

HAMIL: Hierarchical Aggregation-Based Multi-Instance Learning for Microscopy Image Classification

Yanlun Tu, Houchao Lei, Wei Long, and Yang Yang

Abstract—Multi-instance learning is common for computer vision tasks, especially in biomedical image processing. Traditional methods for multi-instance learning focus on designing feature aggregation methods and multi-instance classifiers, where the aggregation operation is performed either in feature extraction or learning phase. As deep neural networks (DNNs) achieve great success in image processing via automatic feature learning, certain feature aggregation mechanisms need to be incorporated into common DNN architecture for multi-instance learning. Moreover, flexibility and reliability are crucial considerations to deal with varying quality and number of instances.

In this study, we propose a hierarchical aggregation network for multi-instance learning, called HAMIL. The hierarchical aggregation protocol enables feature fusion in a defined order, and the simple convolutional aggregation units lead to an efficient and flexible architecture. We assess the model performance on two microscopy image classification tasks, namely protein subcellular localization using immunofluorescence images and gene annotation using spatial gene expression images. The experimental results show that HAMIL outperforms the state-of-the-art feature aggregation methods and the existing models for addressing these two tasks. The visualization analyses also demonstrate the ability of HAMIL to focus on high-quality instances.

Index Terms—Multi-instance learning, biomedical image processing, hierarchical aggregation.

1 INTRODUCTION

A lot of computer vision tasks exhibit multi-instance property, i.e. each input data sample is represented by a bag of instances and learning is performed at the bag level instead of the instance level [1], [2]. For example, video-based identity recognition can be treated as a multi-instance learning (MIL) task, as each video input is a time-series of frames and each frame can be regarded as an instance [3]. Another example is the automatic diagnosis system based on magnetic resonance imaging (MRI) data, which consists of multiple slices presenting the area of interest being scanned [4]. In these two cases, there is either temporal or spatial dependency between images; whereas in a lot more MIL scenarios, instances within each bag are unordered or exhibit no dependency among each other. Such unordered MIL tasks are especially common in biomedical image processing [5], [6], [7], [8].

The imaging of biological samples, plays a key role in current life science research, enabling scientists to analyze from tissue level down to the cell level. Unlike natural images, biological microscopy images are often much harder to obtain due to the difficulty in preparing the specimen or stringent experimental conditions. Therefore, during the imaging procedure, it is common to capture multiple images for a specimen in a single trial of experiments and perform multiple trials for repeatability. To infer the functions of genes or characteristics of molecules, all the captured images should be considered comprehensively to give

- Y. Tu, H. Lei, W. Long, Y. Yang are with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai, SH 200240, China. E-mail: {tuyanlun, leihouchao}@sjtu.edu.cn, yilang2007lw@gmail.com, yangyang@cs.sjtu.edu.cn.

Manuscript received April 19, 2005; revised August 26, 2015.
(Corresponding author: Y. Yang).

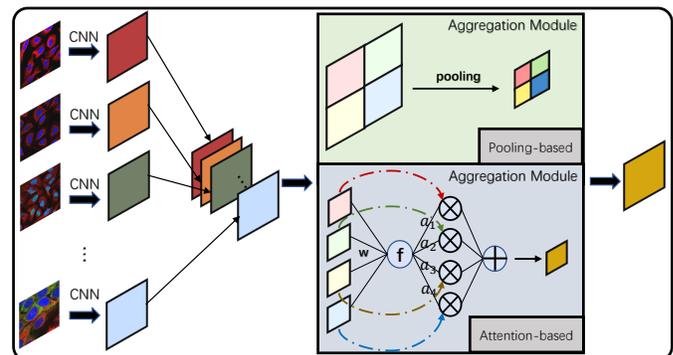


Fig. 1. Two typical aggregation mechanisms in deep learning models. CNNs are commonly used to extract features from input images, and then the extracted features are fed into the feature aggregation module to obtain the final aggregation output.

a more accurate judgement on the final output, as single images may only contain partial information.

Compared with video or 3D data, dealing with unordered input requires the model to have the permutation-invariant property [7], i.e., the model should be insensitive to the presenting order of input instances. Moreover, the models also need to address the issues coming from variable-sized inputs and different image qualities. Thus, developing MIL models is a very challenging job.

Traditional MIL methods fall into two categories, focusing on feature extraction and classification, respectively.

The first category of methods aggregates multiple instances into a fixed-length input before feeding the classifiers. The aggregation operation can be performed before or after feature

extraction. FlyIT [9] is an example of the former type, which first stitches images belonging to the same bag into large images and then extracts features from the large images. The latter type first extracts features for raw instances and then performs the aggregation operation (as illustrated in Fig. 1). For instance, MIMT-CNN [10] employs 10 VGG models to learn feature representation from 10 images, respectively, and concatenates them into a bag representation for the following classification. The feature aggregation can be regarded as a part of feature engineering, and the subsequent learning procedure remains unchanged. Alternatively, common classifiers can be modified to handle inputs of multiple instances. Till now, a lot of MIL methods have been developed, like multi-instance KNN [11], multi-instance neural networks [12], [13], multi-instance decision tree [14], and multi-instance support vector machines (SVMs) [15]. The core idea is to compute pair-wise similarity or loss function at the bag level instead of instance level and define the loss function according to the prediction accuracy on bags.

Note that before using these traditional learning models, image features should be extracted separately, while deep neural networks perform feature learning automatically. For the past decade, MIL tasks have also benefitted a lot from deep learning [16]. However, the performance of the existing multi-instance DNN models has been limited due to the high complexity of input data. For one thing, some methods set a fix-length of input [9], [10], which lacks flexibility of handing varying-sized input. For the other thing, since the number of instances per sample is often large while high-quality instances are few, the methods that are unable to identify informative instances may not work well.

To address these challenges, we propose a model, called HAMIL (Hierarchically Aggregation for Multi-Instance Learning). The model is featured by a hierarchical aggregation protocol with simple but effective aggregation units. It can not only learn input bags with varying size but also give preference to informative instances. We assess the performance of HAMIL on two microscopy image classification tasks. On both tasks, HAMIL achieves significant improvement in prediction accuracy compared with other feature aggregation methods and state-of-the-art predictors for these two tasks.

2 RELATED WORK

2.1 Traditional feature aggregation

In traditional image processing, feature aggregation aims to fuse features extracted from multiple images into a comprehensive feature representation before feeding into a learning model. There are three typical feature fusion methods, the bag of visual words (BOV) [17], Fisher vector [18], and vector of locally aggregated descriptors (VLAD) [19]. BOV regards image features as words, builds a vocabulary of local image features and generates a vector of their occurrence counts. The Fisher vector method stores the mixing coefficients of the Gaussian mixture model (GMM) as well as the mean and covariance deviation vectors of the individual components. The VLAD method computes the distance of each feature point to the cluster center closest to it. All of these three methods produce a fixed-length feature vector for the input image set, which can work with traditional machine learning models, like SVMs.

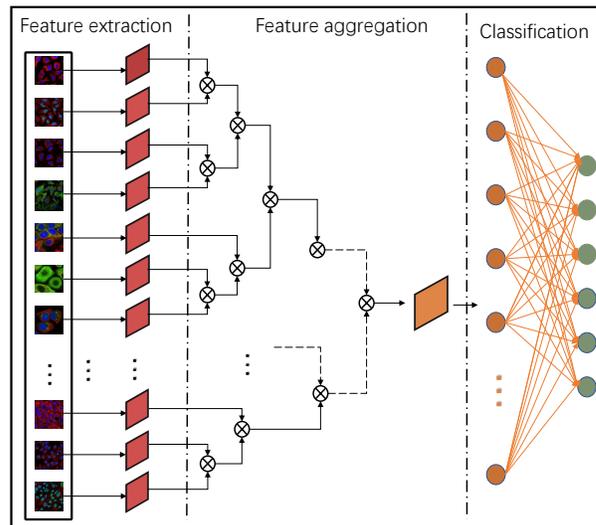


Fig. 2. Model architecture of HAMIL. The raw images are first fed to convolutional layers to get feature embeddings and then aggregated based on the hierarchical tree to get a unified encoding for final classification.

2.2 Aggregation of deep convolutional features

For the past decade, feature encoding via deep neural networks has almost replaced traditional feature engineering in image processing [20] [21] [22]. Various feature aggregation methods working with deep neural networks have also emerged. There are two typical aggregation schemes, i.e. pooling-based and attention-based.

As a straightforward aggregation method, the pooling function is permutation-invariant and requires no (or very few) parameters, thus has been widely used in MIL tasks. Babenko *et al.* showed that a simple sum-pooling method working on deep convolutional features is superior to traditional aggregation methods working on shallow features [23]; Su *et al.* proposed the multi-view CNN (MVCNN) model, which adopts max pooling to combine information of photos captured from multiple orientations for 3D recognition [16]; Wang *et al.* introduced three popular pooling functions in their proposed mi-net, i.e. max pooling, mean pooling, and log-sum-exp (LSE) pooling (a smooth version and convex approximation of the max pooling) [24]; and Kraus *et al.* proposed the Noisy-AND pooling function to handle outliers [8].

A major flaw of pooling is its incapability to focus on important instances, while the attention mechanism is a good choice, as it can assign scores/weights to instances. Till now, various attention-based feature aggregation methods have been developed. Yang *et al.* proposed a neural aggregation network for video face recognition, NAN [3], which aggregates multiple features by an attention module. Ilse *et al.* proposed a general attention-based deep MIL framework, which introduced both attention and gated-attention functions [7]. Due to their good interpretability, attention-based methods have gained more popularity in biomedical image processing, e.g. Yao *et al.* used the attention-based MIL pooling for whole slide image analysis. In addition, the attention-based MIL can be implemented in various types of neural networks. For instance, Annofly uses long short-term memory model as the backbone network [25], while Set Transformer [26] and ImPLoc [27] are based on Transformer models. Besides these two mechanisms, there are some other deep MIL models with special

designs, like DeepMIML [28] utilizing a pluggable sub-concept layer to learn the latent relationship between input patterns and output semantic labels. And Tu *et al.* proposed a graph neural network-based model, which regards each instance as a node in the graph. Learning the graph embedding is essentially an information aggregation process [29].

The proposed HAMIL model has a different aggregation mechanism compared with the existing methods. It has trainable and non-linear convolution aggregation units, and designs a hierarchical clustering protocol for permutation-invariance. Different from the hierarchical graph clustering in GNN-based model [29], both the numbers of clusters and aggregation times are automatically determined by the size of input bag rather than fixed hyperparameters.

3 METHODS

3.1 Problem description

In this study, we discuss a more complex MIL problem, namely multi-instance multi-label learning. Formally, let \mathcal{X} denote the sample set, i.e. $\mathcal{X} = \{X_i\}$, where $i \in \{1, 2, \dots, n\}$, n is the number of samples, and X_i is a sample; $X_i = \{x_{i,1}, x_{i,2}, \dots, x_{i,m}\}$, where m is the number of instances of the i th sample, $x_{i,j}$ ($j \in \{1, 2, \dots, m\}$) denotes an instance of X_i . Let $\mathcal{Y} = \{Y_i\}$ be the output space, and $Y_i = \{y_1, y_2, \dots, y_k\}$ corresponds to the label set of X_i . The goal is to learn a mapping function $f : \mathcal{X} \mapsto \mathcal{Y}$. Especially, as we focus on the image processing tasks, here each instance is an image and each sample is represented by a bag of images.

3.2 Model architecture

To process input bags of any size and varying sizes, HAMIL sets no limit on the number of instances for a bag and implements a hierarchical aggregation protocol. The model architecture is shown in Figure 2. There are three main components, namely feature extraction, feature aggregation, and classification. The first component consists of several convolutional layers serving as the feature extractor. The CNN layers are followed by a hierarchical aggregation procedure to produce a unified representation for the input set of instances, which is further fed to the fully connected layer for classification.

Specifically, by regarding the input bag as a set of images, HAMIL constructs a hierarchy of images/instances within each bag, thus determines an aggregation order according to the tree structure, i.e. from the leaf node to the root node of the hierarchical tree. The construction of instance hierarchy is described in Algorithm 1. In Line 11, the distance between two clusters is defined as the minimum distance of all cross-cluster instance pairs, i.e. the single-link method. And the distance between instances is Euclidean distance computed based on features yielded by the convolutional layers in the first component of HAMIL.

The hierarchy of instances is a binary tree in nature. During the construction of the hierarchy, we record the merging history of clusters. Assume that there are a total of K merging steps (note that K is not a hyperparameter but the height of the binary tree determined by the number of instances), then we will keep a queue \mathcal{T} of length K . Each element in the array is a triplet denoted by t_k , i.e.,

$$t_k = \langle I_{C_1}, I_{C_2}, I_C \rangle, k \in \{1, \dots, K\} \quad (1)$$

Algorithm 1 Construction of the hierarchy of instances

Input: A bag of images, X .
 $X = \{x_1, x_2, \dots, x_m\}$, where m is the number of instances of X .

Output: A queue recording the aggregation history.

- 1: Let \mathcal{S} be a set of clusters.
- 2: $\mathcal{S} = \{\{x_1\}, \{x_2\}, \dots, \{x_m\}\}$.
- 3: $I_{\{x_i\}} = i$, $maxI = m$, where I denotes the index of clusters and $maxI$ is the maximum index of clusters in \mathcal{S} .
- 4: Let \mathcal{T} be a queue of triplets as defined in Eq. (1). Initialize \mathcal{T} to an empty queue.
- 5: **while** $|\mathcal{S}| > 1$ **do**
- 6: $n = |\mathcal{S}|$.
- 7: Let $minL$ be the minimal pairwise distance between clusters. $minL = +\infty$.
- 8: Let $C1, C2$ be the two clusters to be merged.
- 9: **for** $i = 1$ to $n - 1$ **do**
- 10: **for** $j = i + 1$ to n **do**
- 11: $L_{i,j} = dis(S_i, S_j)$, where dis is a distance function and S_i is the i th element in \mathcal{S} by the ascending order of indexes.
- 12: **if** $L_{i,j} < minL$ **then**
- 13: $minL = L_{i,j}$.
- 14: $C1 = S_i, C2 = S_j$.
- 15: **end if**
- 16: **end for**
- 17: **end for**
- 18: $C = C1 \cup C2$
- 19: $S = (S \cup C) \setminus C1 \setminus C2$.
- 20: $maxI = maxI + 1, I_C = maxI$.
- 21: Define a triplet $t = \langle I_{C1}, I_{C2}, I_C \rangle$.
- 22: Push t into \mathcal{T} .
- 23: **end while**
- 24: Return \mathcal{T} .

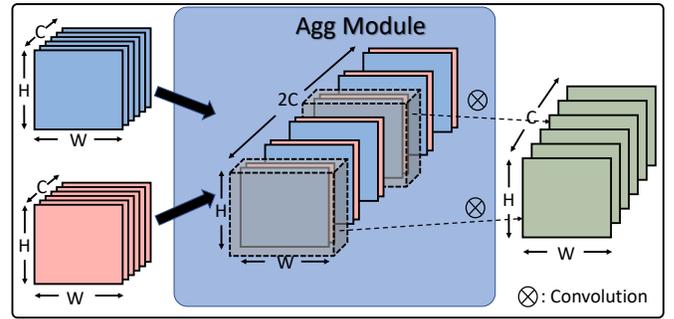


Fig. 3. Illustration of feature aggregation in HAMIL.

t_k consists of three indexes of clusters in \mathcal{S} . The first two are the indexes of the two clusters being merged in the k -th step, and the last one is the index of the newly generated cluster. Then, the aggregation order is determined by the records in \mathcal{T} , and the specific aggregation operations are defined by kernel functions as described in the next Section.

3.3 Feature aggregation unit

As the convolution operation can be regarded as a weighted average of inputs, we design the aggregation units via convolution.

In detail, given two input feature maps of size $H \times W \times C$, which can be considered as C pairs of $H \times W$ matrices along the channels, i.e. for each channel there is a pair of matrices. We first input each pair into the feature aggregation unit to obtain the aggregated output of size $H \times W$. Then, the C outputs are concatenated into the final output of size $H \times W \times C$. Figure 3 shows an illustration of the aggregation units.

Formally, let x_1 and x_2 be the feature maps to be aggregated. The aggregation is defined in Eq. (2),

$$\begin{aligned} \mathbf{X} &= [x_1, x_2], \\ \mathbf{O} &= \mathbf{W} * \mathbf{X} + b, \end{aligned} \quad (2)$$

where \mathbf{X} is a tensor composed by the feature maps, $*$ is the convolution operator, \mathbf{W} is the convolutional filter, b is the bias, and \mathbf{O} is the aggregated feature map. This is a one-layer aggregation unit, We call it L1Agg. To meet the needs of more complex tasks, the aggregation units can be extended to deeper versions with multiple layers. Eqs. (3) and (4) formulate the two-layer and three-layer aggregation units which nest the basic aggregation operation twice and three times, respectively.

$$\mathbf{O} = \mathbf{W}' * f(g(\mathbf{W} * \mathbf{X} + b)) + b', \quad (3)$$

where $g(\cdot)$ is a normalization function, $f(\cdot)$ denotes the ReLU function.

$$\mathbf{O} = \mathbf{W}'' * f(g(\mathbf{W}' * f(g(\mathbf{W} * \mathbf{X} + b)) + b')) + b'', \quad (4)$$

No matter how many layers are there, the aggregation units are shared by all aggregation operations. Thus, the aggregation module introduces only a slight increase in the number of parameters compared to the backbone CNN model.

4 EXPERIMENTS

We assess the model performance on two large-scale image classification tasks, involving two typical kinds of bio-images, i.e. microscopic cell images and gene expression images. The two tasks are described below.

Task I: Prediction of protein subcellular location using immunofluorescence (IF) images. Each protein corresponds to a bag of microscopy images captured from multiple tissues. The labels, i.e. cellular locations, are predicted based on all localization patterns implied in these images. A protein may exist in multiple locations.

Task II: Gene function annotation using gene expression images. The *in situ* hybridization (ISH) imaging technology visualizes spatial distribution patterns of gene expression in tissues and help to reveal gene functions. Each gene corresponds to a bag of expression images captured in different angles or experimental trials. The labels are functional annotation terms. A gene may have more than one annotation terms.

Apparently, both of these two tasks are multi-instance multi-label classification. We compare HAMIL with both single-instance learning models and the existing feature aggregation models as listed in the following. To assess the performance of HAMIL, we compare it with 7 baseline models, including:

- A single-instance learning model (SI)¹;

1. It has the same backbone network as HAMIL but its inputs are single images, which have the same labels as the bags they belong to. The prediction results of single images are combined per bag to yield the bag labels.

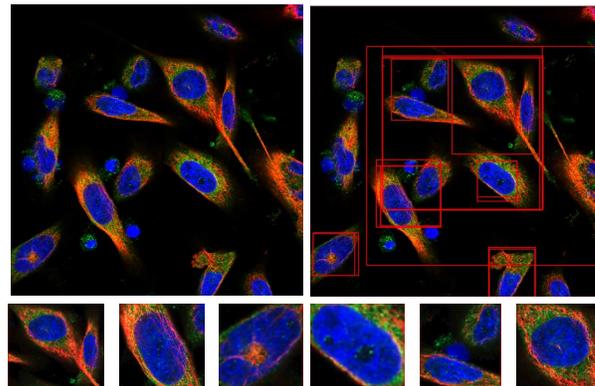


Fig. 4. Extracted patches from HPA raw images. The top left figure is a raw image, the top right figure shows the results of selective search, and the bottom figures are the cropped and scaled patches.

- Three pooling-based methods, MI with mean pooling, MVCNN [16] with max pooling, and SPoC [23] with sum pooling;
- Two attention-based methods, NAN [3] and Attention [7];
- A specially designed deep MIL model, DeepMIML [28].

All the baseline models have the same backbone network (ResNet18) for extracting features from raw images. The prediction performance is evaluated by three metrics, AUC (the Area Under ROC Curve)², macro F_1 , and micro F_1 .

4.1 Task 1- Protein microscopy images

4.1.1 Data source

The data set was collected from the cell atlas of the human protein atlas (HPA) database [30]. We download 1600 bags of IF images, each of which corresponds to a protein. The location annotations of proteins are hand-crafted in the database. The total number of images is 19,777. The sizes of raw images include 800×800 , 1728×1728 , and 2048×2048 . Obviously, the large size of images and small numbers of bags bring great computation difficulties. Considering that there are usually multiple cells within an image, the selective search algorithm [31] is adopted to pick the cellular patches from raw images as shown in Figure 4. We scale the selected patches into a fixed size 512×512 . Proteins with too many patches are split into multiple bags. The top-10 frequent classes are selected for our experiment. The training and test sets are partitioned at the protein level. Details about the data set can be found in suppl. Table S1.

4.1.2 Experiment details

Adam optimizer [32] is adopted to optimize the network. The learning rate is set to 0.0001, and the parameters β_1 and β_2 of the Adam optimizer are set to 0.9 and the default value, respectively. The batch size is 8 and the loss function is standard BCE loss. We select the L1Agg function as our aggregation function, and the kernel size is set to 7, in addition, we use the BN layer after the convolution operation to normalize the feature map. The selected feature extraction network is ResNet18 [33], and the final output layer consists of 10 nodes (to predict 10 major cellular organelles).

2. The averaged AUC value over labels, i.e. macro AUC

TABLE 1
Performance comparison on Task 1*

Method		AUC	MacroF ₁	MicroF ₁
SI		0.938	0.572	0.734
Pool	MI	0.927	0.577	0.745
	MVCNN [16]	0.938	0.675	0.763
	SPoC [23]	0.941	0.714	0.770
Attn	NAN [3]	0.928	0.705	0.784
	Attention [7]	0.910	0.559	0.729
	Gated-Attention [7]	0.900	0.522	0.724
DeepMIML [28]		0.915	0.662	0.746
HAMIL (ours)		0.944	0.733	0.784

*‘Pool’ denotes pooling-based aggregation and ‘Attn’ denotes attention-based aggregation. Gated-Attention used gated-attention mechanism for aggregation. The result of DeepMIML is obtained by resizing the images to 224×224 , which leads to much higher accuracy than using the original size. All methods use ResNet18 as feature extractor.

4.1.3 Results

The prediction accuracy of HAMIL and 7 baseline models on Task I is shown in Table 1. As can be seen, HAMIL achieves the best performance. SI treats the task as single-instance learning, which assumes that every single instance has the same label set as the whole bag, thus would introduce a lot of noisy samples.

The three pooling methods, MI, MVCNN and SPoC, adopt mean, max, and sum pooling, respectively. SPoC achieves the best performance, perhaps due to the ‘centering prior’ incorporated in the sum pooling, which assigns larger weights to center area of feature maps. As preprocessed by selective search (Fig. 4), the central regions of most images contain more cell information.

As for the attention-based methods, NAN performs much better than ‘attention’ and ‘gated-attention’. Despite different model settings and equations for computing attention scores, the major difference is that we apply NAN on feature maps while the latter two methods on feature vectors. Following the implementation in [7], the feature maps yielded by CNN layers are turned into lower-dimensional feature vectors via fully connected layers, thus the attention and gated-attention mechanisms work on feature vectors. This FC transformation may result in information loss.

DeepMIML does not obtain a comparable result to HAMIL, perhaps because it was designed to identify multiple concepts in single-image input, which is another kind of MIL problem different from ours.

4.2 Task 2 - Gene expression images

4.2.1 Data source

The data set consists of standard gene expression images of *Drosophila* embryos from the FlyExpress database [34]. All the images were extracted from the Berkeley *Drosophila* Genome Project (BDGP) (www.fruitfly.org) [35], [36] and preprocessed to a uniform size 180×320 . The data set is in bag-level. Each bag of images presents the expression distribution on the *Drosophila* embryo for a single gene, in which the images were captured from different orientations, i.e. dorsal, lateral and ventral. We use the same dataset as FlyIT [9], where the dataset is divided (at bag-level) into training, validation, and test sets according to the ratio of 4:1:5 (suppl. Table S1 shows data statistics). The total number of labels is 10, each of which is an ontology term describing anatomical and developmental properties. This task is also a multi-

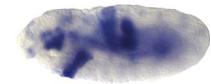
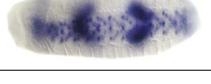
Gene	exec	fkf
Image		
		
CV Term	brain primordium salivary gland primordium hindgut proper primordium	brain primordium salivary gland primordium hindgut proper primordium

Fig. 5. Examples of the gene annotation data of *Drosophila* embryos. The first row shows gene names, the second row shows gene expression images for each gene, and the last row shows the corresponding CV terms of the gene, i.e. labels.

label classification problem, as each gene may have multiple developmental terms. Figure 5 shows some example images.

4.2.2 Experiment details

The experimental settings in Task 2 are almost the same as in Task 1, except that the batch size is set to 32, because the gene expression images are smaller than protein microscopy images, so we can add more train samples during one train process. Here we compare HAMIL with additional four state-of-the-art MIL methods proposed for the *Drosophila* embryo image analysis [5], [9], [25]. E-MIMLSVM is an MIL algorithm that adapts the kernel function of support vector machines to the MIL scenario [5]. FlyIT [9] adopts image stitching to combine raw images. AnnoFly [25] first extracts features via ResNet, then employs RNN to deal with the multi-instance input.

4.2.3 Results

The experimental results are shown in Table 2. In this task, the image quality is much higher than that of Task I and the number of instances per bag is fewer. That could explain why pooling-based methods achieve better results on this Task. By contrast, Task I has much more instances per bag and highly varying image quality (see examples in Fig. 7), the attention mechanism which identifies high-quality instances can help make better decision based on the important instances.

Interestingly, the pooling and attention-based methods have comparable and even better performance than the previously proposed MIL methods for this task. E-MIMLSVM is a shallow learning model. Annofly uses the pre-trained (on Imagenet) CNN to extract feature embeddings and then feeds them into an LSTM model. FlyIt conducts aggregation at the raw-image level. Neither of them fully exploit the feature learning ability of deep CNNs during the aggregation process.

In summary, the advantages of HAMIL over these state-of-the-art methods can be attributed to the end-to-end training of ResNet18 for feature extraction and the hierarchical aggregation operations on the feature maps.

4.2.4 Visualization analysis

We visualize the features representations with and without aggregation by projecting them onto a 2D space via tSNE algorithm [37]. Figure 6 (a) shows the 2D distribution of features learned by ResNet without aggregation, and Figure 6 (b) shows the 2D distribution of aggregated features. Different colors denote

TABLE 2
Performance comparison on Task 2*

Method		AUC	MacroF ₁	MicroF ₁
SI		0.936	0.691	0.685
Pool	MI	0.938	0.688	0.706
	MVCNN [16]	0.942	0.741	0.743
	SPoC [23]	0.944	0.729	0.737
Attn	NAN [3]	0.927	0.726	0.734
	Attention [7]	0.935	0.692	0.707
	Gated-Attention [7]	0.935	0.701	0.701
DeepMIML [28]		0.922	0.681	0.700
E-MIMLSVM [5]		0.846	0.598	0.640
AnnoFly [25]		0.937	0.702	0.713
FlyIT [9]		0.936	0.718	0.711
HAMIL (ours)		0.944	0.755	0.746

*The results of E-MIMLSVM⁺, AnnoFly, and FlyIT are from [9].

different classes. The features of different categories before aggregation have a small margin and large overlap, which brings great difficulties to subsequent classifiers. After aggregation, samples of different categories are much more separated, with small overlap, which is beneficial for subsequent classification.

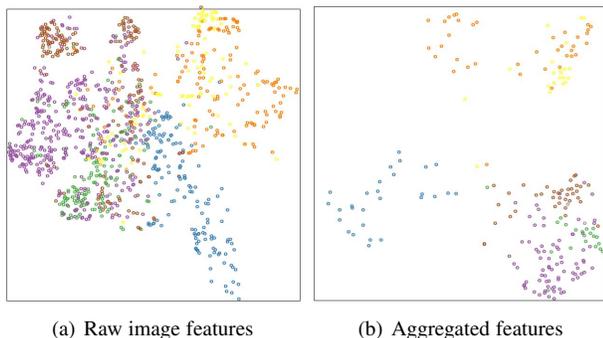


Fig. 6. Visualization of feature representations before and after aggregation.

4.3 Ablation Study and Analysis

4.3.1 Ablation of the hierarchical aggregation

In order to verify the contribution of hierarchical design to the performance of HAMIL, we remove the hierarchical clustering in HAMIL and aggregate the instances in random order using the same convolutional aggregation functions, which is denoted by RAMIL. As RAMIL is not permutation-invariant, during its training process, we randomly shuffle the order of instances within bags in each epoch. The comparison results of HAMIL and RAMIL are shown in Table 3. HAMIL outperforms RAMIL on all the three metrics, suggesting that the hierarchical aggregation scheme is able to improve the performance of the model.

In addition, to get more insight on the aggregation operation, we compute the cosine similarity between aggregated features and the features of input images, and take the similarity value as the score of the input images, which can be regarded as a kind of preference that model assigns to the instances. Figure 7 visualizes the images with different scores. Each row shows images from the same bag, the two numbers below each image represent the weights assigned to the image by HAMIL and RAMIL, respectively. As can be seen, the images with high

TABLE 3
Performance comparison on Task 2*

Task	Method	AUC	MacroF ₁	MicroF ₁
Task I	RAMIL	0.935	0.720	0.778
	HAMIL-A	0.939	0.695	0.773
	HAMIL	0.944	0.733	0.784
Task II	RAMIL	0.935	0.735	0.730
	HAMIL-A	0.932	0.736	0.732
	HAMIL	0.944	0.755	0.746

*RAMIL and HAMIL are two variants of HAMIL. RAMIL performs aggregation on random-ordered instances but with the same aggregation units as HAMIL; HAMIL-A adopts the same hierarchical aggregation protocol but replaces the convolutional aggregation units with mean pooling.

scores contain relatively obvious local patterns corresponding to the target labels, suggesting that the feature aggregation retains more information from the high-quality images, i.e. the features aggregated by HAMIL will automatically focus on important features. By contrast, it can be found that RAMIL pays nearly equal attention to different quality images, as the scores differ slightly. This is because RAMIL aggregates each input in a random order, i.e., each time it randomly selects two feature embeddings to aggregate, thus there is no mechanism for attention allocation. As a result, when the quality of images varies a lot, the performance of RAMIL cannot be guaranteed.

4.3.2 Ablation of the convolutional aggregation units

To verify the effectiveness of convolutional aggregation units in our design, we use a simple average operation to replace the non-linear convolutional aggregation unit. HAMIL-A is used to denote this alteration of aggregation unit. In Table 3 we compare it with HAMIL in Task I and II. Because convolutional aggregation unit is a trainable module, its performance is better than simple average operation. Suppl. Table S4 also shows the results on localized content-based image retrieval.

5 COMPARISON ON AGGREGATION MECHANISMS

HAMIL is a new aggregation framework designed for addressing MIL in image processing. It allows variable sizes of input bags without a restriction on the bag size, and it has a light-weight architecture with simple aggregation operations. The hierarchical aggregation protocol makes the model invariant to the order of instances in bags. Actually, many aggregation methods also have these advantages, like pooling and attention-based [7], [24]. A comparative analysis of these three mechanisms is listed below.

i) Traditional pooling-based methods (max, sum or mean) are non-trainable, thus have limited learning ability; while both attention-based aggregation function and HAMIL are trainable.

ii) The existing attention-based aggregation methods assign weight scores to instances, i.e. they treat an image/instance as a whole, no matter how the pairwise instance similarity is computed (based on feature maps or feature embedding vectors). By contrast, both pooling-based and HAMIL can operate on local regions of multiple instances. Although attention methods have good interpretability (the weights directly reflect the importance of instances), the weighted-sum aggregation function is incompetent to express complex aggregation mechanisms.

iii) According to the experimental results, the pooling-based methods are very efficient and effective on data sets with high quality, and attention-based methods have more advantages in

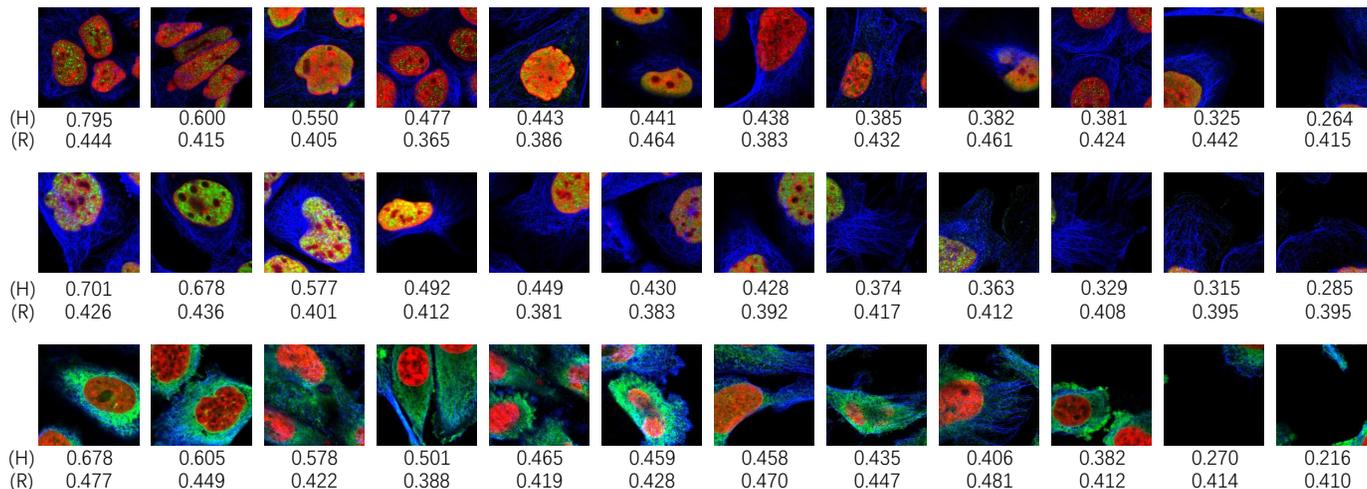


Fig. 7. Visualization of scores for image instances. The three rows denote three bags of images, and the numeric values are the scores assigned by HAMIL and RAMIL to the corresponding images.

handling data with much noise or varying quality; while HAMIL is applicable to both cases, due to the powerful learning ability of its aggregation units.

Although HAMIL is designed for handling raw image data, we perform an additional experiment to investigate its performance on traditional MIL data sets (as shown in suppl. Table S2). As the input features are provided (preprocessed via feature engineering), we remove the first component of HAMIL, i.e. the CNN layers for learning image features. Following the practice in [7], We run the experiments using 10-fold cross-validations for 5 times (each time the data set is randomly partitioned into 10 folds). Experimental details can be found in suppl. Section B.

Suppl. Table S4 shows the results on three localized content-based image retrieval tasks, and suppl. Table S5 shows the results for two drug activity data sets. Considering that the input consists of statistical features of images, we not only replace the original 2D conv to 1D conv but also add HAMIL-A (as mentioned in Section 4.3.2) into comparison. As can be seen, HAMIL achieves a little advantage on Fox and Elephant data sets while no advantage on others. Besides, HAMIL-A and HAMIL have very close performance.

From these experimental results, it can be seen that as the aggregation units of HAMIL are convolution-based and they directly function on feature maps, HAMIL is more suitable for processing raw image data. Nevertheless, the proposed HAMIL model could be an alternative feature aggregation protocol applied to the data sets with traditional features, which obtains comparable results with the existing aggregation methods. The contributions of HAMIL can be summarized as follows.

i) The convolution-based aggregation function of HAMIL is able to capture local neighborhood information in images during the aggregation process and learn complex non-linear aggregation patterns.

ii) The hierarchical design makes the model insensitive to the order of instances and improves model accuracy.

iii) The end-to-end architecture further improves the feature learning, as the hierarchical clustering can be dynamically adjusted to yield better clustering of instances based on the updated feature representations.

6 CONCLUSION

We propose a hierarchical aggregation network for multi-instance learning in image classification. Different from the mainstream pooling and attention-based methods, the proposed HAMIL utilizes convolution operations for aggregation on feature maps and a hierarchical architecture to ensure permutation-invariance. HAMIL achieves better performance than the existing aggregation methods on two microscopy image classification tasks. Moreover, as a general feature aggregation network, HAMIL can be easily applied to other MIL image processing tasks.

APPENDIX A DETAILS FOR TASKS I AND II

TABLE S1
Data Overview of Tasks I and II

Dataset	Protein microscopy image			Gene expression image		
	Train	Val	Test	Train	Val	Test
Bag #	5496	1320	6901	2714	678	3393
Instance #	69943	16687	86964	10237	2426	12872

The datasets can be accessed at <https://bioimagestore.blob.core.windows.net/dataset?restype=container&comp=list>

APPENDIX B DETAILS FOR CLASSICAL MIL DATA SETS

TABLE S2
Overview of five classical MIL data sets

Dataset	# of bags	# of instances	# of features
Musk1	92	476	166
Musk2	102	6598	166
Tiger	200	1220	230
Fox	200	1302	230
Elephant	200	1391	230

TABLE S3
Model architecture of HAMIL and HAMIL-A working on the five traditional data sets*

Layer	Type
1	fc-256+ReLU
2	dropout
3	fc-128+ReLU
4	dropout
5	fc-64+ReLU
6	dropout
7	hierarchical clustering
8	conv/average aggregation
9	fc-1+sigm

*To compare with previous studies, Layers 1-6 and 9 have exactly the same settings as Ilse et al's attention-based method [7].

TABLE S4
Performance comparison of classification accuracy (mean±std) on localized content-based image retrieval¹

Method	Fox	Tiger	Elephant
mi-SVM [38]	0.582	0.784	0.822
MI-SVM [38]	0.578	0.840	0.843
MI-Kernel [39]	0.603	0.842	0.843
EM-DD [40]	0.609±0.101	0.730±0.096	0.771±0.097
mi-Graph [41]	0.620±0.098	0.860±0.083	0.869±0.078
miVLAD [42]	0.620±0.098	0.811±0.087	0.850±0.080
miFV [42]	0.621±0.109	0.813±0.083	0.852±0.081
mi-Net [24]	0.613±0.078	0.824±0.076	0.858±0.083
MI-Net [24]	0.622±0.084	0.830 ±0.072	0.862±0.077
MI-Net with DS [24]	0.630±0.080	0.845 ±0.087	0.872±0.072
MI-Net with RC [24]	0.619±0.104	0.836 ±0.083	0.857±0.089
Attention [7]	0.615±0.043	0.839 ±0.022	0.868±0.022
Gated-Attention [7]	0.603±0.029	0.845 ±0.018	0.857±0.027
HAMIL-A ²	0.631±0.117	0.806 ±0.095	0.878±0.091
HAMIL ³	0.647±0.105	0.815 ±0.091	0.865±0.079

¹Results of the first 11 methods are from [24]; and results of the Attention and Gate-Attention methods are from [7].

² Replace the conv-based aggregation of HAMIL with mean pooling.

³ Change the aggregation units of HAMIL from 2D conv to 1D conv.

TABLE S5
Performance comparison of classification accuracy (mean±std) on drug activity prediction*

Method	MUSK1	MUSK2
mi-SVM [38]	0.874	0.836
MI-SVM [38]	0.779	0.843
MI-Kernel [39]	0.880	0.893
EM-DD [40]	0.849±0.098	0.869±0.108
mi-Graph [41]	0.889±0.073	0.903±0.086
miVLAD [42]	0.871±0.098	0.872±0.095
miFV [42]	0.909±0.089	0.884±0.094
mi-Net [24]	0.889±0.088	0.858±0.110
MI-Net [24]	0.887±0.091	0.859±0.102
MI-Net with DS [24]	0.894±0.093	0.874±0.097
MI-Net with RC [24]	0.898±0.097	0.873±0.098
Attention [7]	0.892±0.040	0.858±0.048
Gated-Attention [7]	0.900±0.050	0.863±0.042
HAMIL-A	0.892±0.120	0.857±0.102
HAMIL	0.866±0.121	0.820±0.107

*Results of the first 11 methods are from Wang et al.(2018) [24]; and results of the Attention and Gate-Attention methods are from Ilse et al. (2018) [7].

The experimental settings are almost the same for the five data sets. We use the SGD optimizer, the momentum is 0.9, and the

weight decay is 0.005. We set the learning rate to 0.0001. Our models are trained for 100 epochs on Fox, MUSK1, and MUSK2. As for Elephant and Tiger, the number of training epochs is 50, due to their faster convergence.

REFERENCES

- [1] Z.-H. Zhou, "Multi-instance learning: A survey," *Department of Computer Science & Technology, Nanjing University, Tech. Rep.*, 2004. 1
- [2] J. Foulds and E. Frank, "A review of multi-instance learning assumptions," *The Knowledge Engineering Review*, vol. 25, no. 1, pp. 1–25, 2010. 1
- [3] D. F. H. J. Yang, P. Ren and G. Hua., "Neural aggregation network for video face recognition." *arxiv preprint arxiv:1603.05474*, 2016. 1, 2, 4, 5, 6
- [4] M. Liu, J. Zhang, E. Adeli, and D. Shen, "Landmark-based deep multi-instance learning for brain disease diagnosis," *Medical Image Analysis*, vol. 43, pp. 157 – 168, 2018. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1361841517301524> 1
- [5] Y.-X. Li, S. Ji, S. Kumar, J. Ye, and Z.-H. Zhou, "Drosophila gene expression pattern annotation through multi-instance multi-label learning," *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, vol. 9, no. 1, pp. 98–112, 2012. 1, 5, 6
- [6] "Whole slide images based cancer survival prediction using attention guided deep multiple instance learning networks," *Medical Image Analysis*, vol. 65, p. 101789, 2020. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1361841520301535> 1
- [7] M. Ilse, J. M. Tomczak, and M. Welling, "Attention-based deep multiple instance learning," *arXiv preprint arXiv:1802.04712*, 2018. 1, 2, 4, 5, 6, 7, 8
- [8] O. Z. Kraus, J. L. Ba, and B. J. Frey, "Classifying and segmenting microscopy images with deep multiple instance learning," *Bioinformatics*, vol. 32, no. 12, pp. i52–i59, 06 2016. [Online]. Available: <https://doi.org/10.1093/bioinformatics/btw252> 1, 2
- [9] W. Long, T. Li, Y. Yang, and H.-B. Shen, "Flyit: Drosophila embryogenesis image annotation based on image tiling and convolutional neural networks," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 2019. 2, 5, 6
- [10] T. Zeng and S. Ji, "Deep convolutional neural networks for multi-instance multi-task learning," in *Data Mining (ICDM), 2015 IEEE International Conference on*. IEEE, 2015, pp. 579–588. 2
- [11] J. Wang and J.-D. Zucker, "Solving multiple-instance problem: A lazy learning approach," 2000. 2
- [12] M.-L. Zhang and Z.-H. Zhou, "Improve multi-instance neural networks through feature selection," *Neural Processing Letters*, vol. 19, no. 1, pp. 1–10, 2004. 2
- [13] Z.-H. Zhou and M.-L. Zhang, "Neural networks for multi-instance learning," in *Proceedings of the International Conference on Intelligent Information Technology, Beijing, China*, 2002, pp. 455–459. 2
- [14] Y. Chevaleyre and J.-D. Zucker, "Solving multiple-instance and multiple-part learning problems with decision trees and rule sets. application to the mutagenesis problem," in *Conference of the Canadian Society for Computational Studies of Intelligence*. Springer, 2001, pp. 204–214. 2
- [15] S. Andrews, I. Tschantaridis, and T. Hofmann, "Support vector machines for multiple-instance learning," in *Advances in neural information processing systems*, 2003, pp. 577–584. 2
- [16] E. K. Hang Su, Subhransu Maji and E. Learned-Miller, "Multi-view convolutional neural networks for 3d shape recognition," *In ICCV*, 2015. 2, 4, 5, 6
- [17] J. Yang, Y.-G. Jiang, A. Hauptmann, and C.-W. Ngo, "Evaluating bag-of-visual-words representations in scene classification," 01 2007, pp. 197–206. 2
- [18] J. S. F. Perronnin, Y. Liu and H. Poirier, "Large-scale image retrieval with compressed fisher vectors." *In The Twenty-Third IEEE Conference on Computer vision and Pattern Recognition, CVPR*, pp. 3384–3391, 2010. 2
- [19] C. H. Jegou, M. Douze and P. Perez, "Aggregating local descriptors into a compact image representation." *In The Twenty-Third IEEE Conference on Computer vision and Pattern Recognition, CVPR*, pp. 3304–3311, 2010. 2
- [20] J. A. S. Razavin, H. Azizpour and S. Carlsson, "Cnn features off-the-shelf: An astounding baseline for recognition," *In DeepVision workshop*, 2014. 2
- [21] I. M. Oquab, L. Bottou and J. Sivic, "Learning and transferring mid-level image representatons using convolutional neural networks," *In IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, pp. 1717–1724, 2014. 2

- [22] A. V. K. Chatfield, K. Simonyan and A. Zisserman, "Return of the devil in the details: Delving deep into convolutional nets," *In British Machine Vision Conference, MNVC*, 2014. [2](#)
- [23] A. Babenko and V. Lempitsky, "Aggregating deep convolutional features for image retrieval," *In ICCV*, 2015. [2](#), [4](#), [5](#), [6](#)
- [24] X. Wang, Y. Yan, P. Tang, X. Bai, and W. Liu, "Revisiting multiple instance neural networks," *Pattern Recognition*, vol. 74, pp. 15–24, 2018. [2](#), [6](#), [8](#)
- [25] Y. Yang, M. Zhou, Q. Fang, and H.-B. Shen, "Annofly: annotating drosophila embryonic images based on an attention-enhanced rnn model," *Bioinformatics (Oxford, England)*, vol. 35, pp. 2834–2842, 08 2019. [2](#), [5](#), [6](#)
- [26] J. Lee, Y. Lee, J. Kim, A. Kosiorek, S. Choi, and Y. W. Teh, "Set transformer: A framework for attention-based permutation-invariant neural networks," ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. Long Beach, California, USA: PMLR, 09–15 Jun 2019, pp. 3744–3753. [Online]. Available: <http://proceedings.mlr.press/v97/lee19d.html> [2](#)
- [27] W. Long, Y. Yang, and H.-B. Shen, "Imploc: a multi-instance deep learning model for the prediction of protein subcellular localization based on immunohistochemistry images," *Bioinformatics*, vol. 36, no. 7, pp. 2244–2250, 2020. [2](#)
- [28] J. Feng and Z.-H. Zhou, "Deep miml network." in *AAAI*, 2017, pp. 1884–1890. [3](#), [4](#), [5](#), [6](#)
- [29] M. Tu, J. Huang, X. He, and B. Zhou, "Multiple instance learning with graph neural networks," *arXiv preprint arXiv:1906.04881*, 2019. [3](#)
- [30] M. Uhlen, P. Oksvold, L. Fagerberg, E. Lundberg, K. Jonasson, M. Forsberg, M. Zwahlen, C. Kampf, K. Wester, S. Hober *et al.*, "Towards a knowledge-based human protein atlas," *Nature Biotechnology*, vol. 28, no. 12, pp. 1248–1250, 2010. [4](#)
- [31] J. Uijlings, K. E. Sande, T. Gevers, and A. W. M. Smeulders, "Selective search for object recognition," *International Journal of Computer Vision*, vol. 104, no. 2, pp. 154–171, 2013. [4](#)
- [32] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *Arxiv preprint arxiv:1312.6117*, 2013. [4](#)
- [33] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," pp. 770–778, 2016. [4](#)
- [34] S. Kumar, C. Konikoff, B. Van Emden, C. Busick, K. T. Davis, S. Ji, L. W. Wu, H. Ramos, T. Brody, S. Panchanathan *et al.*, "Flyexpress: visual mining of spatiotemporal patterns for genes and publications in drosophila embryogenesis," *Bioinformatics*, vol. 27, no. 23, pp. 3319–3320, 2011. [5](#)
- [35] P. Tomancak, A. Beaton, R. Weiszmann, E. Kwan, S. Shu, S. E. Lewis, S. Richards, M. Ashburner, V. Hartenstein, S. E. Celniker *et al.*, "Systematic determination of patterns of gene expression during drosophila embryogenesis," *Genome Biology*, vol. 3, no. 12, pp. 1–14, 2002. [5](#)
- [36] P. Tomancak, B. P. Berman, A. Beaton, R. Weiszmann, E. Kwan, V. Hartenstein, S. E. Celniker, and G. M. Rubin, "Global analysis of patterns of gene expression during drosophila embryogenesis," *Genome Biology*, vol. 8, no. 7, pp. 1–24, 2007. [5](#)
- [37] L. V. Der Maaten and G. E. Hinton, "Visualizing data using t-sne," *Journal of Machine Learning Research*, vol. 9, pp. 2579–2605, 2008. [5](#)
- [38] S. Andrews, I. Tsochantaridis, and T. Hofmann, "Support vector machines for multiple-instance learning," *Advances in Neural Information Processing Systems*, vol. 15, pp. 561–568, 01 2002. [8](#)
- [39] T. Gärtner, P. A. Flach, A. Kowalczyk, and A. J. Smola, "Multi-instance kernels," in *Proceedings of the Nineteenth International Conference on Machine Learning*, ser. ICML '02. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2002, p. 179–186. [8](#)
- [40] Q. Zhang and S. Goldman, "Em-dd: An improved multiple-instance learning technique." 01 2001, pp. 1073–1080. [8](#)
- [41] Z.-H. Zhou, Y.-Y. Sun, and Y.-F. Li, "Multi-instance learning by treating instances as non-i.i.d. samples," 2009. [8](#)
- [42] X. S. Wei, J. Wu, and Z. H. Zhou, "Scalable algorithms for multi-instance learning," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 4, pp. 975–987, 2017. [8](#)