

Interventional Bag Multi-Instance Learning On Whole-Slide Pathological Images

Tiancheng Lin^{1,2} Zhimiao Yu^{1,2} Hongyu Hu^{1,2} Yi Xu^{1,2*} Chang Wen Chen³

¹ Shanghai Key Lab of Digital Media Processing and Transmission, Shanghai Jiao Tong University

² MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University

³ The Hong Kong Polytechnic University, Hong Kong, China

{ltc19940819, carboxy, mathewcrespo, xuyi}@sjtu.edu.cn, changwen.chen@polyu.edu.hk

Abstract

Multi-instance learning (MIL) is an effective paradigm for whole-slide pathological images (WSIs) classification to handle the gigapixel resolution and slide-level label. Prevaling MIL methods primarily focus on improving the feature extractor and aggregator. However, one deficiency of these methods is that the bag contextual prior may trick the model into capturing spurious correlations between bags and labels. This deficiency is a confounder that limits the performance of existing MIL methods. In this paper, we propose a novel scheme, **Interventional Bag Multi-Instance Learning (IBMIL)**, to achieve deconfounded bag-level prediction. Unlike traditional likelihood-based strategies, the proposed scheme is based on the backdoor adjustment to achieve the interventional training, thus is capable of suppressing the bias caused by the bag contextual prior. Note that the principle of IBMIL is orthogonal to existing bag MIL methods. Therefore, IBMIL is able to bring consistent performance boosting to existing schemes, achieving new state-of-the-art performance. Code is available at <https://github.com/HHHedo/IBMIL>.

1. Introduction

The quantitative analysis of whole-slide pathological images (WSIs) is essential for both diagnostic and research purposes [13]. Beyond complex biological structures, WSIs are quite different from natural images in the gigapixel resolution and expensive annotation, which is thus formulated as a multi-instance learning (MIL) [9] problem: treating each WSI as a labeled bag and the corresponding patches as unlabeled instances. Such a *de facto* paradigm has been demonstrated in extensive tasks on WSIs, *e.g.*, classification [7, 15, 18, 46], regression [40, 41, 48] and segmentation [37]. The prevailing scheme for WSI classification — bag-level MIL — is depicted in Fig. 1a. Given the patchified

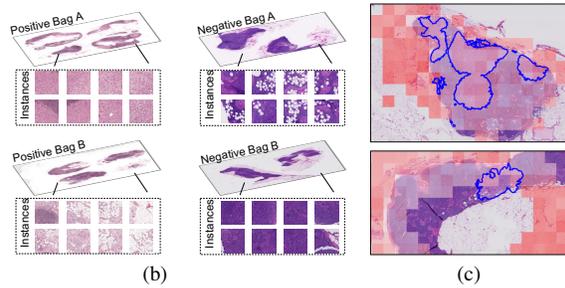
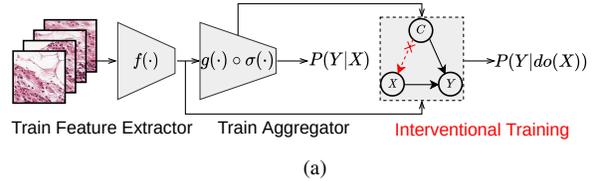


Figure 1. (a) Traditional scheme and our interventional training. (b) Dataset bias. (c) Unreasonable attention maps with right predictions.

images as instances, each instance is embedded in vectors by a feature extractor in the first stage. Second, for each bag, their corresponding instance features are aggregated as a bag-level feature for classification.

More and more new frameworks are proposed to improve the two stages following this scheme [19, 28, 31, 44]. It is convinced that learning better instance features and modeling more accurate instance relationships can bring better performance of MIL. While we have witnessed the great efforts, they still leave the “bag contextual prior” issue unsolved: the information shared by bags of the same class but irrelevant to the label, which may affect the final predictions. For example, in Fig. 1b, due to the dataset bias, most of the instances in the positive bags are stained pink but purple in the negative bags. The co-occurrence of specific color patterns and labels may mislead the model to classify bags by color statistics instead of the key instances — the more

*Corresponding author.

pink instances a bag contains, the more likely it is a positive bag. Fig. 1c illustrates another example: even if the prediction is correct, the underlying visual attention is not reasonable, where the high attention scores are put on the disease-irrelevant instances outside the blue curves in the bags. From the causal lens, the bag contextual prior is a confounder that opens up a backdoor path for bags and labels, causing spurious correlations between them. To suppress such a bias, we need a more efficient mechanism for the actual causality between bags and labels, *i.e.*, the bag prediction is based on the bag’s content (*e.g.*, key instances), which can not be fully achieved only by above mentioned new frameworks.

In fact, it is challenging to achieve unbiased bag predictions as such a bias happens in the data generation – the tissue preparations, staining protocols, digital scanners, *etc.* In this paper, we propose a novel MIL scheme, **Interventional Bag Multi-Instance Learning (IBMIL)**, to tackle this challenge. In particular, we propose a structure causal model (SCM) [24] to analyze the causalities among bag contextual prior, bags and labels. The key difference of IBMIL is that it contains another stage of interventional training (see Fig. 1a right). Given the aggregator trained in the second stage, instead of directly using it for inference via **likelihood**: $P(Y|X)$, we apply it for the approximation of confounders. With the confounders observed, we eliminate their effect via the backdoor adjustment formulation [23], where the intuitive understanding is: if a WSI model can learn from “purple” and “pink” positive/negative bags, respectively, then the bag context of color will no longer confound the recognition. Therefore, our IBMIL is fundamentally different from the existing scheme as we use a **causal intervention**: $P(Y|do(X))$ for bag prediction.

We conduct experiments on two public WSI datasets, *i.e.*, Camelyon16 [1] and TCGA-NSCLC. Experimental results show that IBMIL is agnostic to both feature extractors and aggregation networks, *i.e.*, it brings consistent performance boosting to all compared state-of-the-art MIL methods in the WSI classification tasks. Further ablation studies and analyses demonstrate the effectiveness of interventional training.

2. Related Work

2.1. Instance-level MIL on WSIs

Instance-level MIL represents each instance by a score and aggregates instance scores into a bag score. One widely used baseline is SimpleMIL [5, 8], which directly propagates the bag label to its instances. When applying SimpleMIL for WSIs, the unbalanced dataset could result in noisy instance-level supervision since a WSI (*e.g.*, Camelyon16) might only contain a small portion of a disease-positive tissue in clinic [19]. The following works in

this line improve this baseline via various modifications. *Cleaner annotations*: SemiMIL [34] directly introduces cleaner annotations for partial instances with the help of pathologists, where these annotated regions are assigned with larger weights as they carry higher confidence. *Instance selection*: PatchCNN [15] selects instances via a delicate thresholding scheme at both WSI and class levels. Similarly, Top-*k* MIL [6] only uses the top-*k* instances for each bag, but the fixed number of selected instances fails to make a trade-off between preserving clean instances and discarding noisy instances. RCEMIL [2] proposes rectified cross-entropy (RCE) loss to select instances in a softer manner, while the loss requires statistics of possible abnormal tissues among all WSIs. More recently, IMIL [20] summarizes the previous works from a causal lens and propose IMIL to select instance via causal intervention and effect. However, the performance of instance-level MIL methods is usually inferior to bag-level counterparts [35].

2.2. Bag-level MIL on WSIs

The instances are represented as embedding vectors and classified by bag-to-bag distance/similarity or a bag classifier [35]. Conducting bag-level MIL on WSI is non-trivial, because the intermediate results of all patches still need to be stored in memory for backpropagation. Therefore, some recently proposed frameworks separate the training of instance-level feature extractors and aggregation networks, resulting in a two-stage modeling approach [19, 28, 31, 44]. They contribute differently at both stages. For the feature extractor, they introduce different *architectures* from convolutional neural networks (CNNs) to transformer-based models [26], and *training paradigms* from ImageNet pre-training [27, 40, 41] to self-supervised learning [21, 31, 46]. Simultaneously, many works pay attention to new designs of aggregation networks, from non-parametric poolings, *e.g.*, max/mean-poolings [35], to learnable ones, *e.g.*, graph convolution networks [46] and attention mechanisms [17, 19, 28, 44]. Our work lies in this line but aims at empowering these existing methods. Thus the contributions are orthogonal.

2.3. Causal Inference in Computer Vision

Causal inference, a general framework, has been introduced to various computer vision tasks, including classification [16, 30, 42], semantic segmentation [4, 43], unsupervised representation learning [22, 32, 33] and so on. In MIL problems, StableMIL [45] takes “adding an instance to a bag” as a treatment for bag-level prediction, while IMIL [20] uses inverse probability weighting and causal effects for instance-level tasks. Unlike them, our IBMIL is based on backdoor adjustment formulation and works as a general framework to empower existing bag-level MIL for WSI classification tasks.

3. Method

3.1. Preliminaries

MIL formulation. Due to the gigapixel resolution and lacking fine-grained labels, performing downstream tasks on WSIs is formulated as an MIL problem, such that each WSI is treated as a *labeled bag* with corresponding patches as *unlabeled instances*. Take binary classification as an example, let $X = \{(x_1, y_1), \dots, (x_n, y_n)\}$ be a WSI bag, which contains n instances of x_i . The instance-level labels $\{y_i, \dots, y_n\}$ are unavailable. Under the standard MIL assumption, the bag label Y is further given by:

$$Y = \begin{cases} 0, & \text{iff } \sum y_i = 0 \\ 1, & \text{otherwise} \end{cases} \quad (1)$$

which can be modelled by max-pooling [15]. A general three-stage approach goes like 1) *Instance transformation*: a feature extractor $f(\cdot)$ is trained for instance-wise features b , 2) *Instance combination*: the pooling operation $\sigma(\cdot)$ is targeted for bag feature B , 3) *Bag transformation*: a downstream classifier $g(\cdot)$ is used for prediction, which can be formulated as:

$$b_i = f(x_i), B = \sigma(b_1, \dots, b_n), \hat{Y} = g(B), \quad (2)$$

where the pooling $\sigma(\cdot)$ should be a permutation-invariant function [17] for the spatial-invariant MIL method. Some works further absorb the classifier $g(\cdot)$ into the pooling operation $\sigma(\cdot)$, referred to as aggregator/aggregation networks. When applying the MIL methods for WSIs, it should be noted that 1) the diagnosis for WSI analysis can be based on different tissue regions with multiple concepts — the collective MIL assumption, 2) the bag length n for a WSI can be extremely large, *e.g.*, about 8,000 on average [19]. Therefore, the bag MIL methods for WSIs are with *learnable aggregators* and trained in a *two-stage* procedure, *i.e.*, training the feature extractor and aggregator stage by stage. Current works mainly follow this formulation and improve the framework from both feature extractor and aggregator, while our proposed method aims to empower existing works from a causal perspective.

Analysis MIL through causal inference. As shown in Fig. 2a, we formulate the MIL framework as a causal graph (*a.k.a.*, Pearl’s structural causal model or SCM [24]), which contains three nodes: X : whole-slide pathological image (bag), Y : bag label, C : bag contextual information.

$X \rightarrow Y$: This path indicates that the MIL model can learn to predict the bag label on the bag content, *e.g.*, key instances.

$C \rightarrow X$: This path indicates the generation of the whole-slide pathological image. Due to the differences in tissue preparations, staining protocols and digital scanners, the appearance of WSIs can be significantly affected, potentially introducing biases.

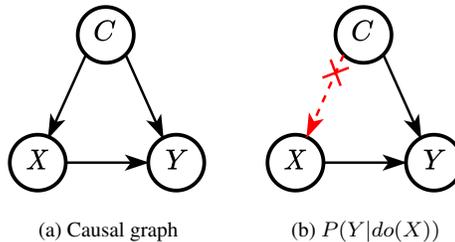


Figure 2. An illustration of causal graph for bag MIL framework.

$C \rightarrow Y$: This path indicates that the bag prediction is affected by the contextual prior information in the training dataset. For example, in Fig. 1b, an MIL model predicts all bags with purple color as positive regardless of content information related to the real label.

In the causal graph, C confounds X and Y via the backdoor path $X \rightarrow C \rightarrow Y$ and causes a spurious correlation between them, which prevents learning robust bag MIL models. For example, the model may wrongly predict the bags when the data are out-of-distribution, *i.e.*, with different context prior. An ideal MIL method should capture the true causality between X and Y , but the conventional correlation of $P(Y|X)$ fails to do so, as such a spurious correlation is inevitable. Therefore, we instead seek to use the causal intervention $P(Y|do(X))$, where the *do*-operation $do(\cdot)$ means forcibly assigning a specific value to the variable X . As shown in Fig. 2b, it can be considered as a modification of the graph — cutting off the backdoor path, thus mitigating the bias caused by confounders. The ideal way of $do(\cdot)$ is the random controlled trials [25] — enumerating each bag with all possible contexts, which is impossible in practice. Next, we propose a practical intervention method to remove the confounding effect caused by the bag contextual prior.

3.2. Interventional Bag Multi-Instance Learning

We propose to use the backdoor adjustment formulation to achieve the causal intervention: $P(Y|do(X))$ for bag-level prediction. Formally, we have the backdoor adjustment for the graph in Fig. 2a as:

$$P(Y | do(X)) = \sum_i P(Y | X, h(X, c_i))P(c_i), \quad (3)$$

where c_i loops over the confounder set and $h(\cdot)$ is a function defined later in Eq. (5). Different from Bayes rule, in Eq. (3), c_i is no longer affected by X but subject to its prior $P(c_i)$, since the causal intervention forces X to incorporate each c_i fairly. Now, we are ready to introduce our interventional bag multi-instance learning stage by stage. Fig. 3 illustrates the overview of IBMIL.

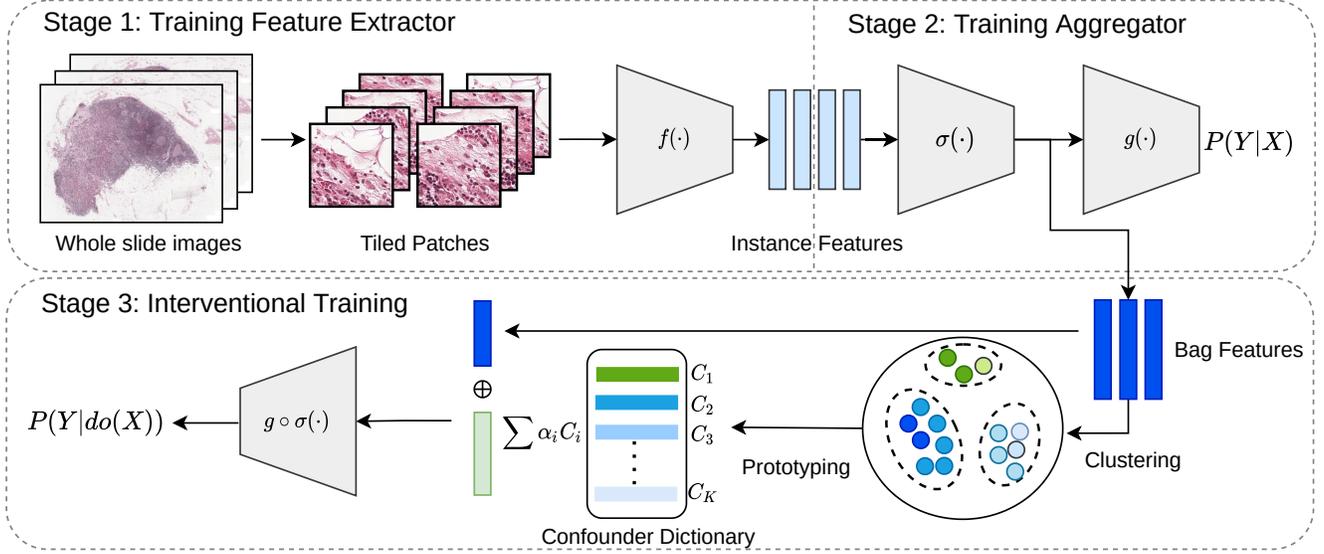


Figure 3. Overview of our proposed Interventional Bag Multi-Instance Learning (IBMIL). Our contribution is to introduce interventional training to the traditional two-stage scheme.

Stage 1: Training feature extractors. We learn a feature extractor $f(\cdot)$ on the patchified images of WSIs $\{x_1, \dots, x_n\}$, aiming at encoding each instance as a discriminative feature vector.

Stage 2: Training aggregators. Given the features of instances $\{b_1, \dots, b_n\}$, the aggregator employs MIL pooling $\sigma(\cdot)$ to assemble them into a bag feature B sequentially or simultaneously, and a classifier $g(\cdot)$ for discrimination. Formally, the loss for training aggregator is defined as:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N Y_i \log \hat{Y}_i + (1 - Y_i) \log (1 - \hat{Y}_i), \quad (4)$$

where N is number of bags in the training set. Note our IBMIL is no not limited to specific feature extractor or aggregators, including the architectures and training paradigms. Please refer to Sec. 4 for our choices.

Stage 3: Causal intervention via backdoor adjustment.

The traditional two-stage bag MIL stops at stage 2 and uses the trained models for inference directly. Instead, we introduce another stage of interventional training, which needs the practical implementation of Eq. (3). Note that backdoor adjustment assumes that we can observe and stratify the confounders of a bag context. Thanks to the powerful ability of deep MIL models, context information is naturally encoded in the higher-level layers [20, 43]. To constitute the confounder set, we use a confounder dictionary $C = [c_1, \dots, c_K]$ for approximation, as collecting all confounders is impossible. Given the trained feature extractor and aggregator, we use K -means over all the bag features in the training set, partitioning the bags into clusters. We average the bag features of each cluster to represent

a confounder stratum c_i , resulting in a confounder dictionary with the shape of $d \times K$, where d is the dimension of bag features. Note that our approximation is reasonable in that these global clusters are susceptible to the visual biases [29], which is exactly the confounders. Then, we define:

$$h(X, c_i) = \alpha_i c_i, \quad (5)$$

$$[\alpha_1, \dots, \alpha_K] = \text{softmax} \left(\frac{(W_1 B)^T (W_2 C)}{\sqrt{I}} \right),$$

where $B = \sigma(f(X))$ is the bag feature, $W_1, W_2 \in \mathbb{R}^{l \times d}$ are two learnable projection matrices to project bag feature B and confounder C into a joint space, and \sqrt{I} is used for feature normalization [33]. Since the prediction comes from both bag X and confounder C (see Fig. 2a), we further define

$$P(Y | X, h(X, c_i)) = P(Y | B \oplus h(X, c_i)), \quad (6)$$

where \oplus denotes vector concatenation, and other implementations can be found in ablation studies. We assume $P(c_i)$ is a uniform prior of $1/K$ for a safe estimation, and a more reasonable assumption, *e.g.*, incorporating expert knowledge, will be our future work. Plugging Eq. (5), Eq. (6) and defined $P(c_i)$ into Eq. (3), we are ready to calculate $P(Y|do(X))$ via passing the network multiple times. In practice, to avoid the expensive cost, we further apply Normalized Weighted Geometric Mean [38] to move the outer sum into the Softmax:

$$P(Y | do(X)) \approx P \left(Y | B \oplus \sum_{i=1}^K \alpha_i c_i P(c_i) \right). \quad (7)$$

Thus, backdoor adjustment can be achieved by one feed-forward of the network.

3.3. Justification

In our implementation of the causal intervention, there are some aspects we need to discuss further.

Compatible with large-scale unlabelled datasets. We constitute the confounder set in an unsupervised fashion. One alternative implementation is to use the available bag labels for guidance, preserving the intra-class variation and capturing the class-relevant characteristics. There are two main reasons for our choice. 1) The unsupervised fashion makes our scheme compatible with large-scale unlabelled datasets, *e.g.*, The Cancer Genome Atlas (TCGA), for better approximation of confounders. 2) The confounder could be irrelevant to the class identity, *e.g.*, the stain color of positive and negative instances can be the same. We explore the other implementations in Sec. 4.3.

One possible more elegant scheme. As we need the trained aggregator to generate the bag features (the stage 2), one more stage is needed to retrain the aggregator (the stage 3). We are thus motivated to further simplify our scheme to avoid extra computational cost. Specifically, we can achieve the bag features by applying the traditional non-parametric aggregators, *e.g.*, max/mean-pooling, to the instances in a bag. It is inspired by the fact that these non-parametric aggregators serve as strong baselines, and we conjecture that statistic bag information they provide can be used for a reasonable approximation of confounders. Therefore, we can omit the stage 2. The experiment results in Sec. 4.3 support that our scheme can be more elegant.

Connection to other methods. *Embedding-based MIL:* As we approximate the confounder set based on bag features, these confounders can be seen as a denoised abstraction of bag features. From this perspective, we share the same spirit with the embedding-based MIL [35], *i.e.*, exploring the relations between bags. That means our IBMIL also explains the effectiveness of embedding-based MIL. *Color Normalization:* Some works [47] propose color normalization methods for H&E stained WSIs. However, color is just one of the confounders, and some confounders are even unobserved. Our method does not focus on color only, and thus is the more reasonable partially observed children of the unobserved confounder [11]; *Instance augmentation:* IMIL [20] uses strong instance augmentation to train the feature extractor for instance prediction. However, the augmentation may affect the statistical information in the bag. Therefore, our method is more suitable for bag MIL. Remix [39] proposes data augmentations for MIL by exploring the relations of instances, but our method explores the bag-level relations based on the causal theory.

4. Experiments

Dataset and evaluation protocol. We conduct the experiments on two public WSI datasets, *i.e.*, Camelyon16 [1] and TCGA-NSCLC. Camelyon16 is a dataset of H&E stained slides for metastasis detection in breast cancer, consisting of 399 WSIs. Following [20], we crop each WSI into 256×256 non-overlapping patches, and remove the background region. There are roughly 2.8 million patches at $20\times$ magnification in total, with about 7,200 patches per bag. TCGA-NSCLC includes two subtypes in lung cancer, *i.e.*, Lung Squamous Cell Carcinoma (LUSC) and Lung Adenocarcinoma (LUAD). The dataset consists of 1,054 WSIs. We directly used the patches released by [19], which are about 5.2 million patches at $20\times$ magnitude, with an average of 5000 patches for each bag. Following the evaluation protocol of [19], we use 270 training images and 129 test images for Camelyon16, and 836 training images and 210 test images for TCGA-NSCLC (some corrupted slides are discarded). We report the class-wise precision, recall, accuracy and area under the curve (AUC) scores.

Feature extractor. We adopt different network architectures with different training paradigms to thoroughly evaluate our IBMIL. **ResNet-18** [14] is a widely used CNN-based model in our community, and we adopt the ImageNet pre-trained one released by PyTorch. **ViT-small** [10] is a typical transformer-based model, which is good at modeling the long-range dependencies in the data. **CTransPath** [36] is a hybrid CNN and transformer architecture, customized for WSIs. We adopt the ViT pre-trained with MoCo V3 [3] and CTransPath pre-trained with semantically-relevant contrastive learning (SRCL), where the used data is about 15 million images from 9 datasets [36]. Please refer to the Supplementary for more details.

Aggregators for MIL models. We build our proposed method upon 4 SOTA methods. **ABMIL** [17] is a classic attention-based MIL, where the attention scores are predicted by a multi-layer perceptron (MLP). **DSMIL** [19], a dual-stream framework, jointly learns an instance and a bag classifier. The highest-score instance is further used to re-calibrate other instances into a bag feature. **TransMIL** [28] is a correlated MIL framework built on transformer to explore both morphological and spatial information, where self-attention is used for bag aggregation. **DFTD-MIL** [44] proposes to virtually enlarge the number of bags by introducing the concept of pseudo-bags, resulting in a double-tier MIL framework. To align with DSMIL, we use the maximum attention score selection (MaxS) for the feature distillation strategy. For more results of DTFD-MIL (MaxMinS), please refer to the Supplementary.

We use DSMIL’s code base for implementation and evaluation, and build other models based on their officially released codes. Since the feature extractors we use are all pre-trained, we can directly transform instances into feature

Table 1. Main results (%) on Camelyon16 and TCGA-NSCLC.

Performance Method			Camelyon16				TCGA-NSCLC				
			Precision	Recall	Accuracy	AUC	Precision	Recall	Accuracy	AUC	
ResNet-18 ImageNet pretrained	ABMIL	+IBMIL	86.71	81.71	84.50	84.07	82.75	85.84	81.43	88.95	
		Δ	+1.87	+5.43	+3.87	+6.36	+2.67	-0.67	+3.81	+2.31	
	DSMIL	+IBMIL	84.56	82.95	84.50	87.16	80.56	85.78	77.62	86.88	
		Δ	+5.61	+3.25	+3.87	+0.53	+1.42	+0.47	+2.38	+0.31	
	TransMIL	+IBMIL	85.43	81.06	83.72	81.29	85.46	85.31	85.24	90.70	
		Δ	-2.29	+1.87	0.00	+7.42	+0.34	+1.75	0.00	+1.84	
	DTFD-MIL (MaxS)	+IBMIL	84.85	80.09	82.95	82.77	82.29	83.77	81.90	88.91	
		Δ	+4.68	+6.42	+5.42	+6.74	+0.96	-0.81	+0.96	+1.59	
	CTransPath SRCL	ABMIL	+IBMIL	91.84	89.09	90.70	92.33	90.76	90.40	90.48	95.87
			Δ	+3.45	+3.22	+3.10	+1.50	+1.16	+1.53	+1.42	+1.04
DSMIL		+IBMIL	89.24	88.10	89.15	93.26	92.11	92.79	90.95	97.13	
		Δ	+1.81	+3.20	+2.31	+1.94	-0.06	+1.03	+0.48	+0.38	
TransMIL		+IBMIL	95.83	93.27	94.57	95.88	92.05	93.82	91.90	95.55	
		Δ	+0.35	+2.33	+1.55	+1.12	+1.89	-0.06	+1.91	+1.69	
DTFD-MIL (MaxS)		+IBMIL	95.87	94.54	95.35	96.18	90.41	88.40	88.57	94.88	
		Δ	+1.08	+0.65	+0.77	+0.10	+0.76	+2.49	+2.38	+1.69	
ViT MoCo V3		ABMIL	+IBMIL	89.95	83.97	86.82	83.94	88.72	88.51	88.57	92.71
			Δ	-2.01	+2.86	+0.78	+7.37	+0.37	+0.51	+0.48	+0.82
	DSMIL	+IBMIL	86.72	77.24	81.40	82.27	90.26	91.37	90.00	95.40	
		Δ	-1.64	+1.58	+0.77	+1.50	+1.26	-0.95	+0.48	+0.80	
	TransMIL	+IBMIL	94.21	92.93	93.80	94.38	94.26	93.83	93.81	96.67	
		Δ	-0.37	+0.31	0.00	+0.82	+0.08	+1.42	+0.48	+1.31	
	DTFD-MIL (MaxS)	+IBMIL	90.71	88.44	89.92	90.96	89.45	90.89	89.05	94.95	
		Δ	+1.94	+4.47	+3.10	+5.39	+0.89	-0.97	+0.95	+1.40	

vectors. For stages 2 and 3, all MIL models are optimized for 50 epochs with learning rate of 0.0001, and other settings are followed their official code. We set the number of confounder $K = 8$ and project dimension $l = 128$ by default for all the main experiments. See Supplementary for more details.

4.1. Experimental Results

We present the results on two benchmark WSI datasets, Camelyon16 and TCGA-NSCLC, covering binary class MIL with unbalanced bags and multiple class MIL with balanced bags, respectively. By “unbalanced”, it means only a small portion of positive instances in positive bags, *e.g.*, roughly <10% in Camelyon16 [19]. From Tab. 1, we observe that 1) IBMIL consistently improves all feature extractors with all aggregators (12 possible combinations) on

both datasets, which suggests that IBMIL is agnostic to feature extractors, aggregators and datasets. 2) In particular, we find the improvement on the ImageNet pre-trained ResNet is larger than others. For example, the average gain of AUC is 5.4% in Camelyon16 and 1.5% in TCGA-NSCLC. This is mainly because ResNet is more likely to learn context patterns as it is supervised trained on ImageNet [12], while the other two are self-supervised trained with strong data augmentations — the “physical intervention”. 3) Our IBMIL improves more on Camelyon16 than TCGA-NSCLC in most cases. The main reason is that the former is a binary class MIL with unbalanced bags, which suffer more severe bag contextual prior — learning the key instances is much harder than context information. Note that the performance could be further improved by tuning the number of confounders for each setting.

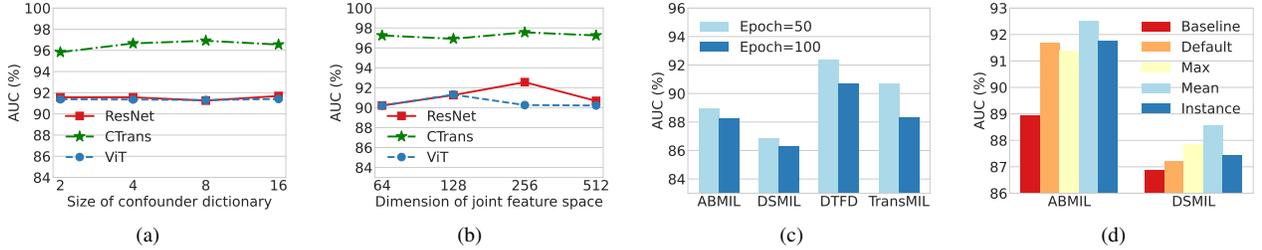


Figure 4. Ablation studies of (a) Size of confounder dictionary, (b) Dimension of joint feature space, (c) More training epochs of baselines, and (d) Means to achieve bag features. “Default” in (d) denotes our default 3-stage scheme.

4.2. Ablation on Model Design Variants

In Sec. 4.2 and Sec. 4.3, experiments are conducted on TCGA-NSCLC dataset with feature extractor of ResNet-18 and aggregator of ABMIL, unless specified otherwise.

Size of confounder dictionary. We ablate size K of the confounder dictionary on three feature extractors, including ResNet-18, CTransPath and ViT. From Fig. 4a, the performance of IBMIL is relatively robust to the size of confounder dictionary. Therefore, we need not elaborately tune this hyper-parameter and an arbitrary size within a wide range is able to boost the performance.

Dimension of joint feature space. As mentioned above, the confounders and bag features are projected into a joint feature space with a dimension of l and attention scores are calculated subsequently. We ablate dimension l on three feature extractors. The results in Fig. 4b reveal that performance does not improve monotonically with increased dimension and is saturated at $l = 256$. We choose a dimension of 128 as the default configuration.

Learnable vs. unlearnable confounders. We explore the effect of learnable and unlearnable confounders. For the former, we update them in an end-to-end manner via backpropagation.

Learnable	Precision	Recall	Accuracy	AUC
✓	83.81	83.82	83.81	90.82
✗	85.42	85.17	85.24	91.26

As can be seen, both of them outperform the baseline accuracy of 81.43%. However, freezing confounders during interventional training beats learnable confounders by 1.43% on accuracy. The reason may be that it is challenging to learn both confounders and interventional training with only bag-level labels, and introducing context-level supervision could be a solution [33]. We set confounders unlearnable as the default configuration.

Implementation of backdoor adjustment. We study the effect of different implementations of backdoor adjustment. Given a bag feature $B \in \mathbb{R}^d$ and the combination of confounders $\sum_{i=1}^K \alpha_i c_i P(c_i) \in \mathbb{R}^d$, we explore three variants to combine them, *i.e.* $B \star \sum_{i=1}^K \alpha_i c_i P(c_i)$

and $\star \in \{\oplus, +, -\}$, where $+/-$ is element-wise addition/subtraction.

Method	Precision	Recall	Accuracy	AUC
\oplus	85.42	85.17	85.24	91.26
$+$	84.99	84.68	84.76	89.28
$-$	84.70	84.18	84.29	90.14

We observe all these implementations can lead to performance improvements, which demonstrates the stability and effectiveness of the proposed intervention.

4.3. Analysis and Discussion.

Do improvements come from more epochs? Note that our proposed method requires an extra stage to train the aggregator. A natural question is whether we can improve baseline performance by training the baseline methods for as many epochs as the extra stage. Fig. 4c displays that more epochs do not bring about performance improvement, showing that our proposed method could empower baseline methods by backdoor adjustment rather than more training epochs. In most cases, training longer even brings performance degradation, which can be caused by the over-fitting problem in MIL [44].

Is stage 2 necessary? Recent MIL methods aggregate the instance features into a bag feature via the weighted average operator. The weights, also referred to as attention scores, are generated by parametric networks, which need an extra stage of training aggregators. Alternatively, we turn to three non-parametric settings to skip this stage and efficiently achieve the bag features. We consider:

- “Mean” / “Max” denotes a bag feature is obtained through a mean-pooling / max-pooling layer among a bag of instance features, which is inspired by the strong baseline of non-parametric MIL method [35].
- “Instance” denotes that K -means is directly performed over all the instance features in training set, since each instance can be regarded as a bag with length of one.

Then, interventional training is applied to baseline methods (including ABMIL and DSMIL) and we report the results in Fig. 4d. Notably, even with such simple aggregation

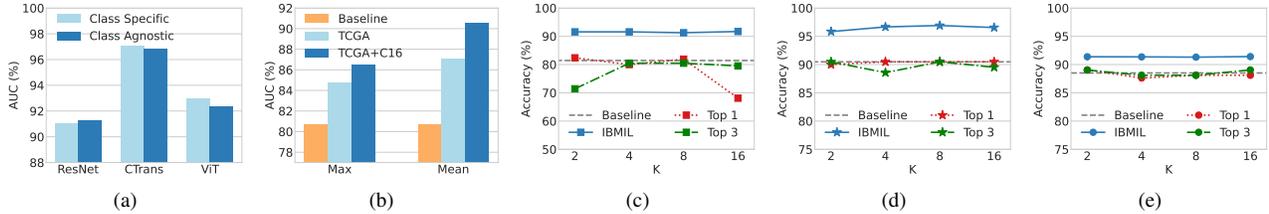


Figure 5. (a) Constitution of confounder set: class-specific vs. class-agnostic. (b) Confounder set from TCGA vs. TCGA+Camelyon16. (c)-(e) Intervention training vs. KNN classifier with the feature extractors of ResNet, CTransPath and ViT, respectively.

Table 2. Performance of non-parametric aggregators.

Aggregator	K	Precision	Recall	Accuracy	AUC
Max	1	77.59	77.35	70.95	82.23
	2	81.58	73.54	72.38	83.44
	4	78.34	79.17	76.67	84.71
	8	80.41	78.72	78.10	84.95
	16	82.32	71.06	70.48	83.79
Mean	1	77.13	71.60	71.43	80.68
	2	81.97	81.49	81.43	85.81
	4	82.04	81.56	81.43	87.10
	8	84.50	80.46	80.48	89.14
	16	85.13	78.16	78.10	89.00

strategy, IBMIL still outperforms the baseline, and remains competitive or even better compared to “Default” setting, which indicates stage 2 in our scheme is unnecessary. By omitting stage 2, our scheme can be more elegant without performance degradation in most cases.

Can IBMIL improve non-parametric baselines? Besides using non-parametric aggregators to generate bag features for confounder set, we further take them as baselines and verify whether IBMIL is also able to improve them (*i.e.*, max/mean-pooling). Surprisingly, in Tab. 2, IBMIL brings significant improvements under all settings, where the best performance is even comparable to these attention-based aggregators. It indicates that IBMIL is indeed compatible with all compared MIL methods, including the non-parametric ones.

Constituting confounder set w/ or w/o bag labels? Given bag labels, we explore the class-specific K -means. In particular, we apply K -means to each class respectively, preserving the intra-class variation and class-relevant characteristics. From Fig. 5a, we observe no obvious performance gap between class-specific and class-agnostic K -means. We conjecture that 1) the confounders could be independent of the class identity, and 2) bag features are already separable by classes. We will explore the way of incorporating bag labels in future work. On the other hand, the unsupervised fashion makes our scheme compatible with large-scale unlabelled datasets. We explore more unlabelled bags via combining the bags of TCGA and Camelyon16,

and constituting the confounder set via the non-parametric aggregators. From Fig. 5b, we observe a clear improvement on AUC under both max- and mean-poolings. That indicates, with more bags, our implementation can achieve better approximation of confounders.

Is IBMIL just post-processing? Since our proposed IBMIL shares some commonalities with the embedding-based MIL [35], one may ask: Do the improvements only come from exploring the bag relations? To answer this question, we make minor modifications on ABMIL. Instead of interventional training with confounders, we obtain the confounder dictionary via the class-specific K -means and treat it as a KNN classifier for evaluation. As can be seen, it brings limited improvements or even degrades the performance, verifying the improvements comes from interventional training, which is not just post-processing.

5. Conclusions

The vast majority of recent efforts in this field seek to enhance the feature extractor and aggregator. This paper addresses MIL from a novel perspective via analyzing the confounders between bags and labels. This leads to the proposed novel Interventional Bag Multi-Instance Learning (IBMIL), a new deconfounded bag-level prediction approach to suppress the bias caused by the bag contextual prior. IBMIL introduces a structure causal model to reveal the causalities and eliminates their effect through the back-door adjustment with practical implementations. Comprehensive experiments have been conducted on various MIL benchmarks and the results show that IBMIL can boost existing methods significantly. In future, we plan to approximate confounder set in a more efficient and elegant manner. As a general method to use causal intervention for bag-level prediction, IBMIL provides fresh insight into MIL problem.

Acknowledgments. Dr. Yi Xu was supported in part by NSFC 62171282, Shanghai Municipal Science and Technology Major Project (2021SHZDZX0102), 111 project BP0719010, and SJTU Science and Technology Innovation Special Fund ZH2018ZDA17 and YG2022QN037.

References

- [1] Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes Van Diest, Bram Van Ginneken, Nico Karssemeijer, Geert Litjens, Jeroen AWM Van Der Laak, Meyke Hermesen, 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. [2](#), [5](#)
- [2] Hanbo Chen, Xiao Han, Xinjuan Fan, et al. Rectified cross-entropy and upper transition loss for weakly supervised whole slide image classifier. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 351–359. Springer, 2019. [2](#)
- [3] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9640–9649, 2021. [5](#)
- [4] Zhang Chen, Zhiqiang Tian, Jihua Zhu, Ce Li, and Shaoyi Du. C-cam: Causal cam for weakly supervised semantic segmentation on medical image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11676–11685, 2022. [2](#)
- [5] Veronika Cheplygina, Lauge Sørensen, David MJ Tax, Marleen de Bruijne, and Marco Loog. Label stability in multiple instance learning. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 539–546. Springer, 2015. [2](#)
- [6] Philip Chikontwe, Meejeong Kim, Soo Jeong Nam, Heounjeong Go, and Sang Hyun Park. Multiple instance learning with center embeddings for histopathology classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 519–528. Springer, 2020. [2](#)
- [7] Nicolas Coudray, Paolo Santiago Ocampo, Theodore Sakellaropoulos, et al. Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning. *Nature Medicine*, 24(10):1559–1567, 2018. [1](#)
- [8] Nicolas Coudray, Paolo Santiago Ocampo, Theodore Sakellaropoulos, Navneet Narula, Matija Snuderl, David Fenyo, Andre L Moreira, Narges Razavian, and Aristotelis Tsirigos. Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning. *Nature Medicine*, 24(10):1559–1567, 2018. [2](#)
- [9] 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. [1](#)
- [10] 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. *arXiv preprint arXiv:2010.11929*, 2020. [5](#)
- [11] Alexander D’Amour. On multi-cause approaches to causal inference with unobserved confounding: Two cautionary failure cases and a promising alternative. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3478–3486. PMLR, 2019. [5](#)
- [12] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018. [6](#)
- [13] Metin N Gurcan, Laura E Boucheron, Ali Can, et al. Histopathological image analysis: A review. *IEEE Reviews in Biomedical Engineering*, 2:147–171, 2009. [1](#)
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. [5](#)
- [15] Le Hou, Dimitris Samaras, Tahsin M Kurc, et al. Patch-based convolutional neural network for whole slide tissue image classification. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2424–2433, 2016. [1](#), [2](#), [3](#)
- [16] Xinting Hu, Kaihua Tang, Chunyan Miao, Xian-Sheng Hua, and Hanwang Zhang. Distilling causal effect of data in class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3957–3966, 2021. [2](#)
- [17] Maximilian Ilse, Jakub Tomczak, and Max Welling. Attention-based deep multiple instance learning. In *International Conference on Machine Learning*, pages 2127–2136. PMLR, 2018. [2](#), [3](#), [5](#)
- [18] Jakob Nikolas Kather, Alexander T Pearson, Niels Halama, et al. Deep learning can predict microsatellite instability directly from histology in gastrointestinal cancer. *Nature Medicine*, 25(7):1054–1056, 2019. [1](#)
- [19] 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*, pages 14318–14328, 2021. [1](#), [2](#), [3](#), [5](#), [6](#)
- [20] Tiancheng Lin, Hongteng Xu, Canqian Yang, and Yi Xu. Interventional multi-instance learning with deconfounded instance-level prediction. *arXiv preprint arXiv:2204.09204*, 2022. [2](#), [4](#), [5](#)
- [21] Ming Y Lu, Richard J Chen, Jingwen Wang, Debora Dillon, and Faisal Mahmood. Semi-supervised histology classification using deep multiple instance learning and contrastive predictive coding. *arXiv preprint arXiv:1910.10825*, 2019. [2](#)
- [22] Jovana Mitrovic, Brian McWilliams, Jacob Walker, Lars Buesing, and Charles Blundell. Representation learning via invariant causal mechanisms. *arXiv preprint arXiv:2010.07922*, 2020. [2](#)
- [23] Judea Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688, 1995. [2](#)
- [24] Judea Pearl, Madelyn Glymour, and Nicholas P Jewell. *Causal inference in statistics: A primer*. John Wiley & Sons, 2016. [2](#), [3](#)
- [25] Judea Pearl and Dana Mackenzie. *The book of why: the new science of cause and effect*. Basic books, 2018. [3](#)

- [26] Ziniu Qian, Kailu Li, Maode Lai, Eric I Chang, Bingzheng Wei, Yubo Fan, Yan Xu, et al. Transformer based multiple instance learning for weakly supervised histopathology image segmentation. *arXiv preprint arXiv:2205.08878*, 2022. **2**
- [27] Alexander Rakhlin, Alexey Shvets, Vladimir Iglovikov, and Alexandr A Kalinin. Deep convolutional neural networks for breast cancer histology image analysis. In *International Conference Image Analysis and Recognition*, pages 737–744. Springer, 2018. **2**
- [28] 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, 2021. **1, 2, 5**
- [29] Yash Sharma, Aman Shrivastava, Lubaina Ehsan, Christopher A Moskaluk, Sana Syed, and Donald Brown. Cluster-to-conquer: A framework for end-to-end multi-instance learning for whole slide image classification. In *Medical Imaging with Deep Learning*, pages 682–698. PMLR, 2021. **4**
- [30] Kaihua Tang, Jianqiang Huang, and Hanwang Zhang. Long-tailed classification by keeping the good and removing the bad momentum causal effect. *arXiv preprint arXiv:2009.12991*, 2020. **2**
- [31] David Tellez, Geert Litjens, Jeroen van der Laak, and Francesco Ciompi. Neural image compression for gigapixel histopathology image analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. **1, 2**
- [32] Nenad Tomasev, Ioana Bica, Brian McWilliams, Lars Buesing, Razvan Pascanu, Charles Blundell, and Jovana Mitrovic. Pushing the limits of self-supervised resnets: Can we outperform supervised learning without labels on imagenet? *arXiv preprint arXiv:2201.05119*, 2022. **2**
- [33] Tan Wang, Jianqiang Huang, Hanwang Zhang, and Qianru Sun. Visual commonsense r-cnn. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10760–10770, 2020. **2, 4, 7**
- [34] Xi Wang, Hao Chen, Caixia Gan, Huangjing Lin, Qi Dou, Efstratios Tsougenis, Qitao Huang, Muyan Cai, and Pheng-Ann Heng. Weakly supervised deep learning for whole slide lung cancer image analysis. *IEEE Transactions on Cybernetics*, 50(9):3950–3962, 2019. **2**
- [35] Xinggang Wang, Yongluan Yan, Peng Tang, Xiang Bai, and Wenyu Liu. Revisiting multiple instance neural networks. *Pattern Recognition*, 74:15–24, 2018. **2, 5, 7, 8**
- [36] Xiyue Wang, Sen Yang, Jun Zhang, Minghui Wang, Jing Zhang, Wei Yang, Junzhou Huang, and Xiao Han. Transformer-based unsupervised contrastive learning for histopathological image classification. *Medical Image Analysis*, 81:102559, 2022. **5**
- [37] Gang Xu, Zhigang Song, et al. Camel: A weakly supervised learning framework for histopathology image segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10682–10691, 2019. **1**
- [38] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057. PMLR, 2015. **4**
- [39] Jiawei Yang, Hanbo Chen, Yu Zhao, Fan Yang, Yao Zhang, Lei He, and Jianhua Yao. Remix: A general and efficient framework for multiple instance learning based whole slide image classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 35–45. Springer, 2022. **5**
- [40] Jiawen Yao, Xinliang Zhu, and Junzhou Huang. Deep multi-instance learning for survival prediction from whole slide images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 496–504. Springer, 2019. **1, 2**
- [41] Jiawen Yao, Xinliang Zhu, Jitendra Jonnagaddala, Nicholas Hawkins, and Junzhou Huang. Whole slide images based cancer survival prediction using attention guided deep multiple instance learning networks. *Medical Image Analysis*, 65:101789, 2020. **1, 2**
- [42] Zhongqi Yue, Hanwang Zhang, Qianru Sun, and Xian-Sheng Hua. Interventional few-shot learning. *arXiv preprint arXiv:2009.13000*, 2020. **2**
- [43] Dong Zhang, Hanwang Zhang, Jinhui Tang, Xiansheng Hua, and Qianru Sun. Causal intervention for weakly-supervised semantic segmentation. *arXiv preprint arXiv:2009.12547*, 2020. **2, 4**
- [44] Hongrui Zhang, Yanda Meng, Yitian Zhao, Yihong Qiao, Xiaoyun Yang, Sarah E Coupland, and Yalin Zheng. Dtdfmil: 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*, pages 18802–18812, 2022. **1, 2, 5, 7**
- [45] Weijia Zhang. Stable multi-instance learning via causal inference. 2019. **2**
- [46] Yu Zhao, Fan Yang, et al. Predicting lymph node metastasis using histopathological images based on multiple instance learning with deep graph convolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4837–4846, 2020. **1, 2**
- [47] Niyun Zhou, De Cai, Xiao Han, and Jianhua Yao. Enhanced cycle-consistent generative adversarial network for color normalization of h&e stained images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 694–702. Springer, 2019. **5**
- [48] Xinliang Zhu, Jiawen Yao, Feiyun Zhu, and Junzhou Huang. Wsisa: Making survival prediction from whole slide histopathological images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7234–7242, 2017. **1**