# Synthetic Data Reveals Generalization Gaps in Correlated Multiple Instance Learning

**Ethan Harvey**[1]

**Dennis Johan Loevlie**[1]

**Michael C. Hughes**[1]

ETHAN.HARVEY@TUFTS.EDU

DENNIS.LOEVLIE@TUFTS.EDU

MICHAEL.HUGHES@TUFTS.EDU

[1]*Department of Computer Science, Tufts University, Medford, MA, USA*

## Abstract

Multiple instance learning (MIL) is often used in medical imaging to classify high-resolution 2D images by processing patches or classify 3D volumes by processing slices. However, conventional MIL approaches treat instances separately, ignoring contextual relationships such as the appearance of nearby patches or slices that can be essential in real applications. We design a synthetic classification task where accounting for adjacent instance features is crucial for accurate prediction. We demonstrate the limitations of off-the-shelf MIL approaches by quantifying their performance compared to the optimal Bayes estimator for this task, which is available in closed-form. We empirically show that newer correlated MIL methods still struggle to generalize as well as possible when trained from scratch on tens of thousands of instances.

**Data and Code Availability:** Synthetic dataset and code available at: https://github.com/tufts-ml/correlated-mil.

**Institutional Review Board (IRB)** Our study uses synthetic data and does not require IRB approval.

## 1. Introduction

Deep neural networks pre-trained on large natural image datasets like ImageNet (Deng et al., 2009) have become a cornerstone of modern computer vision. These models are often adapted to medical imaging tasks via transfer learning, typically by fine-tuning on similar resolution 2D medical images (Mei et al., 2022). However, many medical imaging tasks involve data with different resolution or dimensionality (Quellec et al., 2017). For example, inputs may consist of high-resolution 2D images (e.g., histopathology images) or 3D image volumes (e.g., CT or MRI scans). In these scenarios, a common approach is to divide each image into smaller 2D patches or slices known as *instances*, obtain per-instance representations, and then

aggregate scores or representations across instances to make a whole-image prediction (Ilse et al., 2018; Han et al., 2020; Shao et al., 2021; Harvey et al., 2023). Building predictors that aggregate one coherent prediction from many instance representations is known as *multiple instance learning* (MIL) (Dietterich et al., 1997; Maron and Lozano-Pérez, 1997).

Although MIL offers a practical framework for handling weakly labeled data, conventional MIL approaches broadly treat instances individually and separately. This assumption ignores the spatial and contextual relationships between adjacent patches or slices – relationships that can be critical for accurate prediction in medical applications. To address this problem, recent work has proposed *correlated MIL* (Shao et al., 2021) to model dependencies between instances. Assessing the capabilities and limits of such methods remains an open question.

In this work, we take a probabilistic approach to better understand the importance of spatial and contextual relationships between adjacent instances. Our contributions are:

- We design a novel synthetic dataset to represent key challenges in MIL-for-medical imaging: (1) only some features are discriminative, (2) only a few instances in each bag signal whether it should be positive class, and (3) context from nearby instances matters, as the information from an individual instance may be statistically ambiguous.

- We derive the optimal Bayes estimator for this synthetic dataset and use it as a gold standard to quantify the limitations of MIL approaches.

- We demonstrate that even recent correlated MIL methods designed to account for context do not achieve the best possible performance on our toy task when trained from scratch, as shown in Fig. 1.
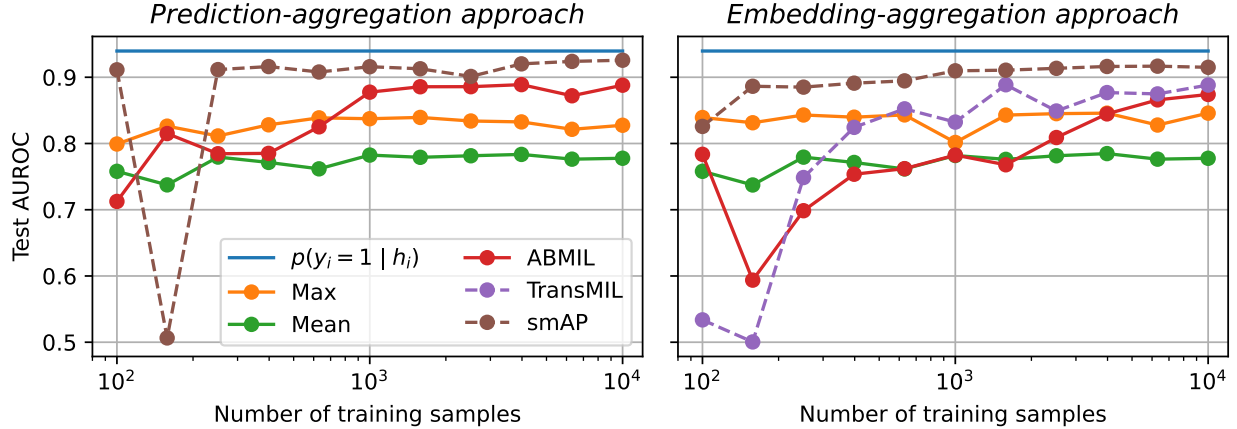
Figure 1: Test AUROC as a function of training set size $N$, for data drawn from our shifted-mean MIL binary classification task with $R=3, \Delta=2$. Conventional MIL approaches (Max, Mean, ABMIL) cannot match the Bayes estimator $p(y_i = 1 \mid h_i)$ as they do not account for dependencies between instances within a bag. Correlated MIL approaches (TransMIL (Shao et al., 2021) or smAP (Castro-Macías et al., 2024)) surprisingly also do not reach the ceiling set by the Bayes estimator even with $N = 10000$. smAP comes close (it induces similar attention for neighboring slices), but paired bootstrap testing reveals the Bayes estimator maintains a statistically significant advantage (mean of AUROC difference is 0.014, 95% CI of [0.007, 0.022] does *not* include zero or any negative values). **Takeaway: Our work reveals a need for data-efficient MIL that better accounts for context between instances.**

These contributions suggest concrete opportunities for future work to improve correlated MIL, perhaps via improved inductive biases or regularization strategies.

## 2. Related Work

**Multiple instance learning.** MIL is a branch of weakly supervised learning where a variable-sized set of instances has a single label. Early MIL approaches used simple, non-trainable operations such as max or mean pooling to aggregate instance representations (Pinheiro and Collobert, 2015; Zhu et al., 2017; Feng and Zhou, 2017). Recent work has proposed attention-based pooling (Ilse et al., 2018). Several works have extended attention-based pooling while maintaining permutation-invariance (Li et al., 2021; Lu et al., 2021; Keshvarikhojasteh et al., 2024). Correlated MIL (Shao et al., 2021) extends traditional MIL by modeling relationships between instances within a bag, allowing the pooling operation to capture morphological and spatial information rather than treating instances as IID. Castro-Macías et al. (2024) proposed a smoothing operator to introduce local dependencies among neighbors. Shao et al. (2025) showed that transfer learning with MIL approaches improves generalization.

Most similar in spirit to our work are the *algorithmic unit tests* for MIL proposed by Raff and Holt (2023). They suggest 3 synthetic classification tasks designed to reveal learning that violates key MIL assumptions, such as a bag is positive if and only if one or more instances have a positive label. Our toy data task instead focuses on cross-instance dependency.

**Adjacent context in deep learning.** Several prior works have introduced architectural modifications to better capture dependencies between adjacent patches in 2D images or slices in 3D images. Shifted windows allow for cross-window connections in vision transformers (ViT) (Liu et al., 2021). Weight inflation transfers pre-trained weights from lower- to higher-dimensional model (e.g., 2D to 3D CNNs) (Carreira and Zisserman, 2017; Zhang et al., 2022).

## 3. Background

### 3.1. Multiple instance learning

To train MIL models, the dataset $\mathcal{D} = \{(x_{i,1:S_i}, y_i)\}_{i=1}^N$ consists of the $N$ labeled bags. Each bag is a set of $S_i$ independent instance feature vectors $\{x_{i,1}, \ldots, x_{i,S_i}\}$ with a single label $y_i$. Among deep MIL pipelines, there are two broad paradigms: embedding-aggregation and prediction-aggregation. Both approaches have neural architectures consisting of three parts: an encoder, a pooling operation, and classifier. They differ in the ordering of these components.
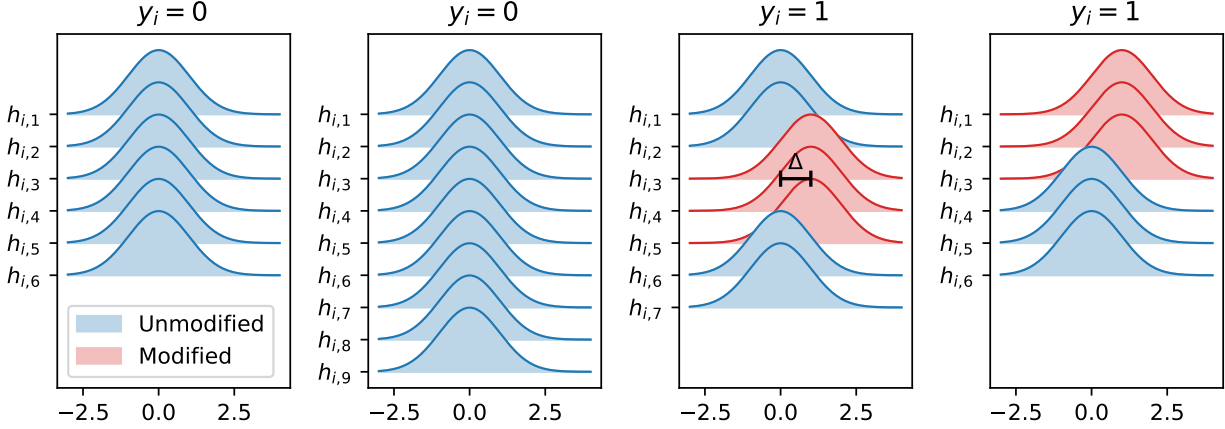
Figure 2: Example data-generating distributions for a discriminative feature for negative ($y_i{=}0$) and positive ($y_i{=}1$) "bags" of $S_i$ instances drawn from our shifted-mean MIL toy data. Setting $R = 3$ means context around signaling instances (in red) can help. Not shown: each bag has $M{=}768$ total features, only $K{=}1$ of which are discriminative.

In the *embedding-aggregation* approach, the order is encode, pool, then classify. First, each instance $x_{i,j} \in \mathbb{R}^{C_{\text{in}} \times W \times H}$ is encoded into a representation vector $h_{i,j} = f(x_{i,j}) \in \mathbb{R}^M$. Second, a pooling operation $\sigma$ (e.g., max, mean, or attention-based pooling) aggregates all $S_i$ instance representations $\{h_{i,1}, \ldots, h_{i,S_i}\}$ into a single representation vector $z_i = \sigma(h_{i,1:S_i}) \in \mathbb{R}^M$. Usually, this pooling is *permutation-invariant*. Finally, the bag level representation vector $z_i$ is classified into a predicted probability vector over $C$ classes, $g(z_i) \in \Delta^C \subset \mathbb{R}^C$. We can denote the ultimate prediction as $\hat{y}_i = g(\sigma(f(x_{i,1:S_i}))$. In this notation, applying $f$ to a set yields another set containing a mapping of each instance.

In the *prediction-aggregation* approach, the ordering of $g$ and $\sigma$ is swapped. A separate prediction score (e.g., a logit or probability vector) is produced for each of the $S_i$ instances separately, and then pooling determines the final prediction, $\hat{y}_i = \sigma(g(f(x_{i,1:S_i}))))$.

Ultimately, in either approach, model parameters for all three parts (encoder, pooling, and classifier) are trained to minimize binary or multi-class cross entropy averaged across all data: $\frac{1}{N} \sum_{i=1}^{N} \ell^{\text{CE}}(y_i, \hat{y}_i)$, where $\hat{y}_i$ is a function of input features and parameters.

### 3.2. Pooling Methods

The design of the pooling layer $\sigma$, which aggregates across instances, is generally most important for understanding how spatial context is incorporated. We describe several architectures below. We focus on embedding-aggregation for concreteness, but translation to prediciton-aggregation is possible. Here, we take as input a set of embeddings $h_i := h_{i,1:S_i}$ for bag $i$, where each instance $j$ in the bag has feature vector $h_{i,j} \in \mathbb{R}^M$.

**Max and mean.** Two simple poolings find the maximum or mean **element-wise** for $M$-dim. vectors:

$$z_i = \max_{j=1,\ldots,S_i} h_{ij}, \quad \text{or} \quad z_i = \operatorname*{mean}_{j=1,\ldots,S_i} h_{ij}. \quad (1)$$

**Attention-based pooling.** Attention-based pooling (ABMIL, Ilse et al. 2018) assigns an attention weight $a_{ij}$ to each instance, then forms bag-level embedding vector $z_i$ via a weighted average:

$$z_i = \sum_{j=1}^{S_i} a_{ij} h_{ij}, \quad a_{ij} = \frac{\exp\left(u^\top \tanh\left(U h_{ij}\right)\right)}{\sum_{k=1}^{S_i} \exp\left(u^\top \tanh\left(U h_{ik}\right)\right)},$$

where the weights $a_{ij}$ are non-negative and sum to one: $\sum_j a_{ij} = 1$; $a_{ij}{\geq}0$ for all $j$. Here, vector $u \in \mathbb{R}^L$ and matrix $U \in \mathbb{R}^{L \times M}$ are trainable parameters.

**Smooth attention pooling.** Smooth attention pooling (smAP) (Castro-Macías et al., 2024) uses a smoothing operation to add local interactions between instance embeddings. The smoothed embeddings $g_i \in \mathbb{R}^{S_i \times M}$ are obtained by minimizing the objective

$$\mathtt{Sm}(h_i) = \operatorname*{argmin}_{g_i} \mathcal{E}(g_i, h_i), \quad (2)$$

$$\mathcal{E}(g_i, h_i) = \alpha \mathcal{E}_D(g_i) + (1-\alpha)\|h_i - g_i\|_F^2, \quad (3)$$

where $\alpha \in [0, 1)$ controls the amount of smoothness, $\|\cdot\|_F$ denotes the Frobenius norm, and

$$\mathcal{E}_D(g_i) = \frac{1}{2} \sum_{j=1}^{S_i} \sum_{k=1}^{S_i} A_{i,j,k} \|g_{i,j} - g_{i,k}\|_2^2. \quad (4)$$

3

Here, $A_i \in \mathbb{R}^{S_i \times S_i}$ is an adjacency matrix defining local relationships between instances and $\|\cdot\|_2^2$ denotes the squared Euclidean norm aka "sum of squares".

**Correlated MIL pooling.** Shao et al. (2021)'s transformer-based correlated MIL (TransMIL) allows instance interactions to inform pooling. First, TransMIL uses convolutions over instances in a pyramidal position encoding generator to model dependencies. Second, interactions between *all pairs* of instances are captured via multi-head self-attention. For layer $\ell$ and head $h$, there's a $S_i{+}1 \times S_i{+}1$ attention matrix, where rows sum to one and weight $j, k$ is:

$$a_{i,j,k}^{(\ell,h)} \propto \exp\left( \left(q_{i,j}^{(\ell,h)}\right)^{\top} k_{i,k}^{(\ell,h)} / \sqrt{d} \right). \qquad (5)$$

Here, each instance $j$ has $L$-dim. embeddings for query $q_{i,j}^{(\ell,h)}{=}W_Q^{(\ell,h)}h_{i,j}^{\ell-1}$, key $k_{i,j}^{(\ell,h)}{=}W_K^{(\ell,h)}h_{i,j}^{\ell-1}$, and value $v_{i,j}^{(\ell,h)}{=}W_V^{(\ell,h)}h_{i,j}^{\ell-1}$. Propagating embeddings via attention-weighted value averages over several layers allows instance features to interact flexibly to inform the ultimate bag-level embedding.

## 4. New Toy Data: Shifted-Mean MIL

We propose a data-generating process designed to mimic several key challenges in real-world multiple-instance medical imaging tasks:

- For each positively-labeled bag, only a few instances are relevant ($R$ of $S_i$) and they are adjacent in a known 1D listing of all $S_i$ instances.

- Across the whole dataset, only a few features of many are discriminative ($K$ of $M$).

- Context matters. Adjacent instances together provide stronger statistical signal than any one relevant instance's discriminative feature value alone.

The generative process for bag $i$ first draws the bag's binary label and the number of instances in the bag

$$y_i \sim \mathrm{Bern}(q_+), \quad S_i \sim \mathrm{Unif}(\{S_{\mathrm{low}}, \dots, S_{\mathrm{high}}\}). \quad (6)$$

Next, for negative bags we sample all features $m$ for all instances $j$ independently from a common Gaussian:

$$h_{i,j,m} \mid y_i{=}0 \sim \mathcal{N}(\mu, \sigma^2). \qquad (7)$$

For positive bags, most instances and features are sampled from this same Gaussian. However, for the $K$ discriminative features, we select $R$ adjacent instances (using $u_i$ to denote the starting index) and sample

these from a Gaussian with *shifted mean*:

$$u_i \mid y_i{=}1 \sim \mathrm{Unif}(\{1, \dots, S_i - R + 1\}), \qquad (8)$$

$$h_{ijm} \mid u_i, y_i{=}1 \sim \begin{cases} \mathcal{N}(\mu + \Delta, \sigma^2), \text{ if } j \in [u_i, u_i{+}R{-}1] \\ \qquad\qquad\qquad \text{ and } m \text{ is discrim.} \\ \mathcal{N}(\mu, \sigma^2), \qquad \text{otherwise.} \end{cases}$$

Here $\Delta > 0$ indicates the magnitude of shift for discriminative features. Setting $R > 1$ indicates that context helps. Given a fixed $\mu$, bags drawn from this process are more challenging to classify (even with knowledge of the true process) when $\Delta$ is smaller, $R$ is smaller, $K/M$ is smaller, and $\sigma$ is larger.

This data-generating process is illustrated in Fig. 2, depicting only one feature that is discriminative. In each positive bag, a different contiguous block of $R{=}3$ instances draw from the shifted-mean Gaussian. If future work wanted to model correlations between features within an instance, the sampling of vector $h_{ij}$ in Eq. (8) could be modified to draw from a multivariate Gaussian with a non-diagonal covariance matrix.

## 5. Bayes estimator for Toy Data

Given a data-generating process, a *Bayes estimator* is a decision rule that minimizes the posterior expected loss with respect to the data-generating distribution (DeGroot, 1970; Murphy, 2022). It is an oracle upper bound on performance. By comparing conventional MIL methods to the Bayes estimator for our synthetic dataset, we can quantify how close they come to the best possible performance.

Given a new bag $h_i$ containing $S_i$ instances and assuming our data-generating process defined above, a Bayes estimator for class label probability is:

$$p(y_i = 1 \mid h_i) = \frac{p(h_i \mid y_i = 1)p(y_i = 1)}{p(h_i)}. \qquad (9)$$

Each term on the right-hand side can be computed in closed-form. The marginal likelihood of $h_i$ in the denominator is given by the sum rule:

$$p(h_i) = p(h_i \mid y_i{=}0)p(y_i{=}0) + p(h_i \mid y_i{=}1)p(y_i{=}1).$$

where we recall $p(y_i{=}1)$ is $q_+$. The class-conditional likelihood for the negative class factors over instances:

$$p(h_i \mid y_i = 0) = \prod_{j=1}^{S_i} \prod_{k=1}^{M} \mathrm{NormPDF}(h_{i,j,k} \mid \mu, \sigma^2).$$

For the positive class, a latent segment of $R$ consecutive slices is modified. The class-conditional likelihood marginalizes over its possible starting indices

$$p(h_i \mid y_i{=}1) = \sum_{u=1}^{S_i-R+1} \left[ p(u \mid y_i{=}1) \prod_{j=1}^{S_i} \prod_{k=1}^{M} p(h_{i,j,k} \mid u, y_i{=}1) \right].$$

with Eq. (8) providing the necessary PDF values to evaluate the right hand side.

## 6. Experimental Results

**Setup.** In our experiments, we sample bag labels uniformly ($q_+ = 0.5$) and the number of instances per bag uniformly between $S_{\text{low}} = 15$ and $S_{\text{high}} = 45$. We fix $M=768$ features to match the size of ViT-B/16 embeddings (Dosovitskiy et al., 2021). Only a single feature ($K=1$) is discriminative. We set $R = 3$ so context matters. For main results in Fig. 1, we fix $\Delta=2$. For results varying $\Delta$, see the Appendix.

For each train set size $N$, we draw $N$ bags from our data-generating process. This data is then split into training and validation sets using a 4:1 ratio.

For each MIL, we try both *prediction-aggregation* and *embedding-aggregation* approaches when possible. We train models for 1000 epochs. We use validation set AUROC for early stopping and to select learning rate from {0.1, 0.01, 0.001, 0.0001} and weight decay from {1.0, 0.1, 0.01, 0.001, 0.0001, 1e-5, 1e-6, 0}.

After selecting the hyperparameters, we report AUROC performance on a common test set of 1000 bags.

Examining the results of our synthetic dataset experiments, key findings are:

- **Conventional MIL trained from scratch cannot match the Bayes estimator when $R = 3$.** Even given $N=10^4$ bags, conventional deep MIL approaches (Max, Mean, ABMIL (Ilse et al., 2018)) deliver test AUROC at least 0.04 below the Bayes estimator in Fig. 1. Trends over $N$ do not suggest this gap will close with more data.

- **TransMIL cannot match the Bayes estimator when $R=3$.** More surprisingly, the same trend is observed with TransMIL, which purports to handle context between instances. As a sanity check, we were able to *hand-craft* neural net parameters for TransMIL that match the Bayes estimator (see Fig. 5). However, training from scratch appears consistently suboptimal on this toy data.

- **smAP cannot match the Bayes estimator when $R=3$.** We conducted an extensive hyperparameter search for smAP (the most comprehensive to our knowledge). At the largest training set size in Fig. 1, the *prediction-aggregation* variant achieved test AUROC that closely approached the Bayes estimator. To determine whether this remaining gap was statistically significant, we performed bootstrap analysis (see App. C), which confirmed that smAP had not yet

reached Bayes-optimal performance. Since smAP smooths instance features with their neighbors via a chain-graph adjacency matrix (appropriate for MRI/CT sequences), it is well-suited to the $R=3$ setting. For $R>3$, we expect the performance gap to increase as the required context exceeds smAP's local aggregation. Adjustments to the graph structure could improve performance for larger $R$, we leave this for future work.

As a sanity check that instance context is what is important, we verified that for the context-free $R = 1$ case, max pooling, ABMIL, and TransMIL all can match the Bayes estimator with handcrafted parameters (see App. A). In App. A, we further show how adding context to ABMIL via convolutions over instances can help when $R=3$.

## 7. Conclusion

We designed an easy-to-implement synthetic dataset designed to mimic key challenges in MIL for medical imaging, especially the need for context from nearby instances. We then demonstrated the limited ability of conventional MIL on such data, by quantifying the performance gap between the optimal Bayes estimator. More recent correlated MIL methods like TransMIL are still notably worse than optimal even with a labeled dataset of 10000 bags. Very recent methods like smAP come closer, but still fall short.

**Outlook.** We hope our synthetic dataset enables the development of correlated MIL methods that can be trained effectively with limited labeled data.

# References

Joao Carreira and Andrew Zisserman. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

Francisco M. Castro-Macías, Pablo Morales-Álvarez, Yunan Wu, Rafael Molina, and Aggelos K. Katsaggelos. Sm: enhanced localization in Multiple Instance Learning for medical imaging classification. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.

Morris H. DeGroot. *Optimal Statistical Decisions*. McGraw-Hill, 1970.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.

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.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations (ICLR)*, 2021.

Ji Feng and Zhi-Hua Zhou. Deep MIML Network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2017.

Zhongyi Han, Benzheng Wei, Yanfei Hong, Tianyang Li, Jinyu Cong, Xue Zhu, Haifeng Wei, and Wei Zhang. Accurate screening of covid-19 using attention-based deep 3d multiple instance learning. *IEEE transactions on medical imaging*, 39(8): 2584–2594, 2020.

Ethan Harvey, Wansu Chen, David M. Kent, and Michael C. Hughes. A Probabilistic Method to Predict Classifier Accuracy on Larger Datasets given Small Pilot Data. In *Machine Learning for Health (ML4H)*, 2023.

Maximilian Ilse, Jakub Tomczak, and Max Welling. Attention-based Deep Multiple Instance Learning. In *International Conference on Machine Learning (ICML)*, 2018.

Hassan Keshvarikhojasteh, Josien P. W. Pluim, and Mitko Veta. Multi-head attention-based deep multiple instance learning. In *Proceedings of the MICCAI Workshop on Computational Pathology*, 2024.

Bin Li, Yin Li, and Kevin W Eliceiri. Dual-stream Multiple Instance Learning Network for Whole Slide Image Classification with Self-supervised Contrastive Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.

Ming Y Lu, Drew FK Williamson, Tiffany Y Chen, Richard J Chen, Matteo Barbieri, and Faisal Mahmood. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature Biomedical Engineering*, 5(6):555–570, 2021.

Oded Maron and Tomás Lozano-Pérez. A framework for multiple-instance learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 1997.

Xueyan Mei, Zelong Liu, Philip M Robson, Brett Marinelli, Mingqian Huang, Amish Doshi, Adam Jacobi, Chendi Cao, Katherine E Link, Thomas Yang, et al. Radimagenet: An open radiologic deep learning research dataset for effective transfer learning. *Radiology: Artificial Intelligence*, 4(5): e210315, 2022.

Kevin S. Murphy. *Probabilistic Machine Learning: An Introduction*, chapter 5.1 Bayesian decision theory. MIT Press, 2022.

Pedro O Pinheiro and Ronan Collobert. From Image-Level to Pixel-Level Labeling With Convolutional Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.

Gwenolé Quellec, Guy Cazuguel, Béatrice Cochener, and Mathieu Lamard. Multiple-instance learning for medical image and video analysis. *IEEE Reviews in Biomedical Engineering*, 10:213–234, 2017.

Edward Raff and James Holt. Reproducibility in multiple instance learning: A case for algorithmic unit tests. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

Daniel Shao, Richard J Chen, Andrew H Song, Joel Runevic, Ming Y Lu, Tong Ding, and Faisal Mahmood. Do Multiple Instance Learning Models Transfer? In *International Conference on Machine Learning (ICML)*, 2025.

Zhuchen Shao, Hao Bian, Yang Chen, Yifeng Wang, Jian Zhang, Xiangyang Ji, et al. TransMIL: Transformer based Correlated Multiple Instance Learning for Whole Slide Image Classification. *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.

Yuhui Zhang, Shih-Cheng Huang, Zhengping Zhou, Matthew P Lungren, and Serena Yeung. Adapting Pre-trained Vision Transformers from 2D to 3D through Weight Inflation Improves Medical Image Segmentation. In *Machine Learning for Health (ML4H)*, 2022.

Wentao Zhu, Qi Lou, Yeeleng Scott Vang, and Xiaohui Xie. Deep Multi-instance Networks with Sparse Label Assignment for Whole Mammogram Classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2017.

# Appendix A. Handcrafted Parameters

We set classifier weights corresponding to each discriminative feature to one, all other weights to zero, and the bias parameter to $-\frac{\Delta}{2}$. For attention pooling, we use the linear part of the tanh function to create attention weights proportional to each feature's linear score. For instance convolutions, we set the center $R$ weights to one and all other weights and biases to zero to sum up each feature's linear score over $R$ instances.
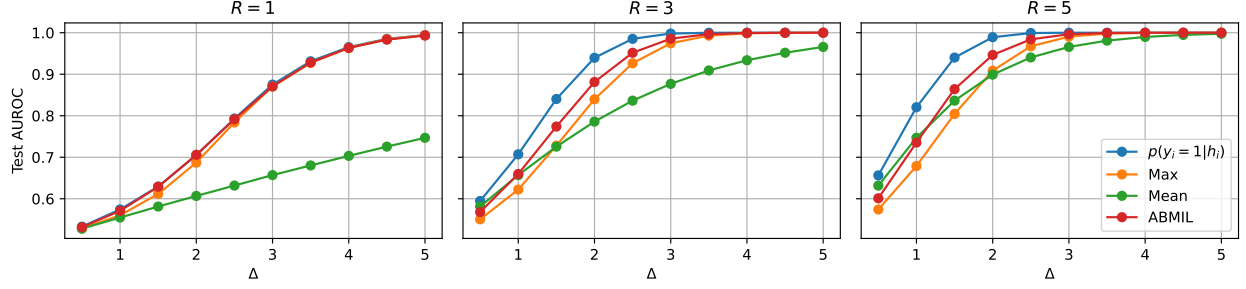


Figure 3: Handcrafted parameters for classification (*prediction-aggregation* approach).



Figure 4: Handcrafted parameters for classification (*prediction-aggregation* approach).
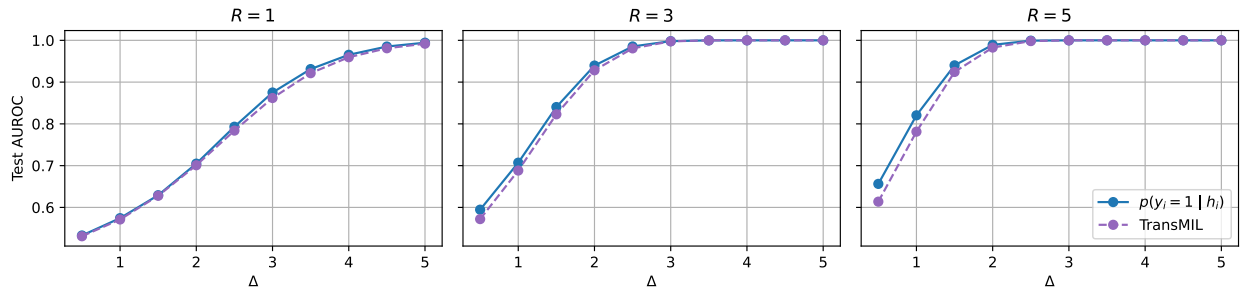


Figure 5: Handcrafted parameters for classification (*embedding-aggregation* approach).

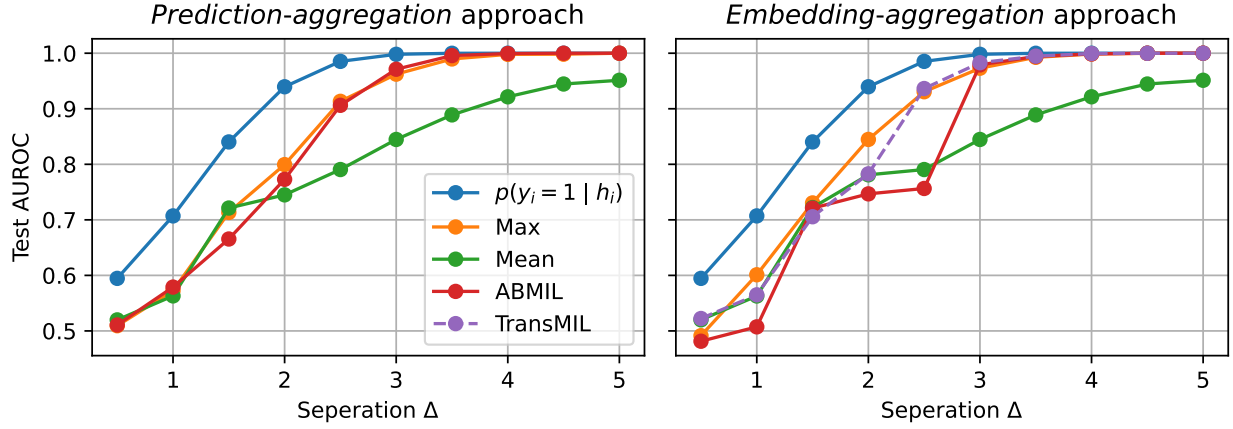## Appendix B. Varying Separation



Figure 6: Test AUROC as a function of separation $\Delta$, for data drawn from our shifted-mean MIL binary classification task with $R = 3,400$ training samples.

## Appendix C. Bootstrap Analysis

We use bootstrapping to access the statistical significance of the AUROC difference between the Bayes estimator and smAP for the *prediction-aggregation* approach trained on 10,000 samples We report the mean and 95% confidence interval over 500 subsamples of the test set.
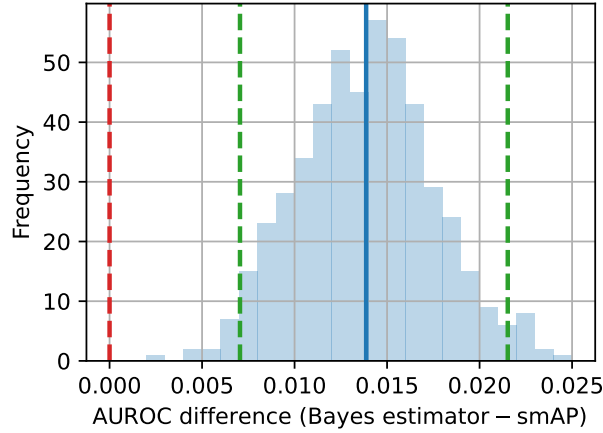


Figure 7: Bootstrap analysis comparing the AUROC difference between the Bayes estimator and smAP for the *prediction-aggregation* approach train on 10,000 samples.

## Appendix D. Distribution over Number of Slices

We have a uniform distribution over $S_i$ which can be factored out of the numerator and denominator in the posterior and cancels out

$$p(y_i = 1 \mid h_i) = \frac{p(h_i \mid y_i = 1, S_i)p(y_i = 1)\cancel{p(S_i)}}{p(h_i \mid y_i = 0, S_i)p(y_i = 0)\cancel{p(S_i)} + p(h_i \mid y_i = 1, S_i)p(y_i = 1)\cancel{p(S_i)}}. \tag{10}$$