

469 **A Proof**

470 *Proof.* Since the bag of instances  $\mathbf{X}$  is sampled from the probability distribution  $\mu(\mathbf{x})$ , we have the  
 471 upper bound for Wasserstein distance between  $\mathbf{X}$  and  $\mu$  [30],

$$\mathbb{E}[\mathcal{W}_p(\mathbf{X}, \mu)] \leq K^{-\frac{1}{d\mu}}, \quad (6)$$

472 where  $K$  is the number of samples. Next we define the function  $\sigma(\mathbf{x})$  as a small perturbation function  
 473  $\sigma(\mathbf{x}) = \mathbf{x} + \delta$  and let  $\tilde{\mathbf{X}} = \{\sigma(\mathbf{x}) : \mathbf{x} \in \mathbf{X}\}$ . Using the triangle inequality, we have

$$\mathbb{E}[\mathcal{W}_p(\mathbf{X}, \tilde{\mathbf{X}})] \leq \mathbb{E}[\mathcal{W}_p(\mathbf{X}, \mu)] + \mathbb{E}[\mathcal{W}_p(\tilde{\mathbf{X}}, \mu)] \leq 2K^{-\frac{1}{d\mu}} + C, \quad (7)$$

474 where  $C$  is a constant as a result of the perturbation. As the function  $S(\cdot)$  is Lipschitz continuous, we  
 475 have

$$|S(\mathbf{X}) - S(\tilde{\mathbf{X}})| \leq L \cdot \mathbb{E}[\mathcal{W}_p(\mathbf{X}, \tilde{\mathbf{X}})] \leq O(L \cdot K^{-\frac{1}{d\mu}}). \quad (8)$$

476 Similar to TransMIL [5], let  $\Phi : \mathcal{X} \rightarrow \mathbb{R}^n$  be any invertible map, where its inverse mapping is  
 477 expressed as  $\Phi^{-1} : \mathbb{R}^d \rightarrow \mathcal{X}$ . Then we have:

$$S(\Phi^{-1}(\Phi_{\mathbf{X} \in \mathcal{X}}(\{\sigma(\mathbf{x}) : \mathbf{x} \in \mathbf{X}\}))) = S(\Phi^{-1}(\Phi_{\tilde{\mathbf{X}} \in \mathcal{X}}(\tilde{\mathbf{X}}))) = S(\tilde{\mathbf{X}}). \quad (9)$$

478 Let  $\gamma = S \circ \Phi^{-1}$ . As  $|S(\mathbf{X}) - S(\tilde{\mathbf{X}})| \leq O(L \cdot K^{-\frac{1}{d\mu}})$ , we have

$$|S(\mathbf{X}) - \gamma(\Phi_{\mathbf{X} \in \mathcal{X}}(\{\sigma(\mathbf{x}) : \mathbf{x} \in \mathbf{X}\}))| \leq O(L \cdot K^{-\frac{1}{d\mu}}). \quad (10)$$

479 □

480 In this proof, the transformation  $\Phi(\cdot) = \mathcal{A}(\cdot)$ . This proof could be easily extended to the represen-  
 481 tations  $\mathbf{H}$  by assuming a probability measure over the instance representations  $\mathbf{h}$  and replacing  $\mathbf{X}$   
 482 with  $\mathbf{H}$ . In this case, the transformation  $\Phi(\mathbf{H}) = \sum_{k=1}^K a_k \mathbf{h}_k$ .

Table 4: Results on general MIL datasets. Experiments were run 5 times and the average classification accuracy ( $\pm$  a standard error of a mean) is reported.

Method	MUSK1	MUSK2	FOX	TIGER	ELEPHANT
Attention	0.892±0.090	0.858±0.106	0.615±0.096	0.839±0.054	0.868±0.054
Attention-Gated	0.900±0.088	0.863±0.094	0.603±0.068	<b>0.845±0.046</b>	0.857±0.064
CLAM	0.900±0.136	0.860±0.128	0.610±0.128	0.805±0.052	0.860±0.080
RAM-MIL	<b>0.911±0.130</b>	<b>0.870±0.142</b>	<b>0.645±0.117</b>	0.820±0.040	<b>0.879±0.096</b>

483 **B General MIL dataset**

484 Table 4 presents the performance of RAM-MIL on general MIL datasets [31, 32], offering a compar-  
 485 ison with baseline methods. The results indicate that OT-based retrieval generally enhances  
 486 the classification performance. The sole exception is observed with the TIGER dataset, where  
 487 both CLAM and RAM-MIL are outperformed. This discrepancy might be attributed to CLAM, as  
 488 RAM-MIL uses CLAM as a pretrained model for attention weights and bag representation extraction.  
 489 Nonetheless, RAM-MIL still improves over its CLAM baseline on TIGER. Note that our primary  
 490 focus lies on the more challenging WSI datasets, hence our models are not extensively optimized for  
 491 general datasets. The data in these general datasets are typically of lower dimensionality and present  
 492 less challenging conditions. Therefore, any potential underperformance in these contexts should not  
 493 detract from the strength of our models in handling the WSI data.

494 **C Experiment Details of WSI Classification**

495 We present the experimental details, ablation studies and analysis step-by-step.

496 **MIL Pre-training.** For the backbone MIL model we use the the same parameter setup as CLAM.  
 497 The model parameters are updated via the Adam optimizer with an L2 weight decay of 1e-5 and a  
 498 learning rate of 2e-4. Each result is obtained with 10-fold splits of training/validation/testing sets.

Table 5: Ablation study for the percentage of instances used on CAMELYON16 and CAMELYON17.

	In-Domain (CAM16)		Out-of-Domain (CAM17)	
	AUC	Accuracy	AUC	Accuracy
10% attention Retr <sub>I</sub>	0.9440±0.037	0.8975±0.052	-	-
10% attention Retr <sub>IO</sub>	0.9365±0.052	<b>0.9200±0.050</b>	<b>0.7974±0.054</b>	0.7433±0.073
10% attention Retr <sub>O</sub>	0.9414±0.046	0.8975±0.056	0.7775±0.050	0.7392±0.063
20% attention Retr <sub>I</sub>	<b>0.9451±0.036</b>	0.8925±0.050	-	-
20% attention Retr <sub>IO</sub>	0.9341±0.051	0.8925±0.053	0.7651±0.056	0.7714±0.030
20% attention Retr <sub>O</sub>	0.9419±0.048	0.9175±0.051	0.7681±0.058	<b>0.7795±0.021</b>

499 **Neighbor Selection.** After pre-training the MIL model, we obtain the slide-level feature and the  
500 attention scores predicted by the network. As computing the optimal transport distance based on all  
501 instances is time-consuming, we approximate the distance with a part of samples in a bag.

$$d_{OT}(\mu, \nu) = \min_{T \in \mathcal{T}(\alpha, \tilde{\alpha})} \sum_{i=1}^{|\alpha|} \sum_{j=1}^{|\tilde{\alpha}|} c(\mathbf{h}_i, \tilde{\mathbf{h}}_j) T_{ij} + \beta \cdot \sum_{ij} T_{ij} \log T_{ij} \quad (11)$$

$$s.t., T^T \mathbf{1}_K = \alpha, T \mathbf{1}_{\tilde{K}} = \tilde{\alpha}, T \geq 0$$

502 where  $\alpha$  and  $\tilde{\alpha}$  are the new attention vector obtained by selecting top  $\eta\%$  from  $\mathbf{a}$  and  $\tilde{\mathbf{a}}$ . In other  
503 words, we approximate a bag with  $\eta\%$  of instances with the highest attention values generated by  
504 pretrained MIL model. As shown in Table 5, we set  $\eta = 10$  and  $\eta = 20$  for the ablation study. In this  
505 experiment, we use Regularization term of 0.5 and Max number of iterations 1000.

506 It is observed from Table 5, improving the percentage of data improves the performance on Retr<sub>I</sub> and  
507 Retr<sub>O</sub>. On Retr<sub>IO</sub>, the performance is saturated when using only 10% of all patches. Differentiating  
508 in-domain and out-of-domain data for retrieval could be easily accomplished by representing a bag  
509 by a few amount of instances.

Table 6: Ablation study for different merge functions on CAMELYON16 and CAMELYON17.

	In-Domain (CAM16)		Out-of-Domain (CAM17)	
	AUC	Accuracy	AUC	Accuracy
Merge <sub>add</sub> (2-feats) Retr <sub>I</sub>	0.9409±0.038	0.9000±0.049	-	-
Merge <sub>add</sub> (2-feats) Retr <sub>IO</sub>	0.9341±0.051	0.8925±0.053	0.7651±0.056	0.7714±0.030
Merge <sub>add</sub> (2-feats) Retr <sub>O</sub>	0.9414±0.046	0.8975±0.056	0.7775±0.050	0.7392±0.063
Merge <sub>add</sub> (3-feats) Retr <sub>I</sub>	0.9383±0.050	0.9175±0.051	-	-
Merge <sub>add</sub> (3-feats) Retr <sub>IO</sub>	0.9313±0.044	0.9000±0.045	0.7641±0.059	0.7553±0.043
Merge <sub>add</sub> (3-feats) Retr <sub>O</sub>	0.9391±0.051	0.9175±0.045	0.7644±0.059	0.7754±0.022
Merge <sub>convex</sub> (2-feats) Retr <sub>I</sub>	0.9451±0.036	0.8925±0.050	-	-
Merge <sub>convex</sub> (2-feats) Retr <sub>IO</sub>	0.9365±0.052	0.9200±0.050	0.7974±0.054	0.7433±0.073
Merge <sub>convex</sub> (2-feats) Retr <sub>O</sub>	0.9419±0.048	0.9175±0.051	0.7681±0.058	0.7795±0.021
Merge <sub>convex</sub> (3-feats) Retr <sub>I</sub>	0.9398±0.043	0.8975±0.052	-	-
Merge <sub>convex</sub> (3-feats) Retr <sub>IO</sub>	0.9435±0.038	0.8975±0.052	0.7652±0.052	0.7714±0.030
Merge <sub>convex</sub> (3-feats) Retr <sub>O</sub>	0.9417±0.048	0.9050±0.050	0.7690±0.056	0.7755±0.021

510 **Merge Function.** Table 6 presents the results using different merge functions. To generate the  
511 bag representations, we employ two merge functions: 1) simple addition, referred to as Merge<sub>add</sub>;  
512 2) convex combination, referred to as Merge<sub>convex</sub>. Additionally, ‘2-feats’ and ‘3-feats’ refer to  
513 bag representation that are merged with 1 nearest neighbor or 2 nearest neighbors, respectively.  
514 For convex combination, ‘2-feats’ uses coefficients of 0.6 and 0.4, while ‘3-feats’ uses coefficients  
515 of 0.6, 0.2 and 0.2, where the greatest coefficient corresponds to the original representation. This  
516 experiment is done with  $\eta = 10$ .

517 It is derived from Table 6 that using 1 nearest neighbor and convex combination presents the best  
518 performance. Using 2 nearest neighbors and addition presents the similar results.

519 **Classification Training.** Finally, we train a single logistic regression classifier using the merged  
 520 representation. The Adam optimizer is used to update the model parameters, with a L2 weight  
 521 decay of  $1e-4$  and a learning rate of  $2e-4$ . The models are trained for a minimum of 40 epochs and  
 522 up to a maximum of 200 epochs if the early stopping criterion is not met. This criterion involves  
 523 monitoring the validation loss each epoch and if it has not decreased from the previous low for over  
 524 15 consecutive epochs, early stopping is used.

## 525 D Patch-level results on tumor slides of CAMELYON16.

Table 7: Patch-level results on tumor slides of CAMELYON16.

	P-Prec.( $\uparrow$ )	FROC( $\uparrow$ )
DSMIL	0.1030	0.4443
CLAM	0.6068	0.4792
TransMIL	0.1726	0.4797
Bayes-MIL	<b>0.8107</b>	0.4919
RAM-MIL	0.6114	<b>0.5281</b>

526 The Tumor-Precision is calculated by the precision of classifying the tumor patches. The Patch-  
 527 Precision is calculated by averaging the precision of classifying both normal and tumor patches. The  
 528 Patch-FROC is defined as the average sensitivity (recall) at 6 predefined false positive rates: 1/4, 1/2,  
 529 1, 2, 4 and 8 FPs per WSI.

530 In Table 7, RAM-MIL presents the second best precision and the best FROC on the patch-level  
 531 segmentation. This indicates that using transport matrix for interpreting the patch-level classification  
 532 achieves the best overall performance in the trade-off of false positive rate and recall. By contrast,  
 533 Bayes-MIL could only obtain a high precision, which reduces the number of false alarm. However,  
 534 for the application of medical WSI, reducing false negative (better recall and FROC) is supposed to  
 535 be more important as classifying a positive instance to be negative is unacceptable in the application  
 536 of medical prognosis or diagnosis.