Supplemental Material: Differentiable Zooming for Multiple Instance Learning on Whole-Slide Images

Kevin Thandiackal^{1,2}*[©], Boqi Chen^{1,2}*[©], Pushpak Pati¹[©], Guillaume Jaume³[†][©], Drew F. K. Williamson³[©], Maria Gabrani¹[©], and Orcun Goksel^{2,4}[©]

 ¹ IBM Research Europe, Zurich, Switzerland {kth,che,pus,mga}@zurich.ibm.com
 ² ETH Zurich, Zurich, Switzerland
 ³ Brigham and Women's Hospital, Harvard Medical School, Boston, USA {gjaume,dwilliamson}@bwh.harvard.edu
 ⁴ Uppsala University, Uppsala, Sweden orcun.goksel@it.uu.se

In this supplemental material, we include additional results, visualizations, and analyses. The contents of the individual sections are:

- Appendix A: Confusion matrices for ZOOMMIL on all datasets. Additional results for the competing methods operating at different magnifications
- Appendix B: Analyzing the limitations of ZOOMMIL
- Appendix C: Analyzing the impact of the number of selected patches in the differentiable Top-K module
- Appendix D: ZOOMMIL attention maps at different magnifications
- Appendix E: Training details of the competing baselines
- Appendix F: Derivation of the Jacobian for differentiable patch selection

A Additional Classification Analysis

In Figure 1, we present the confusion matrices of ZOOMMIL on all datasets. Results are averaged over three runs with different weight initializations. On CRC, ZOOMMIL performs very well and correctly classifies 95.63% and 94.55% of the non-neoplastic and low-grade cases, respectively. Out of all high-grade, cases, our model identifies 83.33% correctly. We can see that BRIGHT is the most difficult dataset due to its challenging pre-cancerous class, which is often confused with either the non-cancerous or the cancerous class. On CAMELYON16, ZOOMMIL accurately identifies 97.92% of non-metastatic cases, while correctly classifying 61.90% of metastatic cases.

Additionally, we report the classification performance of all single-scale baselines operating at different magnifications in Table 1, 2, and 3. It can be observed in Table 1 and 2 that lower magnifications severely impact the performance of the baselines, while ZOOMMIL-EFF performs significantly better. These results indicate the efficacy of our method, and conclude the benefits of zooming and combining information across magnifications.

^{*} Contributed equally.

[†] Work done while at IBM Research Europe



Fig. 1. Confusion matrices of ZOOMMIL on CRC, BRIGHT, and CAMELYON16. An entry represents the corresponding fraction (%) w.r.t. all samples in the same row

Methods	Classification Weighted-F1($\%$) Accuracy($\%$)	
ABMIL [6] $(5\times)$ ABMIL [6] $(10\times)$ CLAM-SB [10] $(5\times)$ CLAM-SB [10] $(10\times)$ TRANSMIL [12] $(5\times)$ TRANSMIL [12] $(10\times)$	$\begin{array}{c} 86.8{\pm}0.7\\ 88.8{\pm}0.7\\ 87.7{\pm}0.5\\ 89.5{\pm}0.5\\ 86.2{\pm}1.1\\ 88.4{\pm}1.3\end{array}$	$\begin{array}{c} 87.0 \pm 0.7 \\ 89.0 \pm 0.6 \\ 87.8 \pm 0.5 \\ 89.6 \pm 0.5 \\ 87.4 \pm 1.1 \\ 89.1 \pm 1.1 \end{array}$
$\frac{\text{ZoomMIL-Eff} (5 \times \to 10 \times)}{\text{ZoomMIL} (5 \times \to 10 \times \to 20 \times)}$	90.3 ± 1.3 92.0 ± 0.6	90.3 ± 1.3 92.1 ± 0.7

Table 1. Classification performances on the CRC dataset [11]

Table 2. Classification performances on the BRIGHT dataset [3]

Methods	Classification Weighted-F1(%) Accuracy(%)	
ABMIL [6] $(1.25 \times)$ ABMIL [6] $(2.5 \times)$ CLAM-SB [10] $(1.25 \times)$ CLAM-SB [10] $(2.5 \times)$ The work MU [12] $(1.25 \times)$	58.4 ± 1.0 58.7 ± 1.1 59.9 ± 1.3 60.1 ± 1.2	58.9 ± 1.6 59.3 ± 1.0 60.3 ± 1.2 60.2 ± 1.6 472 ± 1.6
TRANSMIL [12] (1.25×) TRANSMIL [12] (2.5×) ZOOMMIL-EFF (1.25× \rightarrow 2.5×)	$ \begin{array}{r} 40.1 \pm 3.8 \\ 52.0 \pm 1.3 \\ \hline 66.0 \pm 1.9 \end{array} $	47.3 ± 2.5 54.5 ± 2.7 66.5 ± 1.5
ZOOMMIL $(1.25 \times \rightarrow 2.5 \times \rightarrow 10 \times)$	68.3 ± 1.1	69.3 ± 1.0

B Limitations

We conjecture that our classification performance on CAMELYON16 is limited by the size of the metastatic regions. For the samples including extremely small

Methods	Classification Weighted-F1(%) Accuracy(%)	
ABMIL [6] $(10\times)$ CLAM-SB [10] $(10\times)$ TRANSMIL [12] $(10\times)$	$76.7{\pm}0.8$ $77.5{\pm}0.6$ $76.6{\pm}1.1$	78.3 ± 0.7 79.1 ± 0.6 79.6 ± 1.0
$\frac{1}{\text{ZOOMMIL } (10 \times \rightarrow 20 \times)}$	$83.3{\pm}0.3$	84.2 ± 0.4

Table 3. Classification performances on the CAMELYON16 dataset [1]

Table 4. Classification performances on the validation set of CAMELYON16 [1]. Validation set is further grouped according to the size of the metastatic regions

	Classification	
Size of metastases	Weighted-F1($\%$)	Accuracy(%)
large	96.2 ± 1.0	$96.1 {\pm} 1.1$
small	$86.7 {\pm} 1.0$	87.1 ± 1.0

metastases, it is challenging to optimize the zooming process. To validate our hypothesis, we sub-categorize the metastatic samples in the training set into "small" and "large" metastatic groups, via visual inspection, and create new stratified training and validation sets. In Table 4, we present the classification performances of ZOOMMIL on the validation set, individually for the small and large metastatic samples. The results show that the performance is significantly higher for large metastases, which substantiates our hypothesis.

C Impact of K in Differentiable Top-K Patch Selection

Here, we analyze the impact of the number of selected patches (K) on the classification performance of ZOOMMIL. Figure 2 shows that the performance increases with increasing K. It peaks at K = 12 and then slightly drops for further increment. We reason that this behavior is caused by the average number of patches per Whole-Slide Image (WSI) being 16 in the BRIGHT dataset at the lowest magnification $(1.25\times)$. Almost all patches are selected in this case, which makes it suboptimal to learn to improve the zooming process.

D Interpretability

We have shown the ZOOMMIL attention maps at the lowest magnification, *i.e.*, $1.25\times$, on the BRIGHT dataset in Figure 4 of the main manuscript. Here, in Figure 3, we provide the attention maps at both low and high magnification, *i.e.*, $1.25\times$ and $10\times$, as well as close-up patches with highest attention scores. We can observe in the zooming process from $1.25\times$ to $10\times$, that the model pins down its focus to the most informative regions in the WSI. The observation is



Fig. 2. Classification performance of ZOOMMIL on BRIGHT for different values of K



Fig. 3. Interpreting ZOOMMIL on the BRIGHT dataset: (a) annotated tumor regions, (b) attention maps at $1.25 \times$ magnification, (c) attention maps at $10 \times$ magnification, and (d) a subset of extracted patches with high attention scores at $10 \times$ magnification

further substantiated by the high attention patches that include cancerous content, *i.e.*, invasive tumors and ductal carcinoma in situ tumors, for the presented cancerous WSIs. Figure 4 shows more examples of annotated tumor regions and attention maps for correctly classified and misclassified samples across all the classes in BRIGHT. For both the correct and incorrect classifications, we can notice that the focus of the model is aligning with the focus of the pathologist. However, the morphological ambiguities among the classes finally lead to certain misclassifications.

E Implementation Details

ABMIL: We follow the network architecture proposed in [6]. The model comprises of a gated attention module consisting of three 2-layer Multi-Layer Perceptrons (MLPs), where the first two are followed by Hyperbolic Tangent and Sigmoid activations, and a 2-layer MLP classifier with ReLU activation.

CLAM-SB: We implement CLAM-SB [10] with the code¹ provided by the authors. We use the Adam optimizer [7] with 0.0001 learning rate. The maximum and minimum epochs are set to 100, 50, respectively, and use early stopping with patience=20 epochs based on the validation weighted F1 score (for BRIGHT) and validation loss (for CRC and CAMELYON16). We use cross-entropy loss for both bag loss and instance loss. The weights of bag-level losses are set to 0.7, 0.5, 0.7, and the number of positive/negative patches sampled for instance loss are set to 8, 32, 32 for CRC, BRIGHT and CAMELYON16, respectively. For all three datasets we use a weighted sampler due to class imbalance.

TransMIL: We adopt the original implementation² of TransMIL [12]. We use Lookahead optimizer [13] with learning rate 0.0002 and weight decay 0.00001. The maximum epochs are set to 200 and early stopping is used with a *patience* of 10 epochs based on validation weighted F1 score (for BRIGHT) and validation loss (for CRC and Camelyon). We use cross-entropy loss as training loss.

MSMIL: We adopt the model architecture from the original implementation³ of MSMIL [5]. We construct the multi-scale feature matrix $f \in \mathbb{R}^{(N_0+\dots+N_{n-1})\times d}$ by concatenating the feature matrices f_0, \dots, f_{n-1} at all magnifications m_0, \dots, m_{n-1} where $f_n \in \mathbb{R}^{N_n \times d}$ and $N_n = N_{n-1}(\frac{m_n}{m_{n-1}})^2$ is the number of features at magnification m_n . Here, d is the feature dimension and $n = \{3, 3, 2\}$ is the number of considered magnifications on CRC, BRIGHT, and CAMELYON16, respectively. We use the Adam optimizer [7] with learning rate 0.0001 and plateau scheduler (patience= 5 epochs, decay rate= 0.8). The experiments are run for 100 epochs with a batch size of one. The models with the best weighted F1-score (for BRIGHT) and best loss (for CRC & CAMELYON16) on the validation set are saved for testing.

DSMIL: We adopt the model architecture from the original implementation⁴ of DSMIL [9]. We concatenate the d-dimensional feature vector of each patch

¹ https://github.com/mahmoodlab/CLAM

² https://github.com/szc19990412/TransMIL

³ https://github.com/takeuchi-lab/MS-DA-MIL-CNN

⁴ https://github.com/binli123/dsmil-wsi

6 Thandiackal et al.



 ${\bf Fig.~4.}$ Examples of BRIGHT WSIs with annotated tumor regions and attention maps from ZOOMMIL at 1.25× magnification

at low-magnification m_0 with d-dimensional feature vectors of its corresponding sub-patches at higher magnifications m_1, \dots, m_{n-1} to form an *nd*-dimensional multi-scale feature vector where $n = \{3, 3, 2\}$ is the number of considered magnifications on CRC, BRIGHT, and CAMELYON16, respectively. Each feature vector from a lower magnification patch is replicated and concatenated to each of its higher magnification counterparts. We use the Adam optimizer [7] with learning rate 0.0001 and plateau scheduler (patience = 5 epochs, decay rate = 0.8). The experiments are run for 100 epochs with a batch size of one. The models with the best weighted F1-score (for BRIGHT) and best loss (for CRC & CAMELYON16) on the validation set are saved for testing.

SparseConvMIL: We adopt the original implementation⁵ of SparseConvMIL [8]. Unlike other baselines, we run the model on a V100 GPU with 32GB RAM, as it depends on SparseConvNet^6 . Therefore, we limit the batch size and the number of sampled patches to 8 and 100, respectively. In the model, we set the number of sparseconv channels to 32, the downsampling factor of the sparse map to 128, and the neurons in the MLP classifier to 128. We use ResNet34 for extracting patch features, and finetune it with learning rate 0.00001. The learning rate for the rest of the model is 0.001, and weight decay is 0.0001.

MaxMIL & MeanMIL: We use the formulation presented in [8] for Max and Mean MIL. We set the same values for the hyperparameters as in SparseConvMIL, but freeze the patch feature extractor, *i.e.*, ResNet34.

\mathbf{F} Derivation of Jacobian for Differentiable Patch Selection

Perturbed Maximum As described in [2], given a set of distinct points $\mathcal{T} \subset \mathbb{R}^d$ and its convex hull \mathcal{C} , a discrete optimization problem with inputs $\mathbf{a}_m \in \mathbb{R}^d$ can generally be formulated as:

$$\max_{\hat{\mathbf{t}}\in\mathcal{C}} \langle \hat{\mathbf{t}}, \mathbf{a}_m \rangle \qquad \mathbf{t} = \underset{\hat{\mathbf{t}}\in\mathcal{C}}{\arg \max} \langle \hat{\mathbf{t}}, \mathbf{a}_m \rangle . \tag{1}$$

As per Definition 2.1 in [2], we can obtain a smoothed \mathbf{t} by adding a random noise vector $\sigma \mathbf{Z} \in \mathbb{R}^d$ with distribution $d\mu(\mathbf{z}) \propto \exp(-\nu(\mathbf{z})) d\mathbf{z}$, where $\sigma > 0$ is a scaling parameter. The perturbed version of the maximizer in Eq. (1) then becomes:

$$\mathbf{t} = \mathbb{E}\left[\arg\max_{\hat{\mathbf{t}} \in \mathcal{C}} \langle \hat{\mathbf{t}}, \mathbf{a}_m + \sigma \mathbf{Z} \rangle \right].$$
(2)

According to Proposition 3.1 from [2], the associated Jacobian matrix of t at \mathbf{a}_m can then be computed as follows:

$$J_{\mathbf{a}_m} \mathbf{t} = \mathbb{E} \left[\arg \max_{\hat{\mathbf{t}} \in \mathcal{C}} \langle \hat{\mathbf{t}}, \mathbf{a}_m + \sigma \mathbf{Z} \rangle \nabla_{\mathbf{z}} \nu(\mathbf{Z})^\top / \sigma \right].$$
(3)

⁵ https://github.com/MarvinLer/SparseConvMIL
⁶ https://github.com/facebookresearch/SparseConvNet

8 Thandiackal et al.

We choose our noise to have a normal distribution, *i.e.*, $\mathbf{Z} \sim \mathcal{N}(0, 1)$. We can thus plug $\nabla_{\mathbf{z}} \nu(\mathbf{Z})^{\top} = \mathbf{Z}^{\top}$ into Eq. (3) and obtain:

$$J_{\mathbf{a}_m} \mathbf{t} = \mathop{\mathbb{E}}_{\mathbf{Z} \sim \mathcal{N}(0,1)} \left[\arg \max_{\hat{\mathbf{t}} \in \mathcal{C}} \langle \hat{\mathbf{t}}, \mathbf{a}_m + \sigma \mathbf{Z} \rangle \mathbf{Z}^\top / \sigma \right].$$
(4)

Differentiable Top-K Operator As shown in [4], the Top-K selection with sorted indices can be converted into the same form as Eq. (1). To this end, the constraint set C for the indicator matrix $\hat{\mathbf{T}}$ should first be defined as:

$$\mathcal{C} = \begin{cases} \hat{\mathbf{T}} \in \mathbb{R}^{N \times K} : & \hat{\mathbf{T}}_{n,k} \ge 0 \end{cases}$$
(5)

$$\sum_{j=1}^{K} \hat{\mathbf{T}}_{j,k} = 1 \quad \forall k \in \{1, \dots, K\}$$
(6)

$$\sum_{k=1}^{K} \hat{\mathbf{T}}_{j,k} \le 1 \quad \forall j \in \{1, \dots, N\}$$

$$\tag{7}$$

$$\sum_{i=1}^{N} i \hat{\mathbf{T}}_{i,k} < \sum_{j=1}^{N} j \hat{\mathbf{T}}_{j,k'} \quad k < k' \Big\} , \qquad (8)$$

where Eq. (6) ensures that each column-wise sum in the indicator matrix is equal to one and Eq. (7) constrains each row-wise sum to be at most one. Lastly, Eq. (8) enforces that the indices of the attention weights selected by $\hat{\mathbf{T}}$ are sorted. With these constraints, the general linear program formulation in Eq. (1) can be used to describe the Top-K selection problem:

$$\max_{\hat{\mathbf{T}}\in\mathcal{C}} \langle \hat{\mathbf{T}}, \mathbf{a}_m \mathbf{1}^\top \rangle \qquad \mathbf{T} = \arg\max_{\hat{\mathbf{T}}\in\mathcal{C}} \langle \hat{\mathbf{T}}, \mathbf{a}_m \mathbf{1}^\top \rangle , \qquad (9)$$

where $\mathbf{1}^{\top} = [1 \cdots 1] \in \mathbb{R}^{1 \times K}$ and thus $\mathbf{a}_m \mathbf{1}^{\top} \in \mathbb{R}^{N \times K}$ is a matrix containing the attention vectors \mathbf{a}_m repeated K times as its columns. Note that here, $\langle \cdot \rangle$ computes the scalar product after vectorizing the matrices.

Now, the perturbed maximizer \mathbf{T} and the corresponding Jacobian $J_{\mathbf{a}_m}\mathbf{T}$ can be computed analogously to Eq. (2) and Eq. (4), as presented in the main paper. In contrast to [2], however, for the Top-K selection problem, it is required to apply the same noise vector $\sigma \mathbf{Z}$ to each column in $\mathbf{a}_m \mathbf{1}^{\top}$. Following the insights from [4], we therefore apply in practice the noise directly to \mathbf{a}_m , *i.e.*, we employ $(\mathbf{a}_m + \sigma \mathbf{Z})\mathbf{1}^{\top}$ instead of $\mathbf{a}_m\mathbf{1}^{\top} + \sigma \mathbf{Z}$.

References

- Bejnordi, B., Veta, M., van Dienst, P., van Ginneken, B., Karssemeijer, N., Litjens, G., van der Laak, J., et al.: Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. JAMA **318**, 2199–2210 (2017)
- Berthet, Q., Blondel, M., Teboul, O., Cuturi, M., Vert, J., Bach, F.: Learning with differentiable pertubed optimizers. In: Advances in Neural Information Processing Systems (NeurIPS). vol. 34, pp. 9508–9519 (2020)
- Brancati, N., Anniciello, A., Pati, P., Riccio, D., Scognamiglio, G., Jaume, G., De Pietro, G., Di Bonito, M., Foncubierta-Rodríguez, A., Botti, G., et al.: Bracs: A dataset for breast carcinoma subtyping in h&e histology images. arXiv:2111.04740 (2021)
- Cordonnier, J., Mahendran, A., Dosovitskiy, A.: Differentiable patch selection for image recognition. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2351–2360 (2021)
- Hashimoto, N., Fukushima, D., Koga, R., Takagi, Y., Ko, K., Kohno, K., Nakaguro, M., Nakamura, S., Hontani, H., Takeuchi, I.: Multi-scale domainadversarial multiple-instance cnn for cancer subtype classification with unannotated histopathological images. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3852–3861 (2020)
- Isle, M., Tomczak, J., Welling, M.: Attention-based deep multiple instance learning. In: International Conference on Machine Learning (ICML). vol. 35 (2018)
- 7. Kingma, D., Ba, J.: Adam: A method for stochastic optimization. In: International Conference on Learning Representations (ICLR) (2015)
- Lerousseau, M., Vakalopoulou, M., Deutsch, E., Paragios, N.: Sparseconvmil: Sparse convolutional context-aware multiple instance learning for whole slide image classification. In: MICCAI Workshop on Computational Pathology. pp. 129–139 (2021)
- Li, B., Li, Y., Eliceiri, K.: Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 14318– 14328 (2021)
- Lu, M., Williamson, D., Chen, T., Chen, R., Barbieri, M., Mahmood, F.: Data efficient and weakly supervised computational pathology on whole slide images. Nature Biomedical Engineering 5, 555–570 (2021)
- Oliveira, S., Neto, P., Fraga, J., Montezuma, D., Monteiro, A., Monteiro, J., Ribeiro, L., Gonçalves, S., Pinto, I., Cardoso, J.: Cad systems for colorectal cancer from wsi are still not ready for clinical acceptance. Scientific Reports 11 (2021)
- Shao, Z., Bian, H., Chen, Y., Wang, Y., Zhang, J., Ji, X., Zhang, Y.: Transmil: Transformer based correlated multiple instance learning for whole slide image classification. In: Advances in Neural Information Processing Systems (NeurIPS). vol. 35 (2021)
- Zhang, M., Lucas, J., Hinton, G., Ba, J.: Lookahead optimizer: k steps forward, 1 step back. In: Advances in Neural Information Processing Systems (NeurIPS). vol. 33 (2019)