

# BAYES-MIL: A NEW PROBABILISTIC PERSPECTIVE ON ATTENTION-BASED MULTIPLE INSTANCE LEARNING FOR WHOLE SLIDE IMAGES

Yufei Cui<sup>1</sup>, Ziquan Liu<sup>2</sup>, Xiangyu Liu<sup>3</sup>,  
 Xue Liu<sup>1</sup>, Cong Wang<sup>2</sup>, Tei-Wei Kuo<sup>4,5</sup>, Chun Jason Xue<sup>2</sup>, Antoni B. Chan<sup>2</sup>  
<sup>1</sup>McGill University <sup>2</sup>City University of Hong Kong <sup>3</sup>Bingli AI Research  
<sup>4</sup>National Taiwan University <sup>5</sup>Mohamed bin Zayed University of Artificial Intelligence

## ABSTRACT

Multiple instance learning (MIL) is a popular weakly-supervised learning model on the whole slide image (WSI) for AI-assisted pathology diagnosis. The recent advance in attention-based MIL allows the model to find its region-of-interest (ROI) for interpretation by learning the attention weights for image patches of WSI slides. However, we empirically find that the interpretability of some related methods is either untrustworthy as the principle of MIL is violated or unsatisfactory as the high-attention regions are not consistent with experts’ annotations. In this paper, we propose Bayes-MIL to address the problem from a probabilistic perspective. The induced patch-level uncertainty is proposed as a new measure of MIL interpretability, which outperforms previous methods in matching doctors annotations. We design a slide-dependent patch regularizer (SDPR) for the attention, imposing constraints derived from the MIL assumption, on the attention distribution. SDPR explicitly constrains the model to generate correct attention values. The spatial information is further encoded by an approximate convolutional conditional random field (CRF), for better interpretability. Experimental results show Bayes-MIL outperforms the related methods in patch-level and slide-level metrics and provides much better interpretable ROI on several large-scale WSI datasets.

## 1 INTRODUCTION

In real-world applications of deep learning, data like images or texts are often associated with insufficient labels, due to the expensive annotation cost. For example, the whole slide images (WSI) for medical diagnosis have about  $10^5 \times 10^5$  pixels per image, but are tagged with single categorical labels (Zhang et al., 2019; Campanella et al., 2019). Weakly-supervised learning methods are designed for learning representations and making decision in these cases. Multiple instance learning (MIL) is a popular weakly-supervised learning model for the application of WSI recognition (Ilse et al., 2018; Lu et al., 2021). Concretely, a large WSI slide is sliced into a bag of image patches (instances) with a moderate size.<sup>1</sup> MIL builds an end-to-end parametric model that aggregates the learned features from instances and only learns from bag-level labels. The rule of aggregation is implementing *the key principle of MIL*: for binary classification, a bag is negative when all instances are negative, and a bag is positive when there is one or more positive instance (Ilse et al., 2018).

Recent advances study the *attention-based MIL* for re-weighting the instances for better performance. This attention mechanism for MIL is extensively explored and used as a measure of interpretability in various downstream tasks for medical diagnosis, like prostatic cancer (Zhang et al., 2021), breast cancer (Naik et al., 2020), etc. Specifically, the high attention weights are used to indicate that its associated instances are positive instances, e.g, the cancerous image patches. However, this rule is not formally justified and it is not clear whether the negative instances (i.e., benign) would be assigned a high attention value or the other way around. We first analyze the convergence of attention and provide validity of this rule under binary labels. Based on this rule, we conduct an empirical study on a large scale WSI dataset for how the attention mechanism in the related MIL methods performs. The study shows two clear flaws of the related methods:

<sup>1</sup>We define the following interchangeable terms for simplicity: “bag” and “slide”; “instance” and “patch”.

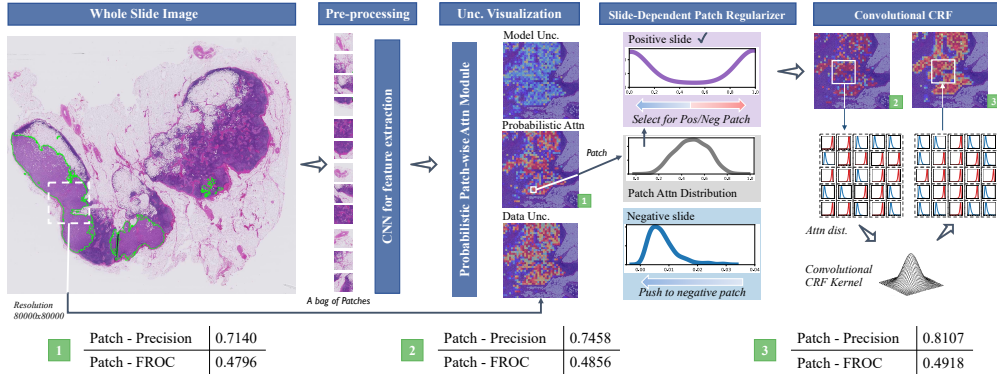


Figure 1: The overview of our Bayes-MIL framework and zoom-in views for clear visualization of interpretability. (1) The basic Bayes-MIL improves patch-level localization performance (Sec. 3.1). (2) The slide-dependent patch regularizer makes attention densely concentrated on the positive area, improving its interpretability (Sec. 3.2). (3) The convolutional CRF improves the localization by smoothing the uncertainty over different patches (Sec. 3.3). (bottom) The ablation results on a few metrics show the improvement of interpretability. The full ablation results are in Sec. 5.

- The interpretability for negative bags is untrustworthy because some methods violate the *key principle of MIL* by placing high attention values on negative bags, thus indicating positive instances.
- The interpretability for positive bags is unsatisfactory because the high attention values could not well match experts’ annotations of positive instances.

In this paper, we address the problems from a probabilistic perspective. First, a basic framework of Bayesian MIL (Bayes-MIL) is proposed, for inducing uncertainty over the attention weights. The uncertainty is potentially an accurate measure for guessing whether the instances are positive or negative, as a replacement of attention. Second, a regularizer is designed by deduction from the MIL principle and implemented via the variational inference framework, which sets specific constraints for the attention distributions of positive and negative bags. Third, to encode the spatial information of instances for medical imaging application, we propose an approximate operation to the convolutional conditional random field, which benefits the localization of the region of interest (ROI). The final classifier is modeled in a Bayesian way, in order to provide calibrated uncertainty of the bag-level prediction. The overview of our proposed method is shown in Fig. 1. The contributions of this paper are listed as follows:

- We analyze the attention-based MIL on the interpretability-critic medical application and point out the flaws by directly using attention for interpretation.
- To address these problems, we propose the first Bayesian MIL for WSI with 3 key components: a probabilistic instance-wise attention module for uncertainty visualization, the slide-dependent patch regularizer for learning the correct attention distribution, and an approximate convolutional conditional random field for encoding spatial information. Our model provides well-calibrated uncertainties, which is crucial for safety in medical applications.
- The evaluation on large-scale MIL datasets shows Bayes-MIL outperforms the related methods in instance-level interpretation and bag-level prediction under various evaluation metrics. The visualized distribution of data uncertainty shows a strong correlation of the designed regularizer, which validates the soundness of regularizer and explains why uncertainty is useful in MIL interpretation.

## 2 FORMULATION AND ANALYSIS OF MULTIPLE INSTANCE LEARNING

**Multiple instance learning formulation** We follow the standard formulation of Attention-based Multiple Instance Learning (MIL) (Ilse et al., 2018; Lu et al., 2021). In MIL, the input is a bag of instances,  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_K\}$ ,  $\mathbf{x}_k \in \mathbb{R}^D$ .  $K$  is the number of instances, which varies for different bags. There is a bag-level label  $Y$ . We further assume the instances also have corresponding instance-level labels  $\{y_1, \dots, y_K\}$ , which are *unknown* during training. There are  $N$  such bag-label pairs constituting the dataset  $\mathcal{D} = \{\mathbf{X}_n, Y_n\}_{n=1}^N$ . The objective of MIL is to learn an optimal function for predicting the bag-level label with the bag of instances as input. To this end, the MIL model should be able to aggregate the information of instances  $\{\mathbf{x}_k\}_{k=1}^K$  to make the final decision. A well-adopted aggregation method is the embedding-based approach which maps  $\mathbf{X}$  to a bag-level representation  $\mathbf{z} \in \mathbb{R}^D$  and use  $\mathbf{z}$  to predict  $Y$ . Ilse et al. (2018) extends the embedding-based

aggregation approach by leveraging the attention mechanism, namely attention-based deep MIL (ABMIL). First, a transformation  $g(\cdot)$  computes a low-dimensional embedding  $\mathbf{h}_k = g(\mathbf{x}_k) \in \mathbb{R}^D$  for each instance  $\mathbf{x}_k$ . The attention module aggregates the set of embeddings  $\{\mathbf{h}_k\}_{k=1}^K$  into a bag level embedding  $\mathbf{z}$ ,  $\mathbf{z} = \sum_{k=1}^K a_k \mathbf{h}_k$ , where the attention for the  $k$ -th instance is computed via a softmax function,

$$a_k = f_\pi(\mathbf{H})_k = \frac{\exp\{\mathbf{m}^T(\tanh(\mathbf{V}_1^T \mathbf{h}_k) \odot \text{sigmoid}(\mathbf{V}_2^T \mathbf{h}_k))\}}{\sum_{j=1}^K \exp\{\mathbf{m}^T(\tanh(\mathbf{V}_1^T \mathbf{h}_j) \odot \text{sigmoid}(\mathbf{V}_2^T \mathbf{h}_j))\}}. \quad (1)$$

where the attention function  $f_\pi : \mathbb{R}^{D \times K} \rightarrow \mathbb{R}^K$  with parameters  $\pi = \{\mathbf{m}, \mathbf{V}_1, \mathbf{V}_2\}$ , and  $f_\pi(\mathbf{H})_k$  denotes the  $k$ -th output of  $f$ .  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_K] \in \mathbb{R}^{D \times K}$  is the matrix of embeddings. The bag embedding  $\mathbf{z}$  is then mapped to the logits  $\mathbf{u}$  with a feed forward layer with parameter  $\mathbf{W}$  for the bag-level classification,  $\mathbf{u} = \mathbf{W}^T \mathbf{z}$ .

**Multiple instance learning for medical imaging** WSI is a type of high-dimensional image data format (up to  $10^5 \times 10^5$  pixels per image) widely adopted in the medical area. Due to the scarcity of experts and the high annotation cost, only class labels (e.g., diagnosis results) are available for most WSIs. The high resolution of data and lack of precise annotations raise challenges for machine learning-assisted classification. To fit the data in modern computation hardware, the high-resolution image slides (*slide*) are partitioned into image patches (*patch*) before further processing. MIL fits the classification task of WSI by corresponding the slides to bags, and patches to instances. The transformation  $g(\cdot)$  is the feature extractor from a pre-trained convolutional neural network. The MIL model only predicts the slide-level label, e.g., whether a slide is cancerous or not. However, the interpretation of the patches is crucial, e.g., which patches indicate the cancer, as users always check the interpretation before trusting the prediction. ABMIL uses the attention weights to tell which patches the MIL model focuses on, and several works DSMIL (Li et al., 2021), CLAM (Lu et al., 2021), TransMIL (Shao et al., 2021) study the variants of attention for better interpretability.

Ilse et al. (2018) suggested that, with binary classification label  $Y \in \{0, 1\}$ , the high attention weights in ABMIL could locate the positive area ( $y_k = 1$ ) in an ideal case. The followup works from ABMIL use the high attention weights to indicate the positive patches in a black-box manner. However, there is no formal justification for this claim, and it is still unclear whether a high attention weight could also be assigned to a negative area during training. Therefore, we analyze the general attention-based MIL framework to provide this justification.

**The convergence of attention** Assume that the  $j$ th input patches in  $i$ th slide is  $\mathbf{h}_{i,j}$ , the classifier weight is  $\mathbf{w} \in \mathbb{R}^D$ , the attention variable is  $\mathbf{a} \in (0, 1)^{K \times 1}$ . The output of the network  $\hat{y}_i$  and loss function  $\mathcal{L}_i$  for the  $i$ -th slide are

$$\hat{y}_i = \sigma(\mathbf{w}^T \mathbf{H}_i \mathbf{a} + b), \quad \mathcal{L}_i = -Y_i \log(\hat{Y}_i) - (1 - Y_i) \log(1 - \hat{Y}_i). \quad (2)$$

The  $\mathbf{H}_i$  is all patches in  $i$ th slide, where the  $l_2$  norm of all  $\mathbf{h}_{i,j}$  is upper bounded by 1, and  $\sigma(x) = 1/(1 + \exp(-x))$  is the sigmoid function. Assume the positive and negative patches are linear separable and there is an optimal  $\mathbf{w}^*$ , where  $\|\mathbf{w}^*\| = 1$  and the margin is  $\gamma = \min_{i,j} |\mathbf{w}^{*T} \mathbf{h}_{i,j}|$ . If a slide  $\mathbf{H}_i$  is negative, then all patches in  $\mathbf{H}_i$  are negative, i.e.,  $\mathbf{w}^{*T} \mathbf{h}_{i,j} < 0, \forall j$  if  $Y_i = 0$ . If a slide  $\mathbf{H}_i$  is positive, then the first  $K_p$  patches are positive and the last  $K_n$  patches are negative, i.e.,  $\mathbf{w}^{*T} [\mathbf{h}_{i,1}, \dots, \mathbf{h}_{i,K_p}] > \mathbf{0}$  and  $\mathbf{w}^{*T} [\mathbf{h}_{i,K_p+1}, \dots, \mathbf{h}_{i,K_p+K_n}] < \mathbf{0}$  if  $Y_i = 1$ . There is an optimal  $\mathbf{a}^*$  so that the first  $K_p$  dimension is  $(1 - \epsilon)/K_p$  and the last dimension is  $\epsilon$ , where  $\epsilon$  is an infinitesimal, e.g.,  $1e-5$ . For the initialization of  $\mathbf{a}$ , we assume  $\mathbf{w}^{*T} \mathbf{H}_i \mathbf{a} > \zeta > 0, \forall i$ . The  $\mathbf{a}$  is output of softmax function,  $\mathbf{a} = s(\mathbf{u})$ , where  $a_j = \frac{\exp(u_j)}{\sum_{k=1}^K \exp(u_k)}$ .

**Lemma 1.** *If we train the  $\mathbf{w}$  and  $\mathbf{a}$  as described in Appendix 1, the  $\mathbf{w}$  converges to the optimal  $\mathbf{w}^*$  in at most  $4/\max(\gamma, \zeta)^2$  steps and  $\mathbf{a}$  converges to the desired  $\mathbf{a}^*$ , where the first  $K_p$  elements are large and the last  $K_n$  elements are small.*

Lemma 1 indicates that MIL is guaranteed to converge to the desired attention variable under the ideal condition. Note that in reality, the position of positive patches is not fixed so we use a parameterized function to make the attention variable  $a_j$  depend on the patch feature  $\mathbf{h}_{i,j}$ . Lemma 1 provides a guarantee for the validity of visualization and the design of our framework. However, in reality, the convergence of weight and attention variables depends on their initialization. Thus, existing methods may not have a good match between high-attention patches and ground-truth positive patches. The following empirical study is performed to illustrate this.

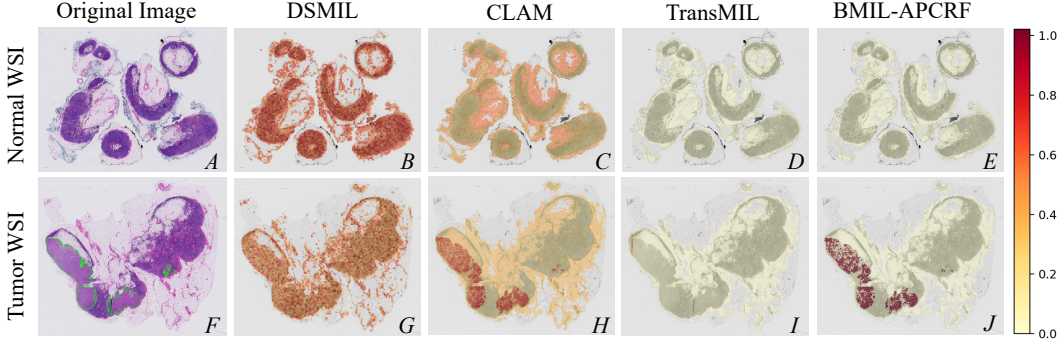


Figure 2: The visualization of normal and tumor slides and the ROIs provided by different models. The patch-level annotations for the tumor image are shown in green color in (F). The attention values  $\mathbf{a}$  are normalized to the same range by  $\frac{\mathbf{a} - \min_{\mathbf{a}}}{\max_{\mathbf{a}} - \min_{\mathbf{a}}}$ . The  $\min_{\mathbf{a}}$  and  $\max_{\mathbf{a}}$  are the same for all methods for better visualization.

**Empirical study of interpretability** We use the attention  $\{a_k\}_{k=1}^K$  as the patch-level prediction confidence and compare with the patch-level ground-truth  $\{y_k\}_{k=1}^K$ , using the related MIL approaches. One key discovery from the related methods is that *high attention values are still generated, indicating positive patches, for negative slides*, as shown in Fig. 2B 2C 2D. Another discovery is that the high attention values are not concentrated on the doctor’s annotation well during inference time, as shown in Fig. 2G 2H 2I. (see numerical results in Tab. 1)

**The need for Bayesian modelling of MIL** The empirical study constitutes one motivation for our probabilistic approach: bring more stochasticity to the optimization process so that the convergence of weight and attention variables does not heavily depend on the initialization. Thus, we study the probabilistic counterpart of MIL, namely Bayes-MIL, which is further potentially beneficial in the following three aspects:

- *Better optimization*: Besides stochasticity in optimization, by turning  $\mathbf{a}$  to stochastic nodes, explicit regularization could be imposed for generating correct attention (see Fig. 2E 2J).
- *New measure of interpretability*: By learning a proper posterior over the parameters  $p(\pi|\mathcal{D})$ , we can induce the uncertainty over patches by  $p(\mathbf{a}|\mathbf{H}^*, \mathcal{D}) = \int p(\mathbf{a}|\pi, \mathbf{H}^*)p(\pi|\mathcal{D})d\pi$ , where  $\mathbf{H}^*$  is the testing data. Patch-level uncertainty can potentially indicate which patches the model is uncertain about, becoming a new measure of interpretability for MIL. In other words, the patch-level uncertainty is leveraged for localizing positive areas.
- *Calibrated uncertainty*: A properly learned posterior  $p(Y|\mathbf{H}^*, \mathcal{D})$  provides a *calibrated uncertainty* on  $Y$ , which is crucial in the application of medical imaging but ignored in existing methods.

### 3 BAYESIAN MULTIPLE INSTANCE LEARNING

#### 3.1 INSTANCE-LEVEL DISENTANGLED UNCERTAINTY

The basic framework of Bayes-MIL is introduced in this section. To obtain uncertainty over the attention, the principled way is to assume a prior distribution  $p(\pi)$  on the parameters of attention function  $f_{\pi}(\cdot)$ , and let the model learn a posterior  $p(\pi|\mathcal{D})$ . The other way is to directly learn an empirical posterior distribution  $p(\pi|\mathcal{D})$ , e.g., ensembles. The posterior induces the uncertainty over the attention by  $p(\mathbf{a}|\mathbf{H}^*, \mathcal{D}) = \int p(\mathbf{a}|\pi, \mathbf{H}^*)p(\pi|\mathcal{D})d\pi \approx \frac{1}{S} \sum_{s=1}^S f_{\pi_s}(\mathbf{H}^*)$ ,  $\pi_s \sim p(\pi|\mathcal{D})$ , where the attention function directly models the conditional distribution, i.e.,  $p(\mathbf{a}|\pi, \mathbf{H}) = f_{\pi}(\mathbf{H}) = f(\pi, \mathbf{H})$ .

For making  $p(\mathbf{a}|\pi, \mathbf{H})$  a strict distribution, the softmax function in (1) could be leveraged for normalization,  $\sum_k a_k = 1$ . Then,  $\mathbf{a}$  is a vector over the simplex, representing one categorical distribution. However, in this case, we could only calculate the uncertainty for the single categorical distribution. To extract the patch-level uncertainty, there should be one probabilistic distribution for the attention of each patch. Therefore, we need to model the distribution for patches,  $p(a_k|\pi, \mathbf{H}) = f_{\pi}(\mathbf{H})_k = f(\pi, \mathbf{H})_k$  and normalize for each patch. To this end, we replace the softmax function in (1) by an element-wise sigmoid function,

$$a_k = \frac{1}{1 + \exp\{-\mathbf{m}^T(\tanh(\mathbf{V}_1 \mathbf{h}_k^T) \odot \text{sigmoid}(\mathbf{V}_2 \mathbf{h}_k^T))\}}, \quad \mathbf{z} = \frac{1}{\sum_k a_k} \sum_k a_k \mathbf{h}_k \quad (3)$$

where the parameters  $\pi = \{\mathbf{m}, \mathbf{V}_1, \mathbf{V}_2\}$  are shared across all patches. The second normalization when computing  $\mathbf{z}$  is for numerical stability.

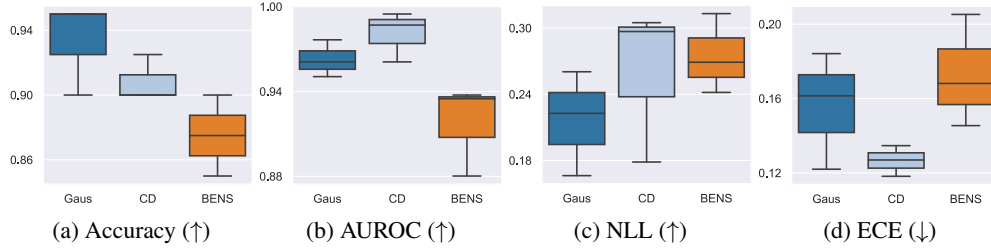


Figure 3: The comparisons of multiplicative Gaussian, concrete dropout and BatchEnsemble for modelling probabilistic MIL weights. The slide-level results are reported. NLL is negative log-likelihood and ECE is the expected calibration error for measuring uncertainty calibration.

With this treatment, the patch-level model uncertainty and data uncertainty for each patch  $a_k$  could be extracted by (Houlsby et al., 2011; Gal et al., 2017b):

$$\underbrace{\mathcal{I}[a_k, \boldsymbol{\pi} | \mathbf{H}^*, \mathcal{D}]}_{\text{Model Uncertainty}} = \underbrace{\mathcal{H}[\mathbb{E}_{p(\boldsymbol{\pi} | \mathcal{D})}[p(a_k | \boldsymbol{\pi}, \mathbf{H}^*)]]}_{\text{Total Uncertainty}} - \underbrace{\mathbb{E}_{p(\boldsymbol{\pi} | \mathcal{D})}[\mathcal{H}[p(a_k | \boldsymbol{\pi}, \mathbf{H}^*)]]}_{\text{Data Uncertainty}} \quad (4)$$

We explore three different methods for modeling the probabilistic weights  $\boldsymbol{\pi}$ : Batch-Ensemble (Wen et al., 2020; Dusenberry et al., 2020) as the empirical posterior, Concrete Dropout (Gal et al., 2017a) and multiplicative Gaussian noise (Kingma et al., 2015; Molchanov et al., 2017; Cui et al., 2021) that assume prior on  $\boldsymbol{\pi}$  and derive the posterior based on variational inference. As there is no direct indication of which set of posterior and prior, we should choose, we validate them empirically with the proposed framework. The experiments are conducted on CAMELYON16 dataset with 5 splits of training and validation sets. For the Batch-Ensemble, we use an ensemble of size 4. For Concrete Dropout and multiplicative Gaussian noise, we take 4 samples from the learned posterior during inference time. The multiplicative Gaussian method is selected for the probabilistic weights due to a high accuracy and a low NLL, shown in Fig. 3.

For neural networks, computing the posterior distribution using the Bayes rule requires computing intractable integrals over  $\boldsymbol{\pi}$ . In this paper, to be consistent with the stochastic attention function modelling, we use variational inference for approximating the posterior. Specifically, the posterior  $p(\boldsymbol{\pi} | \mathcal{D})$  is approximated by  $q_\phi(\boldsymbol{\pi})$ , by minimizing the Kullback-Leibler (KL) divergence  $\text{KL}[q_\phi(\boldsymbol{\pi}) || p(\boldsymbol{\pi} | \mathcal{D})]$ , where  $\phi$  are the variational parameters. This is equivalent to maximizing the evidence lower bound:

$$\max_{\phi} L_\phi = L_{\mathcal{D}}(\phi) - \text{KL}[q_\phi(\boldsymbol{\pi}) || p(\boldsymbol{\pi})], \quad L_{\mathcal{D}}(\phi) = \sum_{i=1}^N \mathbb{E}_{q_\phi(\boldsymbol{\pi})}[\log p(Y_i | \mathbf{H}_i, \boldsymbol{\pi})], \quad (5)$$

where  $L_{\mathcal{D}}(\phi)$  is the expected data log-likelihood and  $p(\boldsymbol{\pi})$  is the prior over  $p(\boldsymbol{\pi})$ .

### 3.2 SLIDE-DEPENDENT PATCH REGULARIZER

Although the proposed model can naturally visualize different types of uncertainty at the patch level and better localize the ROI, we can further *leverage the slide-level information for building a strong regularization on the ROI localization*. The aim is to explicitly encode the underlying logic of MIL into the training process, instead of letting the model explore implicit decision rules during training.

Recall  $Y$  is the slide-level label and  $\{y_k\}_{k=1}^K$  are the labels for patches which are unknown during training. The intuition is based on the logic under the MIL framework that a slide is negative when all patches are negative, while a slide is positive when there are one or more positive patches. For binary label,  $Y = 0$  iff  $\sum_k y_k = 0$ , and  $Y = 1$  otherwise. The following design principle can be drawn by simple deduction: When  $Y = 0$ , the attention distributions  $\{p(a_k | \boldsymbol{\pi}, \mathbf{H})\}_{k=1}^K$  must concentrate on the negative side ( $a_k = 0$ ) to guarantee  $y_k = 0$ . When  $Y = 1$ , the attention distributions are free to select either the positive ( $a_k = 1$ ) or negative sides. However, for precise localization, the attention distributions must concentrate on either the positive or negative side with high confidence.

This design principle is implemented by a variational inference framework by finding a regularizer on patches that is dependent on the slide label. Specifically, we choose the logit-normal distribution for  $a_k$ , due to its expressiveness over the simplex. The regularizer (a non-strict prior) is defined as

$$p(a_k | Y) = (1 - Y) \mathcal{LN}(\mu_0, \sigma_0) + Y \mathcal{LN}(\mu_1, \sigma_1), \quad (6)$$

where  $\mathcal{LN}(a_k | \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \frac{1}{a_k(1-a_k)} e^{-\frac{\text{logit}(a_k) - \mu}{2\sigma^2}}$  is the logit-normal distribution.  $\{\mu_0, \sigma_0\}$  and  $\{\mu_1, \sigma_1\}$  are the pairs of mean and variance for the negative slides and the positive slides, respectively. As shown in Fig. 4, by setting the parameters for the regularizer, we implement the design

principle:  $p(a_k|\mu_0, \sigma_0)$  concentrates on the negative side and  $p(a_k|\mu_1, \sigma_1)$  allows selecting either the negative or the positive side with high confidence.

For the MIL model, we divert  $\{a_k\}_{k=1}^K$  from induced distributions to stochastic nodes. The approximate posterior is defined as  $q(\mathbf{a}|\boldsymbol{\mu}, \boldsymbol{\sigma}) = \mathcal{LN}(\mu_k = f_\mu(\boldsymbol{\pi}, \mathbf{H})_k, \sigma_k = f_\sigma(\boldsymbol{\pi}, \mathbf{H})_k)$  and  $\sum_k \text{KL}[q(a_k|\boldsymbol{\mu}, \boldsymbol{\sigma})||p(a_k|Y)]$  is used as a regularization term during training, which we denote as the slide-dependent patch regularizer (SDPR). When a negative slide is given, the mode of the attention posterior is pushed to the negative side, generating low attention values over all patches. When a positive slide is given, the posterior will be trained to select a positive mode or a negative mode, with only high concentration. This produces dense localization of the ROI. The derivation of KL divergence for the regularization is in the Appendix.

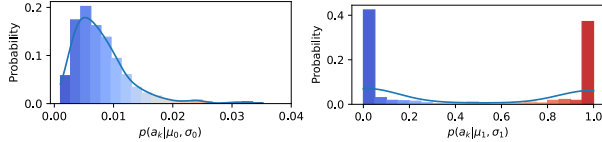


Figure 4: The visualization of density (curves) of the regularizer for the (left) negative and (right) positive slides. Samples bars are visualized with color based on the attention value.

The SDPR is also beneficial for improving the slide-level performance of MIL. The reason is it explicitly constrains the model to generate low attention values for all patches from the negative slides. Note that the major obscurity in MIL comes from the positive slides, where positive patches and negative patches coexist.

There is *no obscurity* for negative slides as all patches are negative by definition. However, the previous MIL models are still free to generate high attention for patches from negative slides, which neglect the underlying logic of MIL.

### 3.3 ENCODING SPATIAL INFORMATION VIA APPROXIMATION TO CONVOLUTIONAL CRF

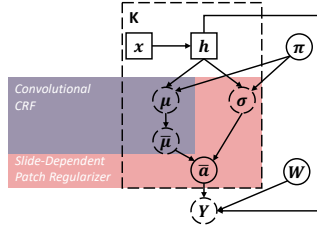


Figure 5: The graphical representation of Bayes-MIL. Boxes are deterministic nodes. Circles are stochastic nodes. Dashed circles are nodes of induced distributions.

The spatial information between patches are important for modelling of MIL for WSI recognition and localization, as there exist patches are spatially correlated (Shao et al., 2021). Here, we study how to encode the spatial information from a Bayesian perspective. The intuition is to let the neighboring attention posteriors  $q(a_k)$  influence each other.

One principled method for encoding the spatial information is to use conditional random field (CRF). Assume  $\bar{\mathbf{a}}$  is the output attention variable of CRF and  $\Theta$  is the input variable in CRF, from data. With CRF, the distribution of  $p(\bar{\mathbf{a}}|\mathbf{m})$  could be directly modelled with  $\mathbf{m}$  as the spatial features of patches (Zheng et al., 2015; Teichmann & Cipolla, 2018).

$$p(\bar{\mathbf{a}}|\mathbf{m}) = \frac{1}{Z(\mathbf{m})} e^{-E(\bar{\mathbf{a}}|\mathbf{m})}, \quad (7)$$

$$E(\bar{\mathbf{a}}|\mathbf{m}) = \sum_k \psi_u(\bar{\mathbf{a}}_k|\mathbf{m}) + \sum_{k \neq j} \psi_p(\bar{\mathbf{a}}_k, \bar{\mathbf{a}}_j|\mathbf{m})$$

where  $\mathbf{m} = [\boldsymbol{\omega}, \boldsymbol{\eta}, \mathbf{a}]$  contains the coordinates of the each patch over the slide,  $\boldsymbol{\omega} = [\omega_1, \dots, \omega_K]^T$  and  $\boldsymbol{\eta} = [\eta_1, \dots, \eta_K]^T$ .  $\psi_u$  is the unary potential that contains information only from the single patch and  $\psi_p$  is the pair-wise potential that captures the pair-wise correlation between patches. However, calculating the pair-wise correlation between patches has  $\mathcal{O}(K^2)$  complexity. For efficiency, we only consider the local dependency around  $a_k$ , by setting the non-local pair-wise correlation to be 0. This is equivalent to applying a softmax normalization and a convolution operation for the provided input  $\mathbf{a}$ . We define an function for the proposed convolutional CRF,  $\bar{\mathbf{a}} = C_{\mathbf{w}, \mathbf{h}}(\mathbf{a})$ . Specifically,  $\hat{\mathbf{a}} = \text{softmax}(\mathbf{a})$ ,  $\hat{\mathbf{a}} = \text{reshape}(\mathbf{w}, \mathbf{h}, \hat{\mathbf{a}})$ ,  $\bar{\mathbf{a}} = \text{convolution}(\hat{\mathbf{a}}, \mathcal{K})$ , where  $\mathcal{K}$  is a convolutional kernel with predefined hyperparameters. The detailed derivations and algorithms for the full convolutional CRF for BayesMIL are in the Appendix. For a precise estimation of the variable  $\bar{\mathbf{a}}$ , the following Monte-Carlo estimator is required,

$$\mathbb{E}[\bar{\mathbf{a}}] = \mathbb{E}_{q(\mathbf{a}|\boldsymbol{\mu}, \boldsymbol{\sigma})}[C_{\mathbf{w}, \mathbf{h}}(\mathbf{a})] \approx \frac{1}{S} \sum_{s=1}^S [C_{\mathbf{w}, \mathbf{h}}(\mathbf{a}_s)], \quad \mathbf{a}_s \sim q(\mathbf{a}|f_\mu(\boldsymbol{\pi}, \mathbf{H}), f_\sigma(\boldsymbol{\pi}, \mathbf{H})) \quad (8)$$

The full convolutional CRF has the most promising results on real-world datasets. However, the repetitive sampling makes the training less efficient. To bypass the heavy-load of sampling, we use

a first-order Taylor expansion to approximate the expectation (see Appendix). The final format of approximation could be written as

$$q(\bar{a}_k | \bar{\mu} = C_{w,h}(\boldsymbol{\mu}), \boldsymbol{\sigma}) = \mathcal{LN}(\bar{\mu}_k = C_{w,h}(f_{\mu}(\boldsymbol{\pi}, \mathbf{H}))_k, \sigma_k = f_{\sigma}(\boldsymbol{\pi}, \mathbf{H})_k), \quad (9)$$

for imposing slide-level patch regularizer, as well as uncertainty disentanglement and visualization.

### 3.4 PUTTING IT ALL TOGETHER

The graphical model for the Bayes-MIL is shown in Fig. 5. For the classifier with parameter  $\mathbf{W}$ , we also use a multiplicative Gaussian modelling (Kingma et al., 2015) as the posterior and posterior for calibrated uncertainty. The final objective is:

$$\max_{\phi} L_{\mathcal{D}}(\phi) - \lambda_0 R_{\pi, \mathbf{W}}(\phi) - \lambda_1 R_{\bar{a}}(\phi) \quad (10)$$

$$L_{\mathcal{D}}(\phi) = \sum_{i=1}^N \mathbb{E}_{q_{\phi}(\boldsymbol{\pi}, \mathbf{W})} [\log p(Y_i | \mathbf{H}_i, \boldsymbol{\pi}, \mathbf{W})], \quad R_{\pi, \mathbf{W}}(\phi) = \text{KL}[q_{\phi}(\boldsymbol{\pi}, \mathbf{W}) || p(\boldsymbol{\pi}, \mathbf{W})],$$

$$R_{\bar{a}}(\phi) = \sum_{i=1}^N \sum_{k=1}^K \text{KL}[q(\bar{a}_k | C_{w,h}(f_{\mu}(\boldsymbol{\pi}, \mathbf{H}_i)), f_{\sigma}(\boldsymbol{\pi}, \mathbf{H}_i)) || p(\bar{a}_k | Y_n)]$$

where  $R_{\pi, \mathbf{W}}(\phi)$  is KL term for the probabilistic weights and  $R_{\bar{a}}(\phi)$  is the SDPR for the attention.  $\lambda_0$  and  $\lambda_1$  are the trade hyperparameters for two terms respectively.

## 4 RELATED WORK

**Attention-based multiple instance learning** For improving interpretability and performance of multiple instance learning (MIL), ABMIL (Ilse et al., 2018) first introduces the attention mechanism for embedding-based MIL. DSMIL (Li et al., 2021) considers contrastive learning for feature extraction and builds the global connection between patch attentions. TransMIL (Shao et al., 2021) proposes a correlated MIL and implements it by multi-head self-attention and spatial information encoding for full global correlation. CLAM (Lu et al., 2021) extends ABMIL to the case of multiple classes and builds integrated toolbox for visualizing the uncertainty. Although the attention has been extensively adopted for the MIL interpretability, wrong attention is still being generated for the negative bags (see Fig. 2 and Fig. 6). This is not tolerable in the application of medical imaging, where the interpretability of model is crucial. The reason is, no methods verify the convergence of attention and impose the constraint on correcting attention from the MIL principle. In this work, these issues are carefully resolved and a new probabilistic framework is proposed.

**Orthogonal works** Chen et al. (2022) proposes a large-scale vision transformer solution to simultaneously learn the feature and classifier for WSI. Bayes-MIL freezes the feature extractor during training following a normal setup in MIL. Zhang et al. (2022) proposes a two-stage feature distillation MIL framework for enhancing the performance. Bayes-MIL studies the fundamental interpretability problem in the one-stage MIL framework. See Appendix for a review of uncertainty in DNNs.

## 5 EXPERIMENTS

The proposed methods are evaluated on two standard WSI datasets: CAMELYON16 (Bejnordi et al., 2017) and CAMELYON17 (Bandi et al., 2018). CAMELYON16 contains 400 hematoxylin and eosin (H&E) stained WSI of sentinel lymph node for breast cancer, labeled as normal or tumor classes. CAMELYON17 contains 1000 WSI of the same type, labeled with normal or different stages (pN-stage) of the breast tumor. Since our paper mainly studies the interpretation of MIL, we treat these stages as one class, generating binary slide-level labels (normal and tumors). We leverage the CLAM testbed for the implementation. A ResNet-50 is used for feature extraction, consistent with previous methods. Each result is obtained with *10-fold* splits of training/validation/testing sets, which is a more thorough evaluation than previous papers. Other hyperparameters are listed in the Appendix. The codes are submitted as supplemental. In our evaluation, we consider 3 variants of BayesMIL: 1) Bayes-IL-Vis is the basic Bayesian MIL in Sec. 3.1, 2) Bayes-MIL-SDPR is the model with slide-dependent patch regularizer from Sec. 3.2, 3) Bayes-MIL-Full is the whole model. The hyperparameters are  $\lambda_0 = 10^{-8}$  and  $\lambda_1 = 10^{-12}$ . Results on CAMELYON17 show our method is better than the CLAM baseline in Patch-localization and Slide-classification, shown in the Appendix.

**Patch-level tumor region localization** We first evaluate the patch-level results on tumor region localization. The tumor region localization uses the Patch-Precision, Patch-FROC and Patch-AUC

Table 1: Results on CAMELYON16: (left) Patch-level localization results using Patch-FROC (P-FROC), Patch-Precision (P-Prec.), Patch-AUROC (P-AUC); (right) Slide-level classification results using Slide-Accuracy (S-Acc.), Slide-AUROC (S-AUC) and Slide-Calibration (S-ECE).

	Patch-level			Slide-level		
	P-Prec. ( $\uparrow$ )	P-FROC ( $\uparrow$ )	P-AUC ( $\uparrow$ )	S-Acc. ( $\uparrow$ )	S-AUC ( $\uparrow$ )	S-ECE ( $\downarrow$ )
DSMIL	0.1030	0.4443	0.7719	0.8682 $\pm$ 0.05	0.8944 $\pm$ 0.06	0.3798 $\pm$ 0.06
CLAM	0.6068	0.4792	0.8839	0.8650 $\pm$ 0.06	0.9177 $\pm$ 0.04	0.1738 $\pm$ 0.06
CLAM-T	0.7068	0.4830	0.8884	—	—	—
TransMIL	0.1726	0.4797	0.8644	0.8837 $\pm$ 0.02	0.9307 $\pm$ 0.04	0.1436 $\pm$ 0.04
Bayes-MIL-Vis	0.7140	0.4797	0.8995	0.8825 $\pm$ 0.05	0.9164 $\pm$ 0.06	0.1702 $\pm$ 0.05
Bayes-MIL-SDPR	0.7458	0.4856	0.9001	0.8875 $\pm$ 0.05	0.9432 $\pm$ 0.05	0.1621 $\pm$ 0.05
Bayes-MIL-APCRF	<b>0.8107</b>	<b>0.4919</b>	<b>0.9129</b>	<b>0.9000</b> $\pm$ 0.04	<b>0.9479</b> $\pm$ 0.05	<b>0.122</b> $\pm$ 0.04

Table 2: The patch-level localization results when using probabilistic attention values, model uncertainty and data uncertainty for our methods, evaluated on CAMELYON16.

	P-Prec. ( $\uparrow$ )			P-FROC ( $\uparrow$ )		
	Attn	Model Unc.	Data Unc.	Attn	Model Unc.	Data Unc.
Bayes-MIL-Vis	0.7107	0.6999	<b>0.7140</b>	0.4796	0.4788	<b>0.4797</b>
Bayes-MIL-SDPR	0.7445	0.7396	<b>0.7458</b>	0.4781	<b>0.4856</b>	0.4849
Bayes-MIL-APCRF	0.8078	0.8033	<b>0.8107</b>	0.4813	<b>0.4919</b>	0.4879

metrics. The Patch-Precision is calculated by averaging the precision of classifying the patches. Each method provides a measure (attention  $a_k$  or normalized uncertainty value  $\mathcal{U}_k$ ) for its ROI. For calculating the precision, we compare the measure with three thresholds (0.1, 0.5, 0.9) for all methods and report the best results. Concretely, if  $a_k$  or  $\mathcal{U}_k$  is greater than the threshold, the patch is predicted as positive, otherwise negative. The Patch-FROC is defined as the average sensitivity (recall) at 6 predefined false positive rates: 1/4, 1/2, 1, 2, 4 and 8 FPs per WSI (Li et al., 2021). The Patch-AUROC evaluates the averaged area under ROC over patches. The ground truth is the doctor’s marking the region of the tumor.

For TransMIL, we take the diagonal of attention map from the last multi-head attention module as the measure. We add another baseline, CLAM-T, on CLAM training with temperature  $T = [0.2, 0.5, 2, 5]$  on softmax function,  $\mathbf{a} = \text{softmax}(\cdot, T)$ , as a method for manually adjusting the density of localization, and select the best results for different metrics. For the Bayes-MIL-Vis, Bayes-MIL-SDPR and Bayes-MIL-APCRF methods, we take the best results from the normalized data uncertainty, the normalized model uncertainty and the probabilistic attention, with MC integration over 16 samples.

The results for patch-level localization of tumor regions are shown in Tab. 1 (left). The Bayesian modelling of MIL (Bayes-MIL-Vis) generally improves the patch-level visualization over other methods. Tuning the temperature for MIL attention (CLAM-T) could marginally improve localization results, however, it requires exhaustively tuning the temperature. For Bayes-MIL, the precision is improved by a large margin (by 0.1072-0.611), showing the advantage of using uncertainty for localizing positive area. Including SDPR and the approximate CRF improve the patch-level results consistently, and the full model Bayes-MIL-APCRF achieves the best performance on the 3 metrics.

Tab. 2 shows the detailed results with probabilistic attention, model uncertainty and data uncertainty for the three proposed methods. The probabilistic attention benefits from the Bayes-MIL model, making it competitive at tumor region localization. The best results for precision are from the disentangled data uncertainty, indicating the data uncertainty is able to precisely localize the tumor region. The reason is that data uncertainty captures the rareness of features of positive patches indicating anomalies, such as cancerous cells, in the whole dataset. The model uncertainty is better on the P-FROC metric, which evaluates the recall at different false positive rates. This shows that the uncertainty induced from parameter distributions in  $\pi$  tends to cover the positive patches better.

**Slide-level evaluation** We next evaluate the slide-level classification performance. Tab. 1 (right) shows the slide-level performance on CAMELYON16, measured by Accuracy, AUC, and ECE. The results of DSMIL and TransMIL are from their papers, as their reported performance is better than our reimplementations. Compared with DSMIL and CLAM, Bayes-MIL-Vis has a higher accuracy, a lower calibration error and a similar AUC. When including our SDPR, Bayes-MIL-SDPR has better accuracy and AUC than TransMIL.

**Why is Bayes-MIL good at slide classification and ROI localization?** To understand the reason behind the good performance of the proposed framework, we visualize attention distribution in Fig. 6. As shown in Fig. 6a and Fig. 6b, DSMIL and CLAM generate similar attention distributions for positive and negative slides. For the negative slides, DSMIL still generates a nearly normal



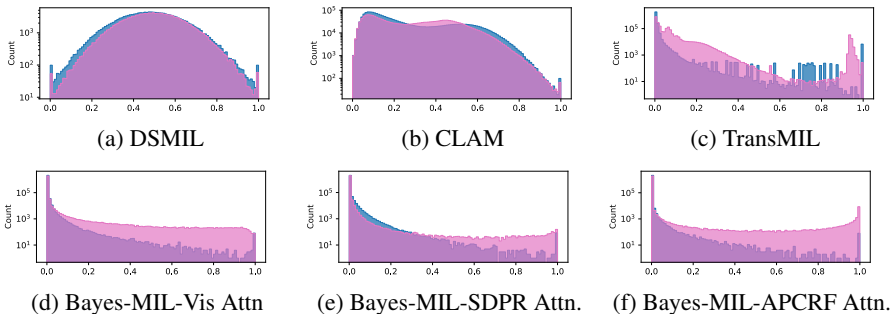


Figure 6: The log-scale histogram of attention for different methods. The pink and blue curves are for positive and negative slides, respectively. All testing images in CAMELYON16 are used, which contain 2M positive and 2M negative patches.

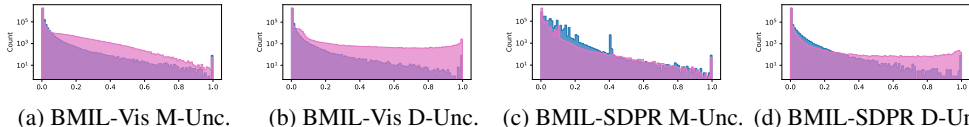


Figure 7: The log-scale histogram of normalized uncertainty for Bayes-MIL. M and D stand for model and data. APCRf has similar histograms as SDPR, thus not shown for saving space.

distribution of attention, which deviates from the MIL principle. For positive slides, the densities on the positive side ( $a \approx 1$ ) of DSMIL and CLAM are low. Note that the vanilla version BMIL-Vis has different attention distributions for positive and negative slides, but does not generate high density on the positive side for positive slides. Bayes-MIL with SDPR is able to push the attention of negative slides to the negative side ( $a \approx 0$ , the blue curve in Fig. 6c), while generating a U-shape distribution for the positive slides (red curve in Fig. 6c), which corresponds to our design of SDPR with the ideal case shown in Fig. 4. This design benefits the following aspects:

- For visualization, Bayes-MIL generates concentrated high attention values for positive slides, while only generates correct values ( $a \approx 0$ ) for negative slides.
- For classification, learning correct attention guided by SDPR for the negative slides will benefit the classification of patches in the positive slides. TransMIL captures this attention distribution to some degree, so it performs well in slide-level classification.

Bayes-APCRF performs a smoothing operation on the mean of the stochastic attention nodes. This pushes the histogram to have a more distinct U-shape (the red U-shape curve in Fig. 6f) by aggregating over neighboring positive patches, which benefits the localization and obtains the best visualization scores in Tab. 1.

**Correctly Visualized ROI** Fig. 2 shows the visualization results on the negative (normal) and positive (tumor) WSI. The related methods all generate high attention for the negative slides, which is against of MIL principle. TransMIL only generates a small area of high attention for negative slides, however the performance on positive slides is poor. *Bayes-MIL is the only method that generate correct values ( $a \approx 0$ ) for the negative slides, while performing the best in localization on positive slides.*

**Measure of uncertainty** The measure of data uncertainty (Fig. 8b) and Fig. 7d) naturally captures the information from the patch distribution. It shows a similar distribution with the attention of Bayes-MIL-SDPR and the ideal case of SDPR in Fig. 4. This might be the reason why data uncertainty has the best performance in localizing the ROI. Furthermore, based on this observation on data uncertainty, the soundness of SDPR is empirically validated.

## 6 CONCLUSION

This paper analyzes the interpretability problem in existing attention-based multiple instance learning (MIL) models. Directly taking attention as a measure of the important instances empirically violates the MIL principle. To address this problem, we propose a probabilistic solution Bayes-MIL that provides new measures for interpretability. To ensure the validity of interpretation in MIL, a regularizer on attention is required. A slide-dependent patch regularizer is proposed for imposing explicit constraints to let the model learn under the MIL principle, which also improves the slide-level performance. The spatial information is further encoded, which improves both slide-level and patch-level performance. The analysis and visualization of data uncertainty distribution further validate the main idea and soundness of SDPR.

## REFERENCES

- Peter Bandi, Oscar Geessink, Quirine Manson, Marcory Van Dijk, Maschenka Balkenhol, Meyke Hermsen, Babak Ehteshami Bejnordi, Byungjae Lee, Kyunghyun Paeng, Aoxiao Zhong, et al. From detection of individual metastases to classification of lymph node status at the patient level: the camelyon17 challenge. *IEEE transactions on medical imaging*, 38(2):550–560, 2018.
- Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes Van Diest, Bram Van Ginneken, Nico Karssemeijer, Geert Litjens, Jeroen AWM Van Der Laak, Meyke Hermsen, Quirine F Manson, Maschenka Balkenhol, et al. Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *Jama*, 318(22):2199–2210, 2017.
- Gabriele Campanella, Matthew G Hanna, Luke Geneslaw, Allen Miraflor, Vitor Werneck Krauss Silva, Klaus J Busam, Edi Brogi, Victor E Reuter, David S Klimstra, and Thomas J Fuchs. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nature medicine*, 25(8):1301–1309, 2019.
- Richard J Chen, Chengkuan Chen, Yicong Li, Tiffany Y Chen, Andrew D Trister, Rahul G Krishnan, and Faisal Mahmood. Scaling vision transformers to gigapixel images via hierarchical self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16144–16155, 2022.
- Yufei Cui, Wuguannan Yao, Qiao Li, Antoni B Chan, and Chun Jason Xue. Accelerating monte carlo bayesian inference via approximating predictive uncertainty over simplex. *arXiv preprint arXiv:1905.12194*, 2019.
- Yufei Cui, Ziquan Liu, Qiao Li, Antoni B Chan, and Chun Jason Xue. Bayesian nested neural networks for uncertainty calibration and adaptive compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2392–2401, 2021.
- Michael Dusenberry, Ghassen Jerfel, Yeming Wen, Yian Ma, Jasper Snoek, Katherine Heller, Balaji Lakshminarayanan, and Dustin Tran. Efficient and scalable bayesian neural nets with rank-1 factors. In *International conference on machine learning*, pp. 2782–2792. PMLR, 2020.
- Yarin Gal, Jiri Hron, and Alex Kendall. Concrete dropout. *Advances in neural information processing systems*, 30, 2017a.
- Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *International Conference on Machine Learning*, pp. 1183–1192. PMLR, 2017b.
- Yarin Gal et al. Uncertainty in deep learning.
- Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011.
- Maximilian Ilse, Jakub Tomczak, and Max Welling. Attention-based deep multiple instance learning. In *International conference on machine learning*, pp. 2127–2136. PMLR, 2018.
- Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in neural information processing systems*, 30, 2017.
- Durk P Kingma, Tim Salimans, and Max Welling. Variational dropout and the local reparameterization trick. *Advances in neural information processing systems*, 28, 2015.
- 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*, pp. 14318–14328, 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.
- Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. *Advances in neural information processing systems*, 31, 2018.

- Andrey Malinin, Bruno Mlodozienec, and Mark Gales. Ensemble distribution distillation. *arXiv preprint arXiv:1905.00076*, 2019.
- Dmitry Molchanov, Arsenii Ashukha, and Dmitry Vetrov. Variational dropout sparsifies deep neural networks. In *International Conference on Machine Learning*, pp. 2498–2507. PMLR, 2017.
- Nikhil Naik, Ali Madani, Andre Esteva, Nitish Shirish Keskar, Michael F Press, Daniel Ruderman, David B Agus, and Richard Socher. Deep learning-enabled breast cancer hormonal receptor status determination from base-level h&e stains. *Nature communications*, 11(1):1–8, 2020.
- 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*, 34:2136–2147, 2021.
- Marvin TT Teichmann and Roberto Cipolla. Convolutional crfs for semantic segmentation. *arXiv preprint arXiv:1805.04777*, 2018.
- Yeming Wen, Dustin Tran, and Jimmy Ba. Batchensemble: an alternative approach to efficient ensemble and lifelong learning. *arXiv preprint arXiv:2002.06715*, 2020.
- Hongrun Zhang, Yanda Meng, Yitian Zhao, Yihong Qiao, Xiaoyun Yang, Sarah E Coupland, and Yalin Zheng. 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.
- Jingwei Zhang, Ke Ma, John Van Arnam, Rajarsi Gupta, Joel Saltz, Maria Vakalopoulou, and Dimitris Samaras. A joint spatial and magnification based attention framework for large scale histopathology classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3776–3784, 2021.
- Zizhao Zhang, Pingjun Chen, Mason McGough, Fuyong Xing, Chunbao Wang, Marilyn Bui, Yuanpu Xie, Manish Sapkota, Lei Cui, Jasreman Dhillon, et al. Pathologist-level interpretable whole-slide cancer diagnosis with deep learning. *Nature Machine Intelligence*, 1(5):236–245, 2019.
- Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip HS Torr. Conditional random fields as recurrent neural networks. In *Proceedings of the IEEE international conference on computer vision*, pp. 1529–1537, 2015.