a guidance to segregate patches with tissue and background patches. Only patches contains tissue were inferred in the tumour classification process.

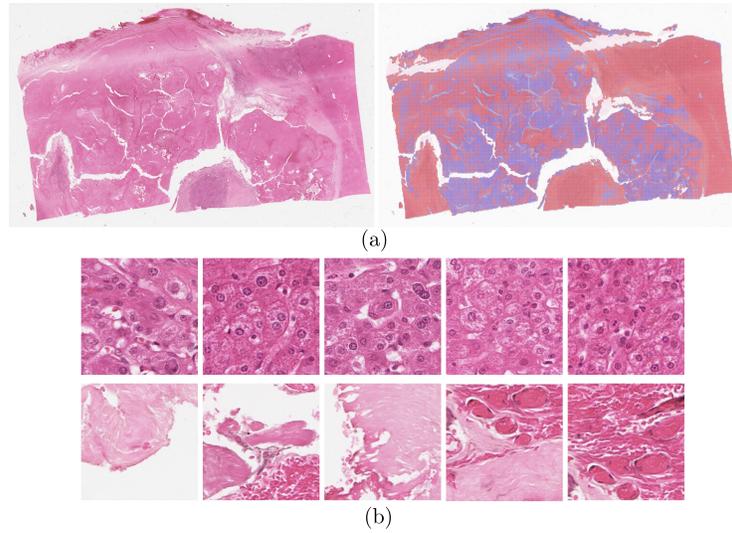


Fig. 2. Tissue Classification Result. (a) Top left image is a thumbnail of one of slide in the dataset and top right image is the corresponding tumour region heatmap. (b) Samples of tumour patches is depicted in the top row and non-tumour region in the bottom row

As can be seen in Figure 2, tumour patches were separated nicely in the WSIs. There was some minor misclassification in some slides due domain shift problem, but the rate was too low and could be ignored. Backbone networks encoded these tumour patches into vector representation and used them as bag of instances for LMP1 status prediction.

3.2 LMP1 Prediction

We argue that the tumour region is more important than any other tissue region for predicting LMP1. The non-tumour region may hinder the performance of

LMP1 prediction and weighs irrelevant regions. To this end, we compared the performance of model with all tissue region as bag of instance vs model which only employs tumour patches as bags of instance. We also compared two backbone networks, namely ResNet-18 and EfficientNet-B0 in this experiment. The result of this experiment can be seen in Table 2. There were differences in performance between the same backbone networks. The model that trained with only the tumour region delivered better performance predicting LMP1 status than the model that trained with all tissue regions. Furthermore, EfficientNet-B0 was better than ResNet-18 in terms of the backbone network. The best performance was achieved by using tumour regions as instances and EfficientNet-B0. This model achieved average accuracy, AUC and F1-score of 0.936, 0.995 and 0.862, respectively.

Table 2. LMP1 Status Prediction Results

Tissue Region	Backbone Network	Accuracy	AUC	F1-score
All region	ResNet-18	0.748 \pm 0.035	0.877 \pm 0.013	0.310 \pm 0.257
All region	EfficientNet-B0	0.926 \pm 0.032	0.978 \pm 0.014	0.825 \pm 0.095
Tumour region	ResNet-18	0.831 \pm 0.082	0.954 \pm 0.023	0.423 \pm 0.364
Tumour region	EfficientNet-B0	0.936 \pm 0.030	0.995 \pm 0.002	0.862 \pm 0.076

We also compared our result with other known methods such as CLAM [13] and MI-Net [18]. In the original pipeline, CLAM uses three sets of data, namely training, validation and testing sets. We modify the CLAM pipeline only to use the training and testing set to allow a fair comparison. Two aggregation methods were examined for MI-Net: the max and mean methods. As for the feature vectors in MI-Net, we use the same as our proposed model, which feature vectors from tumour patches encoded by EfficientNet-B0. The number of outputs of fully-connected layers in the MI-Net are 256, 128 and 64, respectively. MI-Net was trained with Adam optimizer and learning rate of 1×10^{-4} . We also ensure CLAM and MI-Net use the same data as ours in the cross-validation training.

Table 3. Performance Comparison of Our Model

Model	Tissue Region	Accuracy	AUC	F1-score
CLAM [13]	All region	0.723 \pm 0.022	0.548 \pm 0.092	-
MI-Net Max [18]	All region	0.852 \pm 0.011	0.700 \pm 0.018	0.571 \pm 0.037
MI-Net Mean [18]	All region	0.762 \pm 0.002	0.541 \pm 0.016	0.165 \pm 0.076
MI-Net Max [18]	Tumour region	0.832 \pm 0.031	0.661 \pm 0.054	0.477 \pm 0.130
MI-Net Mean [18]	Tumour region	0.767 \pm 0.006	0.559 \pm 0.041	0.207 \pm 0.136
Our model	Tumour region	0.936 \pm 0.030	0.995 \pm 0.002	0.862 \pm 0.076

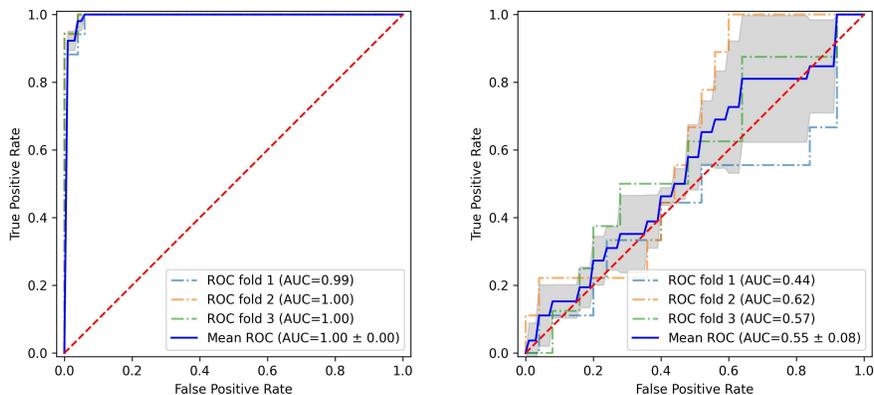


Fig. 3. AUROC curve Left image is AUROC curve for our model and right image is AUROC curve for CLAM

Based on the performance comparison in Table 3 and AUROC in Figure 3, CLAM performed poorly on this particular task. The average accuracy score achieved was only 0.723, with average AUC of 0.548. On the other hand, MI-Net performed well in all experiments compared to the CLAM. Both max and mean aggregation methods surpass CLAM performance in accuracy and AUC. However, aggregation with the max method was better than the mean method in this case. This may occur because the mean method averages from many other irrelevant instances while the max method selects only relevant instances. Our model achieved the best accuracy, AUC and F1-score compared to other models. The second-best model was MI-Net with max-aggregation that trained on tumour regions of the tissue. Both of MI-Net and our models performed better when only tumour regions was used as bag of instances. This indicates that tumour region is more relevant to predicting LMP1 in our case.

Utilising all regions in the tissue may generate an unmeaningful result. This example can be seen in Figure 4. There are two models which trained with a different formulation of the bag of instances. The right image is the attention heatmap of the model, which is trained on all regions of tissue, and the left image is trained on tumour regions. The model on the right side weighs blood cell regions with a high attention score, which is irrelevant in predicting LMP1 status.

3.3 Model Interpretability

The advantage of attention-based multiple instance learning is the interpretability of the model. We can inspect the relative importance among instances in regards to bag labels. Attention heatmaps from both of LMP1 negative and positive can be seen in Figure 5. This figure also visualised the patches with top 10%

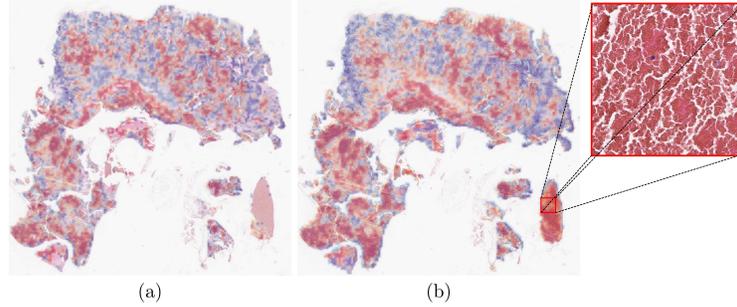


Fig. 4. Prediction comparison on different instance types. LMP1 heatmap prediction with tumour patches as instances(a) and all patches as instances(b). Region with high attention score in(b) was blood cell, which irrelevant region regarding to LMP1.

attention scores from both LMP1 classes. Based on the colour heterogeneity of patches, colour variation from the staining method was uncorrelated to LMP1 status.

We also conducted a simple test to prove that the multiple attention layer chooses the most relevant instances in regards to the bag label. Currently, there is no well-defined explanation regarding LMP1 effect on tumour micro-environment and its morphology. Therefore, we use MNIST as an example. We trained the same MIL architecture as the LMP1 status prediction task on the MNIST dataset. We defined the task into a binary classification. Bags which contained at least one image of digit 1 were defined as the positive class, and bags without image of digit 1 were defined as the negative class. We generated 5000 bags with a balanced class distribution. Each bags contain ten images/instances. The number of digits 1 was randomly selected in the positive class. We trained our model with the same learning rate 1×10^{-4} and optimized it with the Adam algorithm. To encode the image digits in MNIST into a features vector, we flatten the image of 28×28 into a 1-dimensional array of 784 elements. We then changed the scale of the feature vector from 0-255 to 0-1.

Figure 6 depicts an example of inferred results from a positive bag. Bag (a) contains two positive instances/two images of digit 1. These instances have the highest attention scores among other instances in the first aggregation layers. In the second attention layer, the attention scores for negative instances were reduced while the scores for positive instances were increased. This trend continues until the last aggregation layer. Thus all negative instances have zero value of attention. In another case, in bag (b), the attention scores in the last layers were averaged among the positive instances. This experiment proves that multiple attention-based of aggregation layers help the model to select more relevant instances.

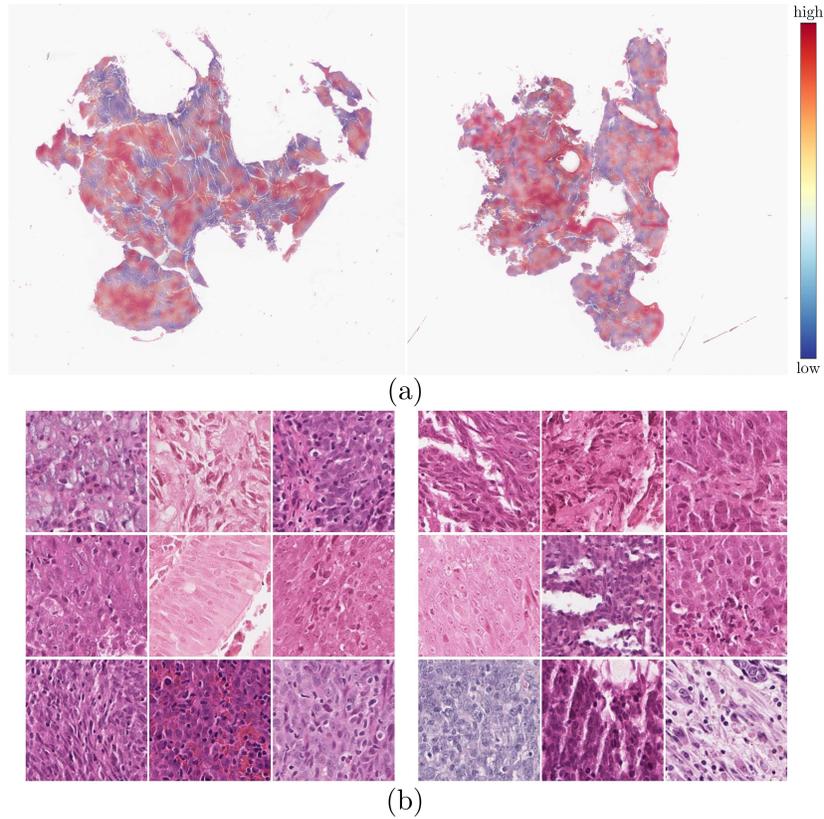


Fig. 5. Attention heatmap of LMP1 prediction model. (a) is attention heatmap, the left image is LMP1 negative and the right image is LMP1 positive. (b) is randomly selected patches from dataset with highest attention score, left images is true negative cases and right images is true positive cases

4 Conclusion

In this paper, we examine the use of deep learning for LMP1 status prediction in NPC patients using the H&E-stained WSIs. LMP1 data in this study was collected from a University Hospital with a total number of 101 cases. We proposed multi attention aggregation layer for MIL in the tumour region to predict LMP1 status. There was an increase in the model performance when using only tumour regions as instances. Despite the simplicity of our proposed method, it outperformed the other known MIL models. To the best of our knowledge, this is the first attempt to predict LMP1 status on NPC using deep learning.

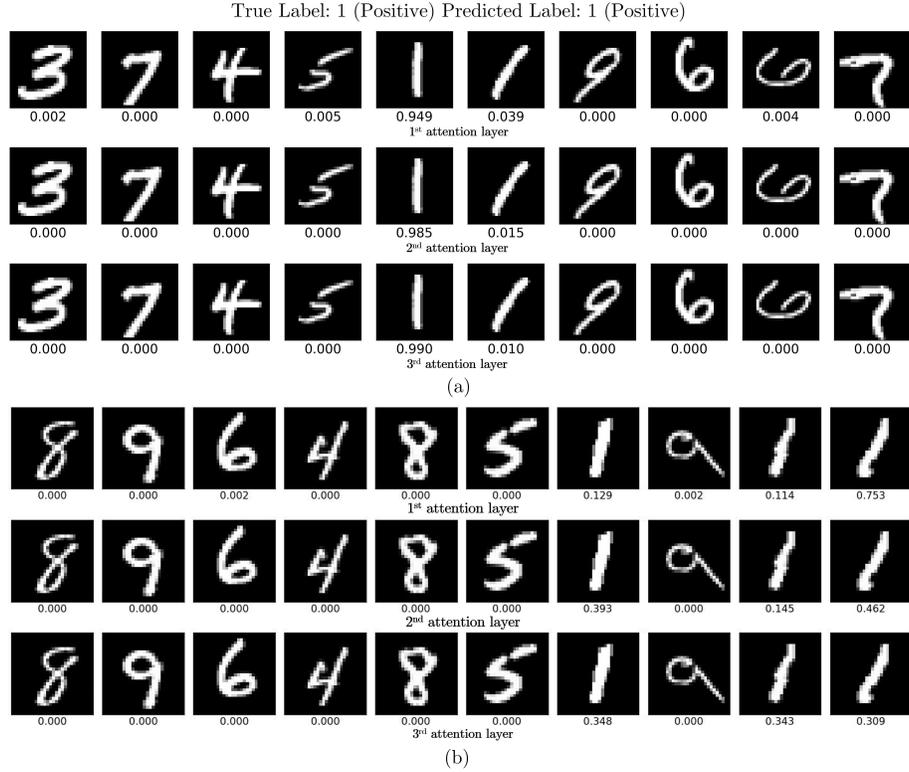


Fig.6. Attention Scores from Three Aggregation Layers These are two examples of bags in MNIST (a,b) with their corresponding attention scores in three aggregation layers. Number below each digits is the attention scores.

5 Acknowledgements

This study is fully supported by a PhD scholarship to the first author funded by Indonesia Endowment Fund for Education (LPDP), Ministry of Finance, Republic of Indonesia under grant number Ref: S-575/LPDP.4/2020.

References

1. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473 (2014)
2. Campanella, G., Hanna, M.G., Geneslaw, L., Mirafior, A., Werneck Krauss Silva, V., Busam, K.J., Brogi, E., Reuter, V.E., Klimstra, D.S., Fuchs, T.J.: Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nature medicine* **25**(8), 1301–1309 (2019)
3. Carbonneau, M.A., Cheplygina, V., Granger, E., Gagnon, G.: Multiple instance learning: A survey of problem characteristics and applications. *Pattern Recognition* **77**, 329–353 (2018). <https://doi.org/https://doi.org/10.1016/j.patcog.2017.10.009>, <https://www.sciencedirect.com/science/article/pii/S0031320317304065>
4. Coudray, N., Ocampo, P.S., Sakellaropoulos, T., Narula, N., Snuderl, M., Fenyö, D., Moreira, A.L., Razavian, N., Tsirigos, A.: Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning. *Nature medicine* **24**(10), 1559–1567 (2018)
5. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: *IEEE conference on computer vision and pattern recognition*. pp. 248–255. IEEE (2009)
6. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *IEEE conference on computer vision and pattern recognition*. pp. 770–778 (2016)
7. Ilse, M., Tomczak, J., Welling, M.: Attention-based deep multiple instance learning. In: *International conference on machine learning*. pp. 2127–2136. PMLR (2018)
8. Kather, J.N., Krisam, J., Charoentong, P., Luedde, T., Herpel, E., Weis, C.A., Gaiser, T., Marx, A., Valous, N.A., Ferber, D., et al.: Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study. *PLoS medicine* **16**(1), e1002730 (2019)
9. Kieser, A., Sterz, K.R.: The latent membrane protein 1 (Imp1). *Epstein Barr Virus Volume 2* pp. 119–149 (2015)
10. Klein, S., Quaas, A., Quantius, J., Löser, H., Meinel, J., Peifer, M., Wagner, S., Gattenlöhner, S., Wittekindt, C., von Knebel Doeberitz, M., et al.: Deep learning predicts hpv association in oropharyngeal squamous cell carcinomas and identifies patients with a favorable prognosis using regular h&e stains. *Deep learning predicts hpv association in opsc. Clinical Cancer Research* **27**(4), 1131–1138 (2021)
11. Lee, A.W., Lung, M.L., Ng, W.T.: *Nasopharyngeal carcinoma: from etiology to clinical practice*. Academic Press (2019)
12. Liang, J., Cui, Z., Wu, C., Yu, Y., Tian, R., Xie, H., Jin, Z., Fan, W., Xie, W., Huang, Z., Xu, W., Zhu, J., You, Z., Guo, X., Qiu, X., Ye, J., Lang, B., Li, M., Tan, S., Hu, Z.: Deepebv: a deep learning model to predict epstein-barr virus (ebv) integration sites. *Bioinformatics* (5 2021). <https://doi.org/10.1093/bioinformatics/btab388>
13. Lu, M.Y., Williamson, D.F., Chen, T.Y., Chen, R.J., Barbieri, M., Mahmood, F.: Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature Biomedical Engineering* **5**(6), 555–570 (2021)
14. Muti, H.S., Heij, L.R., Keller, G., Kohlruss, M., Langer, R., Dislich, B., Cheong, J.H., Kim, Y.W., Kim, H., Kook, M.C., et al.: Development and validation of deep learning classifiers to detect epstein-barr virus and microsatellite instability status in gastric cancer: a retrospective multicentre cohort study. *The Lancet Digital Health* **3**(10), e654–e664 (2021)

15. Qiu, S., Guo, Y., Zhu, C., Zhou, W., Chen, H.: Attention based multi-instance thyroid cytopathological diagnosis with multi-scale feature fusion. In: 2020 25th International Conference on Pattern Recognition (ICPR). pp. 3536–3541. IEEE (2021)
16. Schaumberg, A.J., Rubin, M.A., Fuchs, T.J.: H&e-stained whole slide image deep learning predicts spop mutation state in prostate cancer. *BioRxiv* p. 064279 (2018)
17. Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: International conference on machine learning. pp. 6105–6114. PMLR (2019)
18. Wang, X., Yan, Y., Tang, P., Bai, X., Liu, W.: Revisiting multiple instance neural networks. *Pattern Recognition* **74**, 15–24 (2018)
19. Wong, K.C., Hui, E.P., Lo, K.W., Lam, W.K.J., Johnson, D., Li, L., Tao, Q., Chan, K.C.A., To, K.F., King, A.D., et al.: Nasopharyngeal carcinoma: an evolving paradigm. *Nature Reviews Clinical Oncology* **18**(11), 679–695 (2021)
20. Zhang, H., Meng, Y., Zhao, Y., Qiao, Y., Yang, X., Coupland, S.E., Zheng, Y.: Dtf-d-mil: Double-tier feature distillation multiple instance learning for histopathology whole slide image classification. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18802–18812 (2022)
21. Zheng, X., Wang, R., Zhang, X., Sun, Y., Zhang, H., Zhao, Z., Zheng, Y., Luo, J., Zhang, J., Wu, H., et al.: A deep learning model and human-machine fusion for prediction of ebv-associated gastric cancer from histopathology. *Nature communications* **13**(1), 1–12 (2022)