

Self-Paced Learning for Images of Antinuclear Antibodies

Yiyang Jiang^{ID}, Student Member, IEEE, Guangwu Qian^{ID}, Jiaxin Wu^{ID}, Qi Huang, Qing Li^{ID}, Fellow, IEEE, Yongkang Wu and Xiao-Yong Wei^{ID}, Senior Member, IEEE

arXiv:2511.21519v1 [cs.CV] 26 Nov 2025

Abstract— Antinuclear antibody (ANA) testing is a critical method for diagnosing autoimmune disorders such as Lupus, Sjögren's syndrome, and scleroderma. Despite its importance, manual ANA detection is slow, labor-intensive, and demands years of training. ANA detection is complicated by over 100 coexisting antibody types, resulting in vast fluorescent pattern combinations. Although machine learning and deep learning have enabled automation, ANA detection in real-world clinical settings presents unique challenges as it involves multi-instance, multi-label (MIML) learning. In this paper, a novel framework for ANA detection is proposed that handles the complexities of MIML tasks using unaltered microscope images without manual preprocessing. Inspired by human labeling logic, it identifies consistent ANA sub-regions and assigns aggregated labels accordingly. These steps are implemented using three task-specific components: an instance sampler, a probabilistic pseudo-label dispatcher, and self-paced weight learning rate coefficients. The instance sampler suppresses low-confidence instances by modeling pattern confidence, while the dispatcher adaptively assigns labels based on instance distinguishability. Self-paced learning adjusts training according to empirical label observations. Our framework overcomes limitations of traditional MIML methods and supports end-to-end optimization. Extensive experiments on one ANA dataset and three public medical MIML benchmarks demonstrate the superiority of our framework. On the ANA dataset, our model achieves up to +7.0% F1-Macro and +12.6% mAP gains over the best prior method, setting new state-of-the-art results. It also ranks top-2 across all key metrics on public datasets, reducing Hamming loss and one-error by up to 18.2% and 26.9%, respectively. The source code can be accessed at <https://github.com/fletcherjiang/ANA-SelfPacedLearning>.

Index Terms— Antinuclear antibodies, multi-instance learning, multi-label learning, self-paced learning

This research was supported by the Hong Kong Research Grants Council via the General Research Fund (project no. PolyU 15200023) and the National Natural Science Foundation of China (Grant No.: 62372314). The experimental part of this work was supported by The Centre for Large AI Models (CLAIM) of The Hong Kong Polytechnic University.

Yiyang Jiang, Jiaxin Wu, Qing Li, and Xiao-Yong Wei are with the PolySmart Group, Department of Computing, The Hong Kong Polytechnic University, Hong Kong, China (e-mail: yiyang.jiang@connect.polyu.hk; nikki-jiaxin.wu@polyu.edu.hk; qing-prof.li@polyu.edu.hk; cs007.wei@polyu.edu.hk).

Yiyang Jiang, Guangwu Qian, Qi Huang, and Xiao-Yong Wei are with the College of Computer Science, Sichuan University, Sichuan, China (e-mail: cswei@scu.edu.cn).

Yongkang Wu is with the Department of Laboratory Medicine and Outpatient, West China Hospital, Sichuan University, Sichuan, China (e-mail: vipywk@163.com).

Yiyang Jiang and Guangwu Qian contributed equally to this work.

Xiao-Yong Wei is the corresponding author.

I. INTRODUCTION

ANTINUCLEAR antibody (ANA) test is a widely used diagnostic method for a variety of autoimmune disorders, including systemic lupus erythematosus, Sjögren syndrome, scleroderma, mixed connective tissue disease, polymyositis, dermatomyositis, autoimmune hepatitis, and drug-induced lupus. These disorders can progress to life-threatening conditions, such as cancer, or lead to death, making ANA detection a critical and longstanding focus in the medical field. In this work, a *pattern* denotes the characteristic morphology and subcellular distribution of fluorescence on HEp-2 cells, specifying which compartments (e.g., nucleus, nucleolus, centromeres, cytoplasm) exhibit signal and in what configuration (e.g., homogeneous, speckled, nucleolar), as determined by autoantibody binding to specific antigens. It does not refer to periodic texture in the sense of computer vision.

In the ANA test, patient serum is incubated on HEp-2 cells, and a fluorescent secondary antibody is applied to reveal nuclear binding. The resulting staining pattern and titer indicate ANA positivity, guiding interpretation. Recently, ANA detection has also drawn significant attention from the medical image recognition community, spurred in part by efforts from ICPR (International Conference on Pattern Recognition) since 2012. ICPR introduced a series of ANA datasets between 2012 and 2016 as benchmarks for ANA detection challenges. These datasets typically consist of manually selected and labeled images, each containing either a single cell or multiple cells with the same pattern (as shown in Fig. 1(a) and Fig. 1(b)). Consequently, most prior studies have formulated ANA detection as a multi-class classification problem for either single-instance or multi-instance images [1]–[26]. However, this approach is impractical in real-world clinical or laboratory settings. Technicians typically sample representative sub-regions from microscope images for ANA labeling. These sub-regions often contain hundreds of cells, potentially with different ANA patterns co-occurring (e.g., AC-1 and AC-3 in Fig. 1(c)). The reasoning behind using sub-regions for detection is that antibodies are produced in response to antigens, whose heterogeneous distribution across cells leads to uneven antibody binding; analyzing sub-regions captures this variability. To ensure reliable detection, technicians examine sub-regions to confirm whether antibodies are consistent across most cells, rather than occurring by chance in isolated instances.

Detecting antinuclear antibodies (ANAs) is a challenging

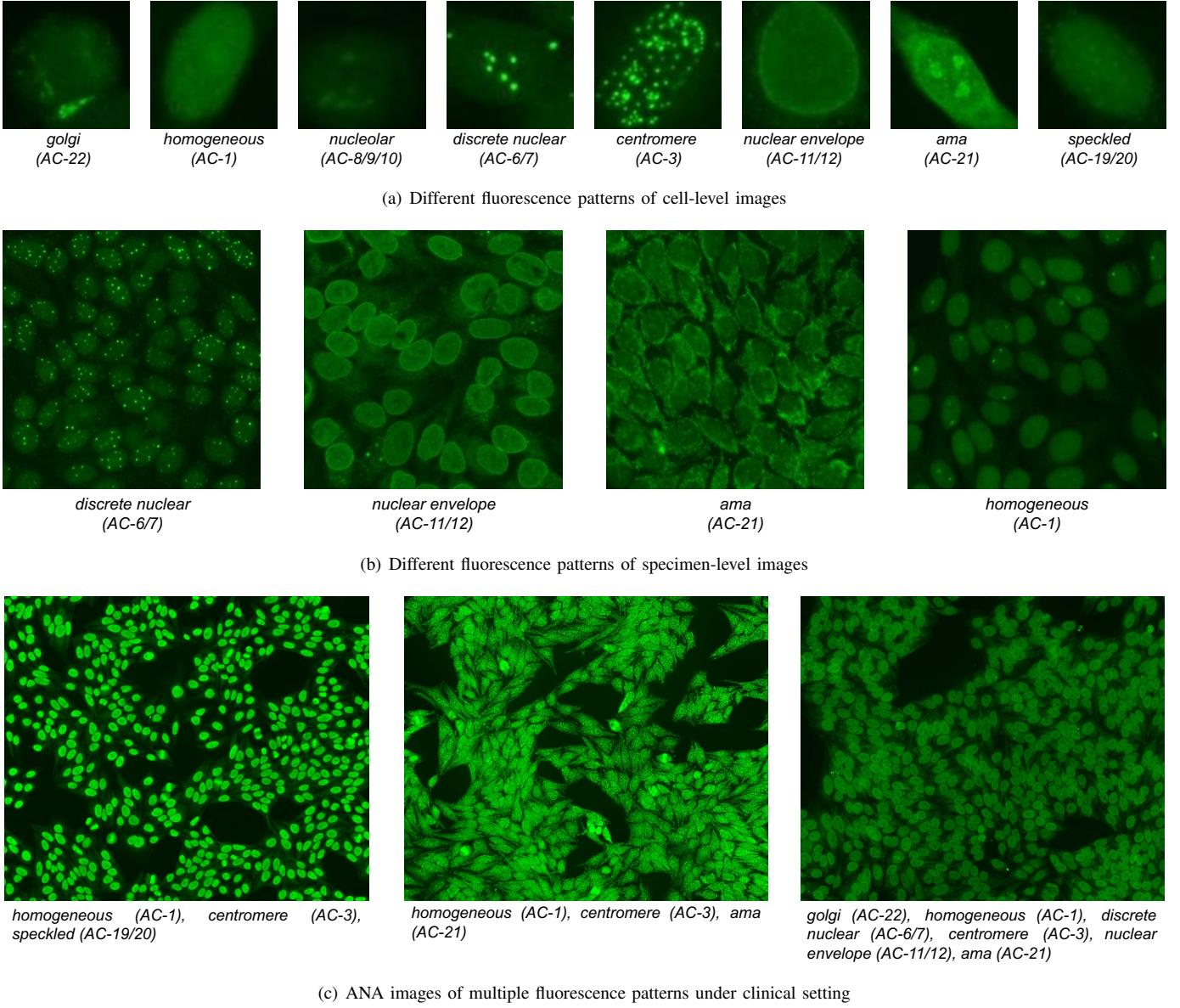


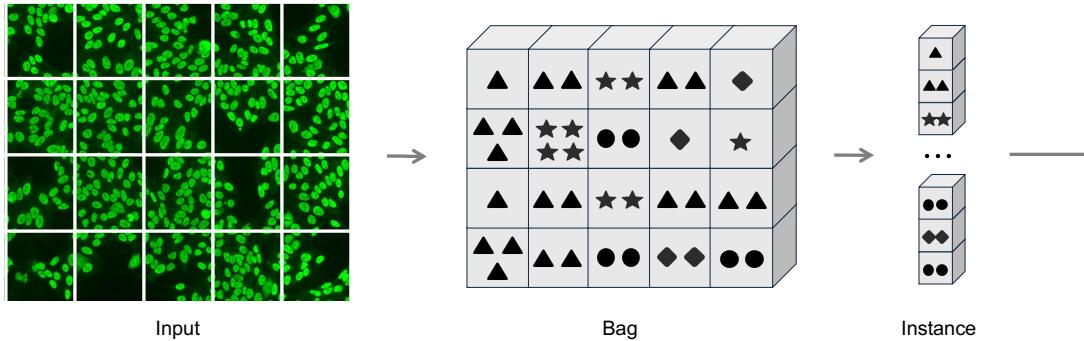
Fig. 1. Examples of ANA images and visualization of the multi-label multi-instance challenge.

task, even for experienced technicians. The difficulty arises from the vast number of antibody types and their tendency to coexist within cells rather than appearing in isolation. This coexistence means that the fluorescent pattern observed is determined by the combination of ANA types. Given roughly 100 ANA pattern types, the space of multi-label assignments is $2^{100} (\approx 10^{30})$, because each pattern is independently present or absent. This combinatorial effect further complicates pattern detection. Several international efforts have been made to standardize and define patterns, such as the recommendations by the European Autoimmunity Standardization Initiative [27] and the International Consensus on ANA Patterns (ICAP) [28]. Despite these efforts, only 29 patterns (labeled AC-1 to AC-29) have been officially recognized and agreed upon by ICAP, which is a small subset of the theoretical pattern space. This discrepancy makes detection prone to distractions caused by the many unrecognized patterns. As a result, it takes years

of rigorous training for a technician to become proficient in detecting and interpreting ANA patterns accurately.

From a machine learning perspective, ANA detection can be framed as a multi-instance multi-label (MIML) problem [29], as ANA samples consist of multiple sub-regions within an image, each associated with mixed labels. This makes it one of the most challenging learning problems. In prior research, multi-instance learning [30], [31] and multi-label learning [32] have typically been studied independently. However, the challenges in MIML stem not only from the combination of these two learning paradigms but also from the incompatibility of ANA data with the common assumptions underlying multi-instance or multi-label learning methods. First, as illustrated in Fig. 3(b), most multi-instance learning approaches are designed for instances that share a consistent label, often framed as a binary classification task (e.g., distinguishing foreground from background) [33]–[35]. The presence of

A Instance Extraction



B Self-Paced Learning

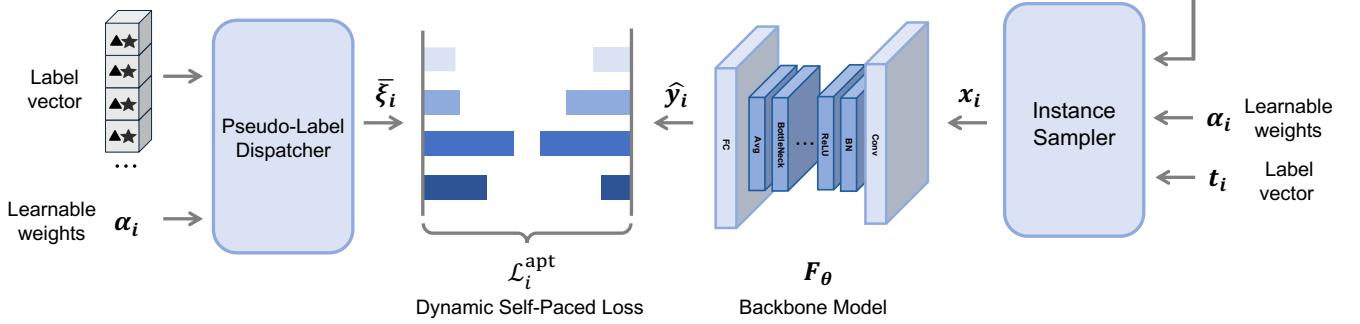


Fig. 2. Overview of the proposed architecture, the framework begins by extracting instances from an input ANA image, different symbols (\blacktriangle , \bullet , \blacklozenge) represent different ANA patterns detected in sub-regions of the ANA image. Each cell in the grid represents a local patch containing one or more ANA patterns. These instances are then evaluated by an Instance Sampler, which leverages learnable confidence weights to select the most representative instances for training. A CNN model provides initial predictions on the chosen instances, and a Pseudo-Label Dispatcher generates soft, continuous pseudo-labels based on both the image-level labels and the learnable confidence weights. The Dynamic Self-Paced Loss uses these pseudo-labels to adaptively emphasize more reliable instances and guide the network's learning process. Finally, all subregion predictions are aggregated for a new image to produce the final multi-label ANA classification result.

mixed instance (sub-region) labels in ANA images renders these traditional multi-instance models unsuitable. Additionally, multi-label learning is typically applied to image datasets (e.g., Core15K/16K [36], NUS-WIDE [37], Scene [38], IMDB [39], MS-COCO [40], WSS4LUAD [41]) where each label is often associated with a few distinct instances. Importantly, these instances typically represent natural objects (e.g., cats, cars) with clear, distinctive appearances. In contrast, ANA sub-regions exhibit arbitrary shapes and highly non-distinctive appearances, making classification far more complex. The difficulty of this task, as highlighted by the international standardization efforts led by ICAP, further underscores the challenges inherent to ANA detection, even for human experts.

In this paper, we present a pilot study that explores the use of multi-instance multi-label (MIML) learning for ANA recognition in a real-world clinical and laboratory setting. Unlike previous studies, the samples we use are unaltered microscope images without any manual preprocessing, such as cropping or resampling. As shown in Fig. 2, we propose a prototype framework inspired by the two-step logic employed by human technicians when labeling: 1) Identify sub-regions with consistent ANA presence (and disregard

the rest). 2) Assign labels to the identified sub-regions and aggregate these labels at the sample level. The first step is implemented using an **instance sampler**, which dynamically models and updates the learner's confidence regarding the presence of ANAs across training epochs. This confidence is used to suppress the influence of less confident instances during both the input and back-propagation stages, leveraging **instance-level attention** and a **dynamic self-paced loss** collaboratively. The second step is implemented through a **probabilistic pseudo-label dispatcher** and **label-aware learning rate coefficients**. The dispatcher integrates instance-level distinguishability, adaptively distributing the aggregated sample-level labels to individual instances. Meanwhile, the label-aware coefficients regulate the learning process by incorporating empirical observations about label distributions and ensuring a controlled and adaptive training procedure. The framework's key advantage lies in addressing MIML challenges in ANA detection through the use of self-paced learning. All parameters are automatically optimized in an end-to-end fashion.

II. RELATED WORK

A. ANA Detection

Autoimmune disorders are diseases that happen when the immune system loses its ability to tolerate the body's own tissues and causes damage to organs and tissues through autoantibodies or sensitized lymphocytes. Antinuclear antibodies (ANAs) are among the most common autoantibodies in these disorders, making their detection necessary for diagnosis. In manual detection, doctors examine ANA images under a fluorescence microscope and identify their patterns. This method is slow, requires much effort, and the results often depend on the doctor's experience. To address these issues, researchers have developed automatic ANA detection methods, which have gained significant attention. Based on datasets, previous studies can be categorized into cell-level [1], [2], [9] and specimen-level ANA detection [42], [43].

Most methods focus on cell-level detection, which involves analyzing single-cell images. In the early stages of machine learning, researchers extracted handcrafted features, combined them, and used machine learning classifiers [1], [2]. For example, Ghosh et al. [44] converted raw images to grayscale, extracted features like shape, texture, and gradient histograms, and used Support Vector Machines (SVM) [45] for classification. Thibault et al. extracted morphological and statistical texture features and classified them using logistic regression, random forest, and neural networks [2]. With the rise of deep learning, many Convolutional Neural Network (CNN)-based methods have been proposed for single-cell ANA detection [6], [9]. Gao developed a 3-layer CNN, which outperformed methods relying on handcrafted features [11]. Building on ResNet [46], DSRN further improved performance by incorporating both feature maps and final outputs into the model [7], [47].

Specimen-level ANA detection requires extra annotations to create a segmentation mask for each cell, which converts the task from specimen-level to cell-level. After predicting results for individual cells, an aggregation algorithm is used to combine these predictions into a specimen-level result. For example, Foggia et al. proposed a method that identifies the most common ANA type among cells and uses it as the specimen-level prediction [48]. Li et al. proposed a different approach in which all cell-level predictions were used to create a pattern histogram, which was then fed into an SVM to predict the specimen-level label [12]. However, these methods, which depend on extra annotations, are not suitable for clinical practice where only sample-level labels are available.

B. Multi-Instance Multi-Label Learning

Most multi-instance learning methods are designed for binary or multi-class classification. However, ANA detection in clinical practice is a typical multi-instance multi-label (MIML) task, which requires specialized algorithms. The differences between supervised learning, multi-instance learning, multi-label learning, and MIML are shown in Fig. 3. Among these, MIML is the most challenging because it combines the complexities of both multi-instance learning and multi-label learning.

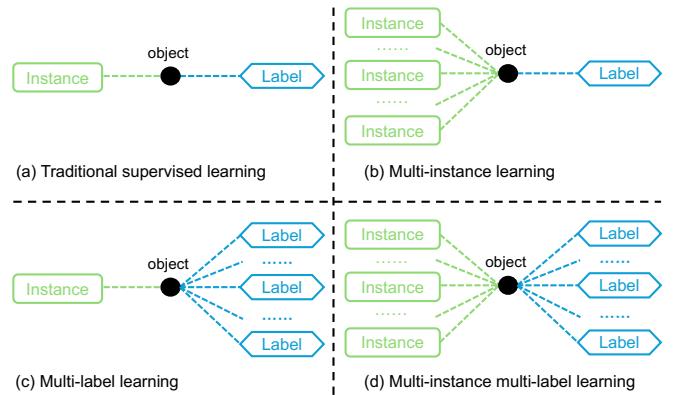


Fig. 3. Four different learning frameworks

Zhou et al. proposed two approaches to solve MIML tasks by converting them into either multi-instance learning or multi-label learning tasks [29]. Based on these approaches, he developed MIMLBoost and MIMLSVM to handle the tasks effectively. Inspired by advancements in deep learning, Feng introduced the DeepMIML network, which eliminates the need for pre-generating instances [49]. Instead, it automatically creates instance representations using a deep neural network structure. DeepMIML also includes a sub-concept layer designed to uncover hidden connections between input patterns and output semantic labels.

C. Self-Paced Learning

In 2010, Kumar et al. introduced Self-Paced Learning (SPL) [50] as an extension of curriculum learning (CL), which Bengio proposed in 2009 [51]. SPL is based on the idea of mimicking how humans learn: starting with simple examples and gradually moving to more complex ones [52]. Specifically, SPL addresses dataset issues such as class imbalance and corrupted labels by adding a pace-controlled regularization term to the learning objective. A minimal form is:

$$\min_{\theta, \mathbf{w} \in [0,1]^n} \sum_{i=1}^n \left[w_i \mathcal{L}(y_i, \hat{y}_i) + r(w_i, c) \right], \quad (1)$$

where \hat{y}_i denotes the model prediction for x_i , \mathcal{L} is any supervised loss, and $r(\cdot; c)$ is the self-paced regularizer with pace $c > 0$. Here, $w_i \in [0, 1]$ is the self-paced weight assigned to sample i and is optimized jointly with θ . Intuitively, w_i increases for “easy” samples (low loss) and decreases for hard/noisy ones, with the pace c gradually relaxing the selection. Our proposed framework is inspired by self-paced learning, but we have developed task-specific components to enhance the performance. There are several variations of the SPL learning framework, but we can only briefly cover those designed for multi-instance learning or multi-label learning, as our focus is on MIML. Zhang et al. applied self-paced learning to multi-instance learning and replaced binary weights with real-valued ones, which were directly calculated using the loss function [53], [54]. Li et al. extended self-paced learning to multi-label learning by considering the relationships between labels [55]. Instead of assigning weights based

on loss values, Zhong introduced a new reweighting function that considers instance complexity, determined by the distance between instance features and their corresponding labels [56]. Ren et al. proposed meta-learning algorithms guided by a small amount of unbiased validation data, assigning weights to training examples based on gradient directions [57], [58]. Our framework, in contrast, does not explicitly calculate weights or rely on unbiased validation data. Instead, it leverages prior knowledge about multi-label reliability and introduces label-level weights, which are used by three enhanced components to build an end-to-end deep model.

III. METHODOLOGY

A. Problem Formulation

We formulate the ANA task as a multi-instance multi-label (MIML) learning problem, and an overview of the proposed framework is shown in Fig. 2. Given a set of unaltered microscope ANA images and their ground truth labels $\{(x_i, y_i)\}_{i=1}^n$, the goal is to learn a function that takes an image x_i as input and predicts \hat{y}_i to best match y_i . Let $L = \{l_1, l_2, \dots, l_K\}$ be the set of predefined ANA pattern labels (e.g., Golgi, Nucleolar), where each l_k denotes a specific pattern.

To address the complexity of image-level prediction in ANA classification, we adopt a divide-and-aggregate strategy. Specifically, each ANA image is decomposed into a bag of instances, denoted as $x_i = \{x_i^1, x_i^2, \dots, x_i^{N_i}\}$, where N_i is the number of extracted instances from image x_i . We denote by i the **image/bag** index, by j the **instance/patch** index within image x_i , and by k the **class/label** index. Let x_i be the input image and x_i^j its j -th instance. Instance-level predictions are first made independently, and then aggregated to produce the final image-level prediction.

Instead of uniform sampling, we introduce ***an instance sampler*** to extract meaningful sub-regions from ANA images. Let t_i^j denote a 0–1 label vector, where the k^{th} element $t_{i,k}^j$ is equal to 1 if l_k is present in x_i^j , and 0 otherwise. The general idea is to find a vector $\alpha_i^j = \{\alpha_{i,k}^j\}$ for each instance x_i^j , where the k^{th} element $\alpha_{i,k}^j$ indicates how confident the learner is about the presence of l_k in the instance. This sampler is designed to assign a confidence score to each candidate region and filter out irrelevant or uninformative areas (e.g., regions without ANA cells). This step ensures that only high-quality instances are used for downstream prediction.

The learning process can thus be formulated as a two-stage inference:

$$\hat{y}_i^j = \mathcal{F}(\mathcal{S}(x_i^j, \alpha_{i,k}^j, t_{i,k}^j)), \quad (2)$$

$$\hat{y}_i = \mathcal{G}(\{\hat{y}_i^j\}_{j=1}^{N_i}). \quad (3)$$

where $\mathcal{S}(\cdot)$ is the instance sampler function, $\mathcal{F}(\cdot)$ is the instance-level prediction function, and $\mathcal{G}(\cdot)$ denotes the aggregation function that synthesizes instance-level outputs into an image-level prediction \hat{y}_i . The learning objective minimizes the instance-level binary cross-entropy between \hat{y}_i^j and the pseudo-label $\bar{\xi}_i^j$ produced by \mathcal{D} (Sec. III-D).

However, only image-level ground-truth labels y_i are available in real-world scenarios, and instance-level annotations

y_i^j are typically not provided. To address this challenge, we introduce ***a probabilistic pseudo-label dispatcher***, which estimates a soft pseudo-label vector $\bar{\xi}_i^j = \{\bar{\xi}_{i,k}^j\}$ for each instance based on the image-level annotation:

$$\bar{\xi}_{i,k}^j = \mathcal{D}(\alpha_{i,k}^j, t_{i,k}^j). \quad (4)$$

Here $y_i \in \{0, 1\}^K$ is the image-level multi-hot ground-truth label. We set $t_i \equiv y_i$ and denote its broadcast to instance j by $t_i^j := t_i$ for all $j \in \{1, \dots, N_i\}$. In other words, t_i^j is the per-instance copy of y_i and the $\alpha_{i,k}^j$ denotes the learnable per-instance, per-class weight vector. We associate each instance x_i^j with a per-class confidence vector α_i^j . Both the instance sampler and the probabilistic pseudo-label dispatcher are built upon α_i^j . The concrete computation of \mathcal{D} is detailed in Eq. 8. The instance-level objective minimizes the binary cross-entropy between the predicted probabilities and the pseudo-labels. Furthermore, we also introduce ***a self-paced loss function*** $\mathcal{L}(\cdot)$ to dynamically adjust the contribution of each instance's loss to the overall objective.

In the following subsections, we illustrate our implementation of the instance-level prediction function $\mathcal{F}(\cdot)$, the aggregation function $\mathcal{G}(\cdot)$, the instance sampler $\mathcal{S}(\cdot)$, the probabilistic pseudo-label dispatcher $\mathcal{D}(\cdot)$, and the self-paced loss function $\mathcal{L}(\cdot)$.

B. Instance Sampler

We design the instance sampler $\mathcal{S}(\cdot)$ to select informative sub-regions from each ANA image, imitating the technician's strategy of identifying the most representative regions for ANA pattern recognition. To construct instances for multi-instance learning, we adopt a grid-based sliding window strategy to partition each input image x_i into N_i non-overlapping sub-regions, denoted as x_i^j . Each sub-region is a fixed-size patch of 448×448 pixels.

The idea is to encourage the selection of sub-regions in which the learner has developed higher confidence (from previous epochs) about ANA presence, while suppressing less confident ones.

$$s_i^j = \max_k [\text{ReLU}(\alpha_{i,k}^j) t_{i,k}^j], \quad \text{with } t_i^j \equiv t_i. \quad (5)$$

where $\text{ReLU}(\cdot)$ avoids negative weights. The per-instance sampling score s_i^j is the maximum (over classes) of the masked weights.

Then, the training instance x_i^j is sampled from the image using a weighted random sampling strategy based on s_i^j , which reflects the model's confidence in that region. In other words, instances x_i^j with higher s_i^j values (indicating higher confidence of containing relevant labels) are more likely to be selected and contribute to training. Instead of employing a summation or mean, the maximum is preferred due to two reasons: 1) Even if only a minority of ANA types are given high confidences, i.e., a distinguishable instance with dominating ANA types, the maximum strategy still ensures such an instance can be easily sampled. 2) When all ANA types of an instance share similar confidences, i.e., a hard instance with non-dominating ANA types, although the maximum strategy

sounds like a random one, the sampling probability is also low. Only confidences on labeled ANA types are considered, because we believe that ground truths labeled by technicians at the sample level are always correct, but distributing labels directly to sub-region instances is inappropriate. Instead of using the full set of sample-level labels, a subset obtained by setting certain elements of t_i^j from 1 to 0 is a more appropriate candidate for instance-level labels.

C. Instance-Level Prediction and Aggregation

After sampling the informative sub-regions x_i^j from image x_i , each instance x_i^j is passed through a shared prediction network to produce instance-level classification scores:

$$\hat{y}_i^j = \mathcal{F}(x_i^j, \theta), \quad (6)$$

where $\hat{y}_i^j \in [0, 1]^K$ denotes the predicted class probabilities for the j -th instance, and θ represents the parameters of a CNN-based backbone (e.g., ResNet).

To produce the final image-level prediction, we aggregate the instance-level outputs using an aggregation function $\mathcal{G}(\cdot)$ that operates in two steps: class-wise max pooling followed by thresholding. Specifically, for each class k , the aggregated prediction is given by:

$$\hat{y}_{i,k} = \mathbb{1} \left(\max_j \hat{y}_{i,k}^j > \tau \right), \quad (7)$$

where $\hat{y}_{i,k}$ is the predicted probability of class k for instance j , τ is a fixed threshold (e.g., 0.5), and $\mathbb{1}(\cdot)$ is the indicator function. This formulation reflects the multi-instance learning assumption: if any sub-region strongly supports a specific ANA pattern, the entire image is considered positive for that label.

D. Probabilistic Pseudo-label Dispatcher

In the absence of instance-level annotations, we introduce a probabilistic pseudo-label dispatcher $\mathcal{D}(\cdot)$ to generate soft supervisory signals for each instance. The probabilistic pseudo-label dispatcher distributes pseudo-labels based on the accumulated (sample-level) labels and the corresponding probabilities.

The weighting factor $\alpha_i^j = \{\alpha_{i,k}^j\}$ is leveraged to generate pseudo-labels, eliminating the need to introduce additional variables. For each instance, a scalar weighting factor is sufficient for the instance sampler; however, such granular weighting factors are incapable of pseudo-label adjustment. This motivates the design of a vector α_i^j for each instance x_i^j from the outset. Furthermore, pseudo-labels for each instance, such as predictions \hat{y}_i^j , are non-negative real numbers in the range $[0, 1]$, while ground truths are 0–1 label vectors t_i^j . Given the learnable weighting factor α_i^j , the pseudo-label $\bar{\xi}_{i,k}^j$ is defined as follows:

$$\mathcal{D}(\alpha_i^j, t_i^j) = \begin{cases} \xi_{i,k}^j = \text{ReLU}(\alpha_{i,k}^j)t_{i,k}^j \\ \bar{\xi}_{i,k}^j = \frac{\xi_{i,k}^j - \min_k \xi_{i,k}^j}{\max_k \xi_{i,k}^j - \min_k \xi_{i,k}^j} \end{cases}, \quad (8)$$

where $\bar{\xi}_{i,k}^j$ is the max–min normalized $\xi_{i,k}^j$. Element-wise multiplication by ground truths constrains generated pseudo-labels in subsets of accumulated labels, while max–min normalization ensures pseudo-labels in the range $[0, 1]$. Rather than 0–1 vectors, pseudo-labels in the range $[0, 1]$ are directly regarded as targets in the loss function, because real numbers can precisely reflect the level of confidence of the generated pseudo-labels. Instead of ground truths, pseudo-labels are generated and used for supervision in the training stage, while only ground truths are employed for evaluation and testing procedures.

E. Self-Paced Loss Function

In case CNN is selected to implement the prediction function, the learning is conducted by adaptive moment estimation (Adam) for minimizing the cross-entropy loss. The baseline image-level binary cross-entropy is:

$$\mathcal{L}_i = - \sum_{k=1}^K \left[t_{i,k} \log \hat{y}_{i,k} + (1 - t_{i,k}) \log (1 - \hat{y}_{i,k}) \right], \quad (9)$$

In real-world settings, only image-level labels are available (no instance-level annotations). We therefore introduce an instance sampler to select informative patches and a pseudo-label generator to enable instance-level training. Both the instance sampler and the pseudo-label generator are built upon α_i^j . Finally, α_i^j is used to control the contribution of the instance in learning, which results in a final dynamic self-paced loss function:

$$\begin{aligned} \mathcal{L}_i^{\text{apt}} &= - \sum_{k=1}^K \left[\alpha_{i,k}^j \bar{\xi}_{i,k}^j \log(\hat{y}_{i,k}^j) \right. \\ &\quad \left. + (1 - \bar{\xi}_{i,k}^j) \log(1 - \hat{y}_{i,k}^j) \right] \\ &= - \left[(\alpha_i^j \odot \bar{\xi}_i^j)^\top \log \hat{y}_i^j + (\mathbf{1} - \bar{\xi}_i^j)^\top \log(\mathbf{1} - \hat{y}_i^j) \right], \end{aligned} \quad (10)$$

where \odot is the Hadamard product. The key challenge is to design strategies for constructing and updating the two vectors α_i^j and $\bar{\xi}_i^j$, which are the results of the instance sampler and the pseudo-label generator, respectively.

IV. EXPERIMENTS

A. Datasets

We evaluate our proposed method from two perspectives: its effectiveness in ANA detection and its generalizability to other multi-label medical tasks. The ANA detection experiments are conducted on a real-world ANA dataset collected through the Integrated Care Organization (WCO), which comprises 686 hospitals across West China. The data aggregation process is led by West China Hospital, Sichuan University. The dataset contains 6,563 ANA immunofluorescence images, each labeled by clinicians during the diagnostic process, covering 8 ANA pattern classes. Labels were subsequently cleaned, corrected, and verified by a panel of three expert technicians. A label was finalized as ground truth only if at least two technicians agreed.

The dataset is categorized into 3,359 single-label samples, 2,685 double-label samples, and 519 samples with more than three labels. From a category perspective, the samples are distributed as follows: 116 Golgi (AC-22), 1,770 homogeneous (AC-1), 682 nucleolar (AC-8/9/10), 279 discrete nuclear (AC-6/7), 613 centromere (AC-3), 373 nuclear envelope (AC-11/12), 2,546 AMA (AC-21), and 4,115 speckled (AC-19/20). Labels from different categories may coexist in the same sample. To divide the dataset, 70% of the samples were assigned to the training set, while the validation and test sets each received 15%. For multi-instance experiments, visual patches were extracted from the samples using a sliding window algorithm, with each window sized 448×448 . This process resulted in 82,240 patches from 4,605 training samples, 19,400 patches from 939 validation samples, and 19,330 patches from 939 test samples.

To validate the generalizability of our approach, we conduct experiments on three real-world multi-label medical image datasets involving *whole-slide images* (WSIs):

- **NuCLS** [59]: 1,358 WSIs from TCGA breast cancer cohort, annotated with seven possible labels.
- **BCSS (Breast Cancer Semantic Segmentation)** [60]: 151 hematoxylin and eosin (H&E) stained WSIs covering 22 histologically-confirmed breast cancer cases.
- **PaNNuke** [61]: 7,904 WSIs spanning 19 tissue types with five nuclear structure labels.

Each dataset is split into 60% training, 10% validation, and 30% testing. Models are trained for 10 epochs with instance shuffling in each epoch.

B. Implementation Details

All experiments are implemented in PyTorch [62] and conducted on a computing platform equipped with an AMD EPYC 7542 CPU (2.9GHz), 94 GB of RAM, and an NVIDIA RTX 4090 GPU. For the real-world ANA dataset, we adopt the Adam optimizer [63] with $\beta = (0.9, 0.999)$, a learning rate of 5×10^{-3} , and a batch size of 32. All backbone networks are pre-trained on ImageNet [64], and input patches of size 448×448 are resized to match the input dimensions of the corresponding models.

For the three real-world multi-label medical image datasets, we follow a unified hyperparameter protocol across all competing methods unless otherwise specified in their original implementations. In our setup, each image is divided into smaller patches of approximately 64×64 pixels, with the total number of instances per image determined by its resolution. All models are trained and evaluated using 10-fold cross-validation to ensure robustness and fair comparison.

C. Evaluation Metrics

For the ANA task, we report four widely adopted metrics: $F1_{Micro}$, $F1_{Macro}$, accuracy, and mean Average Precision (mAP). Notably, $F1_{Micro}$ treats all instances equally, while $F1_{Macro}$ averages over all categories, allowing assessment from both global and per-class perspectives. An overall metric is also reported by averaging the above four scores.

For other medical tasks, we adopt four standard metrics commonly used in MIML literature: *Hamming Loss* (HL), *One Error* (OE), *Ranking Loss* (RL), and *Average Precision* (AP), following the definitions outlined in [65].

TABLE I

EXPERIMENTAL RESULTS OF SOTA METHODS AND OURS ON THE TEST SET. THE BEST RESULTS ARE **BOLD**. MODELS MARKED WITH \dagger ARE SPECIFICALLY DESIGNED FOR ANA DETECTION.

Model	$F1_{Micro}$	$F1_{Macro}$	Acc.	mAP	Overall
Multi-label					
Add-GCN [66]	0.771	0.673	0.440	0.799	0.671
ASL [67]	0.667	0.568	0.159	0.817	0.553
MIMLmiSVM [68]	0.483	0.088	0.243	0.145	0.240
MIMLkNN [69]	0.455	0.079	0.256	0.156	0.237
MIMLBoost [70]	0.461	0.091	0.182	0.168	0.226
DeepMIML [49]	0.493	0.098	0.278	0.199	0.267
MIMLfast [71]	0.513	0.355	0.478	0.677	0.506
Q2L-CvT [72]	0.582	0.275	0.324	0.433	0.404
ML-Decoder [73]	0.692	0.373	0.438	0.482	0.496
CBAM [74]	<u>0.792</u>	<u>0.695</u>	0.471	<u>0.812</u>	<u>0.693</u>
IDA-SwinL [75]	0.685	0.380	0.405	0.443	0.478
RRA † [76]	0.747	0.494	<u>0.505</u>	0.603	0.587
Ours	0.821	0.745	<u>0.573</u>	0.787	0.732
Single-label					
FUS † [77]	0.738	0.354	<u>0.738</u>	0.483	0.578
RRA † [76]	<u>0.835</u>	<u>0.711</u>	0.698	<u>0.872</u>	<u>0.779</u>
Ours	0.863	0.775	0.764	0.932	0.834

D. Comparison to SOTA Methods

To comprehensively evaluate the effectiveness of our approach, we compare it against a wide range of state-of-the-art (SOTA) methods on both the ANA dataset and three publicly available medical multi-label datasets, including NuCLS, BCSS, and PaNuke. The baselines include general-purpose multi-label classification models (e.g., CBAM [74], Add-GCN [66], ASL [67], Q2L-CvT [72], ML-Decoder [73], IDA-SwinL [75]), ANA-specific detection models (e.g., RRA [76], FUS [77]), and representative MIML methods (e.g., MIMLNN [68], MIMLSVM [70], MIMLmiSVM [68], MIMLkNN [69], MIMLBOOST [70], MIMLfast [71], DeepMIML [49], BMIML [65]).

1) *Evaluation on ANA Dataset*: To ensure a fair and comprehensive evaluation, each method adopts its default backbone as specified in the original publications: ResNet-50 for CBAM, ASL, and our method; ResNet-101 for Add-GCN, IDA-SwinL, RRA, and FUS; VGG16 for DeepMIML; CvT-w24 for Q2L-CvT; and TResNet-L for ML-Decoder.

ANA images naturally show multi-label characteristics. Multiple fluorescence patterns can appear in the same cell. However, most previous ANA detection methods were designed and tested under single-label assumptions. To enable a fair comparison, we additionally construct a single-label subset comprising 3,359 images annotated with only one ANA pattern. This allows benchmarking against single-label methods such as RRA and FUS.

As shown in Table I, our method achieves state-of-the-art performance in both settings. In the multi-label scenario, our

TABLE II

COMPARISON RESULTS (MEAN \pm STD.) ON THREE MEDICAL DATA SETS. \uparrow (\downarrow) INDICATES THAT THE LARGER (SMALLER) THE VALUE, THE BETTER THE PERFORMANCE; **BOLD** INDICATES THE BEST PERFORMANCE OF THIS METRIC; UNDERLINE INDICATES THE NEXT BEST PERFORMANCE OF THIS METRIC; N/A REPRESENTS THAT NO RESULT WAS OBTAINED IN 72 HOURS.

Methods	MIMLNN	MIMLSVM	MIMLmiSVM	MIMLkNN	MIMLBOOST	MIMLfast	DeepMIML	BMIML	Ours
<i>NuCLS</i>									
H.L. \downarrow	0.125 \pm 0.004	0.106 \pm 0.008	0.494 \pm 0.017	0.233 \pm 0.005	<u>0.116\pm0.025</u>	0.253 \pm 0.028	0.202 \pm 0.030	0.088 \pm 0.030	0.072\pm0.040
O.E. \downarrow	0.264 \pm 0.010	0.132 \pm 0.027	0.136 \pm 0.043	0.284 \pm 0.022	0.029\pm0.001	0.583 \pm 0.061	0.525 \pm 0.008	0.037 \pm 0.015	<u>0.030\pm0.030</u>
R.L. \downarrow	0.077 \pm 0.002	0.041 \pm 0.020	0.368 \pm 0.017	0.380 \pm 0.023	0.099 \pm 0.005	0.392 \pm 0.004	0.325 \pm 0.019	<u>0.043\pm0.010</u>	0.040\pm0.015
A.P. \uparrow	0.857 \pm 0.041	0.941 \pm 0.006	0.856 \pm 0.028	0.757 \pm 0.007	0.921 \pm 0.009	0.722 \pm 0.011	0.815 \pm 0.046	0.968\pm0.007	0.949 \pm 0.005
<i>Breast</i>									
H.L. \downarrow	0.293 \pm 0.060	0.297 \pm 0.011	0.511 \pm 0.041	0.297 \pm 0.033	0.460 \pm 0.030	0.318 \pm 0.021	0.541 \pm 0.032	0.290 \pm 0.017	0.288\pm0.015
O.E. \downarrow	0.219 \pm 0.013	0.206 \pm 0.032	0.183 \pm 0.003	0.250 \pm 0.062	<u>0.013\pm0.001</u>	0.500 \pm 0.016	0.500 \pm 0.003	0.094 \pm 0.001	0.083\pm0.001
R.L. \downarrow	<u>0.204\pm0.007</u>	0.196 \pm 0.050	0.438 \pm 0.028	0.483 \pm 0.010	0.943 \pm 0.041	0.493 \pm 0.022	0.502 \pm 0.046	<u>0.172\pm0.004</u>	0.169\pm0.010
A.P. \uparrow	0.822 \pm 0.028	.832 \pm 0.064	0.770 \pm 0.071	0.599 \pm 0.025	0.624 \pm 0.019	0.591 \pm 0.016	0.530 \pm 0.026	0.854\pm0.021	0.851 \pm 0.015
<i>Panuke</i>									
H.L. \downarrow	0.299 \pm 0.036	0.285 \pm 0.041	0.510 \pm 0.018	0.299 \pm 0.005	N/A	0.377 \pm 0.011	N/A	<u>0.276\pm0.005</u>	0.272\pm0.015
O.E. \downarrow	0.250 \pm 0.012	<u>0.167\pm0.024</u>	0.182 \pm 0.033	0.200 \pm 0.022	N/A	0.600 \pm 0.032	N/A	0.212 \pm 0.038	0.155\pm0.017
R.L. \downarrow	0.209 \pm 0.030	<u>0.189\pm0.006</u>	0.438 \pm 0.009	0.509 \pm 0.036	N/A	0.465 \pm 0.031	N/A	0.151\pm0.014	0.197 \pm 0.016
A.P. \uparrow	0.806 \pm 0.042	0.823 \pm 0.045	0.770 \pm 0.013	0.441 \pm 0.040	N/A	0.439 \pm 0.060	N/A	0.846\pm0.003	0.836 \pm 0.005

approach consistently outperforms all baselines in accuracy, $F1_{Micro}$, $F1_{Macro}$, and the overall averaged metric. Although ASL slightly surpasses our method in mAP, the difference is marginal. Meanwhile, our model's balanced performance across all metrics highlights its robustness and adaptability. DeepMIML ranks lowest, likely due to limitations in its shallow VGG16 backbone, lacking residual connections and multi-scale representation. Likewise, general-purpose models such as Q2L-CvT and ML-Decoder underperform, possibly because they were not tailored for the unique structural features of ANA data.

In the single-label setting, our model again surpasses RRA and FUS in all evaluation metrics, showcasing its strong generalization ability even under constrained labeling conditions. Given the prevalence of single-label classification in clinical ANA datasets, this result underscores the real-world applicability of our framework.

2) Evaluation on Public MIML Medical Datasets: While many classical MIML methods, such as MIMLNN, MIMLSVM, MIMLmiSVM, MIMLkNN, MIMLBOOST, MIMLfast, and BMIML, rely on pre-extracted and fixed instance-level features, only DeepMIML adopts an end-to-end neural network backbone to learn instance representations jointly with bag-level prediction. Traditional approaches often depend on hand-crafted features or simple image statistics, and subsequently apply shallow models like kernel-based classifiers, boosting frameworks, or k-nearest-neighbor algorithms. Though computationally lightweight and practical for small-scale data, these methods lack the capacity to capture rich visual semantics and are inherently limited when applied to high-resolution medical images or large datasets, where joint optimization is critical.

In contrast, our method fully leverages deep feature learning and achieves significantly stronger results across diverse and complex domains. As summarized in Table II, we evaluate our framework on three challenging public MIML medical datasets: NuCLS, BCSS, and PaNuke. Our method achieves

the best or second-best performance in nearly all evaluation metrics: it attains the lowest Hamming Loss and Ranking Loss on NuCLS, the best One Error and Accuracy on BCSS, and maintains leading performance on PaNuke, despite its high inter-class visual variability. These results highlight the superior generalization ability of our approach and its robustness across varied multi-label medical image classification tasks. By jointly learning instance features and label dependencies, our method consistently outperforms traditional MIML pipelines and demonstrates its scalability and adaptability to real-world medical applications.

E. Ablation Study

To study the importance of different components of the proposed method, we have conducted several ablation experiments.

1) Backbone Network: The backbone network is crucial for feature extraction in the proposed method. We evaluate six widely used architectures: VGG16 [78], Inception-v3 [79], ResNet-50 [80], ResNet-101 [80], CvT-w24 [81], and TResNet-L [82]. Each backbone follows its default input resolution and is assessed under the same experimental conditions. As shown in Table IV, ResNet-50 achieves the best overall performance across key metrics, making it the optimal choice for integrating with our proposed approach.

From Table IV, it is evident that ResNet-50 achieves the best overall performance, surpassing other models in $F1_{Micro}$, $F1_{Macro}$, and Accuracy, making it the most suitable choice for ANA detection. TResNet-L achieves the highest mAP, indicating its strength in ranking-based evaluation, but its overall performance remains slightly behind ResNet-50. ResNet-101 performs competitively but does not surpass ResNet-50, likely due to increased computational complexity without a proportional gain in feature representation.

In contrast, VGG16 performs the worst across all metrics, suggesting that its shallower architecture lacks the depth required to extract meaningful features for ANA detection.

TABLE III

EXPERIMENTAL RESULTS OF ABLATION STUDY ON THE TEST SET. THE BEST RESULTS ARE IN BOLD FONT. “GRANULARITY” REFERS TO THE LEVEL AT WHICH SAMPLING OR LEARNING WEIGHTS ARE APPLIED: \times INDICATES NO WEIGHTS (UNIFORM SAMPLING), \bullet INDICATES INSTANCE-LEVEL GRANULARITY (EACH PATCH HAS A DISTINCT WEIGHT), \circ INDICATES LABEL-LEVEL GRANULARITY (EACH PATCH HAS ONE WEIGHT PER LABEL). \times , \blacktriangle , AND \triangle IN THE COLUMN INITIALIZATION INDICATE THE INITIALIZATION METHODS OF WEIGHTS: NO WEIGHTS, FREQUENCIES AVERAGED ON ALL DATA, AND FREQUENCIES AVERAGED ON EACH SAMPLE, RESPECTIVELY. \times AND \checkmark IN THE COLUMN SAMPLING, PSEUDO-LABEL, AND COEFFICIENTS INDICATE WHETHER THOSE COMPONENTS ARE USED, RESPECTIVELY.

Model	Granularity	Initialization	Sampling	Pseudo-label	Coefficients	$F1_{Micro}$	$F1_{Macro}$	Acc.	mAP	Overall
1	\times	\times	\times	\times	\times	0.717	0.353	0.367	0.445	0.470
2	\bullet	\blacktriangle	\times	\times	\times	0.550	0.200	0.299	0.357	0.352
3	\bullet	\blacktriangle	\times	\times	\checkmark	0.693	0.417	0.402	0.503	0.504
4	\bullet	\blacktriangle	\checkmark	\times	\times	0.774	0.613	0.451	0.692	0.633
5	\bullet	\blacktriangle	\checkmark	\times	\checkmark	0.788	0.693	0.533	0.746	0.690
6	\circ	\blacktriangle	\times	\times	\times	0.242	0.199	0.136	0.441	0.255
7	\circ	\blacktriangle	\times	\times	\checkmark	0.672	0.395	0.373	0.516	0.489
8	\circ	\blacktriangle	\times	\checkmark	\times	0.173	0.085	0.093	0.481	0.208
9	\circ	\blacktriangle	\times	\checkmark	\checkmark	0.753	0.457	0.436	0.538	0.546
10	\circ	\blacktriangle	\checkmark	\times	\times	0.780	0.675	0.547	0.762	0.691
11	\circ	\blacktriangle	\checkmark	\times	\checkmark	0.790	0.678	0.533	0.761	0.690
12	\circ	\blacktriangle	\checkmark	\checkmark	\times	0.770	0.707	0.541	0.804	0.705
13	\circ	\blacktriangle	\checkmark	\checkmark	\checkmark	0.780	0.682	0.485	0.747	0.674
14	\circ	\triangle	\times	\times	\times	0.321	0.197	0.175	0.472	0.291
15	\circ	\triangle	\times	\times	\checkmark	0.719	0.435	0.443	0.529	0.532
16	\circ	\triangle	\times	\checkmark	\times	0.050	0.084	0.031	0.464	0.157
17	\circ	\triangle	\times	\checkmark	\checkmark	0.729	0.461	0.422	0.578	0.547
18	\circ	\triangle	\checkmark	\times	\times	0.757	0.663	0.543	0.770	0.683
19	\circ	\triangle	\checkmark	\times	\checkmark	0.783	0.676	0.524	0.758	0.685
20	\circ	\triangle	\checkmark	\checkmark	\times	0.779	0.685	0.554	0.778	0.699
21	\circ	\triangle	\checkmark	\checkmark	\checkmark	0.821	0.745	0.573	0.787	0.732

TABLE IV

EXPERIMENTAL RESULTS OF DIFFERENT BASE MODELS ON THE TEST SET. THE BEST RESULTS ARE IN BOLD FONT.

Backbone	Input Resolution	$F1_{Micro}$	$F1_{Macro}$	Acc.	mAP
VGG16 [78]	224×224	0.030	0.009	0.004	0.199
Inception-v3 [79]	299×299	0.599	0.185	0.285	0.323
ResNet-50 [80]	224×224	0.717	0.353	0.467	0.445
ResNet-101 [80]	448×448	0.672	0.311	0.337	0.421
CvT-w24 [81]	384×384	0.698	0.297	0.372	0.433
TResNet-L [82]	448×448	0.701	0.335	0.407	0.461

Inception-v3 also struggles, particularly in $F1_{Macro}$, which indicates poor performance on less frequent ANA patterns. This may be attributed to its convolutional layers with a stride of 2, which enlarge the receptive field, inadvertently introducing irrelevant neighboring information and leading to misclassification. Both CvT-w24 and TResNet-L achieve moderate performance, demonstrating that vision transformers and optimized CNN architectures have potential in ANA detection. However, their performance does not surpass ResNet-50, indicating that a well-balanced deep residual network remains the most effective backbone for this task.

2) Ablation Analysis of Core Components and Weighting Strategies: Besides the three components introduced in III, we also evaluate the influence of different initializing strategies for weights and different granularities of weights. There are two types of initializations for weights, namely initializations with frequencies averaged on all data and each sample, respectively. The former counts the frequency numbers of each ANA type over all data, and the weights of all samples share the same

initial values, whereas the latter only counts that over each sample, and the weights within the same sample share the initial values. The formal definitions of the two initializations are written as:

$$w_k^d = \text{softmax}\left(\frac{\sum_{i=1}^N (t_{i,k}^j)}{\sum_{k=1}^K \sum_{i=1}^N (t_{i,k}^j)}\right) \quad (11)$$

$$w_{p,k}^s = \text{softmax}\left(\frac{\sum_{i \in S_p} (t_{i,k}^j)}{\sum_{k=1}^K \sum_{i \in S_p} (t_{i,k}^j)}\right)$$

where $w_k^d, w_{p,k}^s$ are initialized weights averaged on all data and each sample respectively, $\text{softmax}(\cdot)$ denotes the softmax function applied along the label dimension, and S_p is the extracted patch set of p -th sample. Rather than weights at the label-level, weights at the instance-level, namely a scalar weight for each instance instead of a vector weight, are also evaluated in the ablation study. Due to the enlarged granularity, the previous initialization of weights w_k^d is not applicable and changed to:

$$w_{q,k}^d = \text{softmax}\left(\frac{N_{L_q}}{N}\right) \quad (12)$$

where $w_{q,d}^d$ denotes the new initialized weights averaged on all data at the instance-level and N_{L_q} is the number of instances, which share the same label $L_q \in \mathbb{Z}^K$, where K represents the total number of categories of labels. Besides, in the case of weights at the instance level, it is not able to implement the initialization of frequencies averaged on each sample and

the probabilistic pseudo-label dispatcher, the corresponding experimental results of which are not reported in this section.

From Table III, it is evident that the proposed components effectively improve the performance of the model. From the perspective of granularities of weights, all experiments in the ablation study except for the first one can be divided into two groups, one is the case of weights at the instance-level and the other is that at the label-level (\bullet and \circ in column Granularity). For the former group, although the performance is limited due to the granularity of weights, the effectiveness of the instance sampler and label-aware coefficients can be observed. For the latter group, the experiments are further organized into two subgroups (\blacktriangle and \triangle in column Initialization) in terms of the initialization method of weights formulated in Eq. 11. Only a slight performance difference between the two initialization methods can be seen from the two subgroup experiments, where significant performance improvement has been demonstrated for the proposed three components. The instance sampler contributes most to the performance improvement, while the probabilistic pseudo-label and label-aware coefficients usually provide minor improvement. Among all experiments, the instance sampler improves the overall performance by at least 23.44%. The probabilistic pseudo-label contributes positively when models are stable, which ensures the correctness of generated pseudo-labels and avoids wrong leads. Besides improving the performance, label-aware coefficients of weight learning rates also effectively stabilize the learning procedure and prevent it from being stuck in local optima.

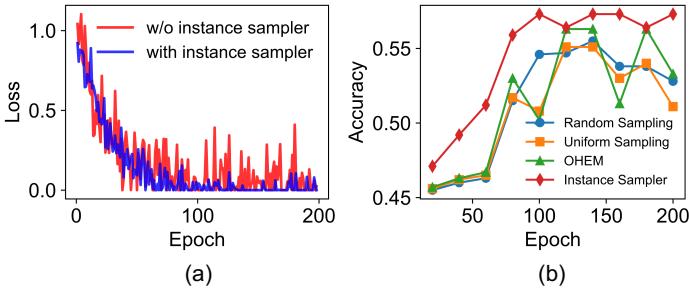


Fig. 4. (a) Training loss curves comparing the model with and without the instance sampler. (b) Performance results under different sampling strategies: Random, Uniform, and OHEM.

3) Effectiveness of Instance Sampler: In Fig. 4 (a), we present the training loss curves comparing the model with and without the instance sampler, showing that our sampler not only accelerates convergence but also stabilizes the training process. Meanwhile, Fig. 4 (b) reports the final accuracy under various sampling strategies, including random sampling, uniform sampling, and OHEM [83]. The results reveal that our instance sampler consistently outperforms these alternatives in the multi-label ANA classification task.

Our instance sampler achieves higher overall performance than both the Random and Uniform baselines, indicating that merely relying on random or evenly distributed sampling fails to exploit the inherent diversity among subregions. By contrast, our method effectively identifies and prioritizes subregions that exhibit representative ANA patterns. Although

TABLE V
PERFORMANCE COMPARISON UNDER DIFFERENT VALUES OF C .

C	Accuracy (%)	mAP (%)
2	51.2	74.5
3	52.4	75.6
4	54.1	76.2
5	55.5	76.9
6	57.3	78.7
7	45.1	62.5

OHEM does focus on high-loss samples to some extent, it relies primarily on loss value and overlooks the interplay or co-occurrence between different ANA types in a multi-label setting. In comparison, our instance sampler leverages self-paced learning and the dynamic weighting of α_i^j , enabling a more comprehensive assessment of subregion value in multi-label tasks and thereby yielding superior results.

4) Sensitivity Analysis of Label-aware Scaling Parameter C : We also explore the impact of different values of C on model performance, as shown in Table V. Since C determines the maximum number of coexisting labels per sample, it directly influences the label-aware coefficient λ , thereby adjusting the learning rate of weights. Intuitively, increasing C allows the model to better handle multi-label complexity by adapting to more challenging cases with higher learning rates. Gradually increasing C from 2 to 6 consistently improves both accuracy and mAP, suggesting that enabling the model to learn from more complex multi-label samples enhances its generalization ability. Notably, performance saturates and peaks at $C = 6$, which aligns with the actual maximum number of coexisting labels in our dataset. However, when C is set to 7, a significant performance drop is observed, likely due to excessive learning rates leading to instability. This finding highlights the importance of setting an appropriate C value to balance learning effectiveness and model stability.

5) Sensitivity to Patch Size: To investigate the effect of patch size on classification performance, we evaluate our method using patch resolutions ranging from 32×32 to 640×640 , along with the original unpatched image. All patches are resized to 224×224 before being fed into the ResNet backbone. As shown in Table VI, extremely small patches lead to massive instance counts and degraded performance due to excessive noise. Medium-sized patches (e.g., 224×224 and 448×448) offer a better balance of localization and context, achieving the best results across metrics. In particular, 448×448 yields the highest $F1_{Macro}$ and mAP. Additionally, we include the average number of subregions (N) per image to account for variation in image resolution, as different patch sizes produce different numbers of instances.

To provide further insight, we visualize Grad-CAM [84] heatmaps under different patch sizes in Figure 5. The results reveal that 448×448 best localizes class-discriminative regions, validating its effectiveness in capturing informative foreground content.

TABLE VI
SENSITIVITY ANALYSIS OF DIFFERENT PATCH SIZES

Patch Size	Input instances	\bar{N}	$F1_{\text{Micro}}$	$F1_{\text{Macro}}$	Acc.	mAP
32 × 32	22,158,220	4811	0.080	0.012	0.011	0.179
64 × 64	5,597,220	1215	0.088	0.025	0.017	0.199
128 × 128	1,390,620	301	0.440	0.215	0.277	0.332
224 × 224	412,250	89	0.840	0.725	0.554	0.760
448 × 448	82,240	17	0.821	0.745	0.573	0.787
640 × 640	40,985	8	0.792	0.729	0.522	0.701
Original	4,605	8	0.672	0.311	0.337	0.421

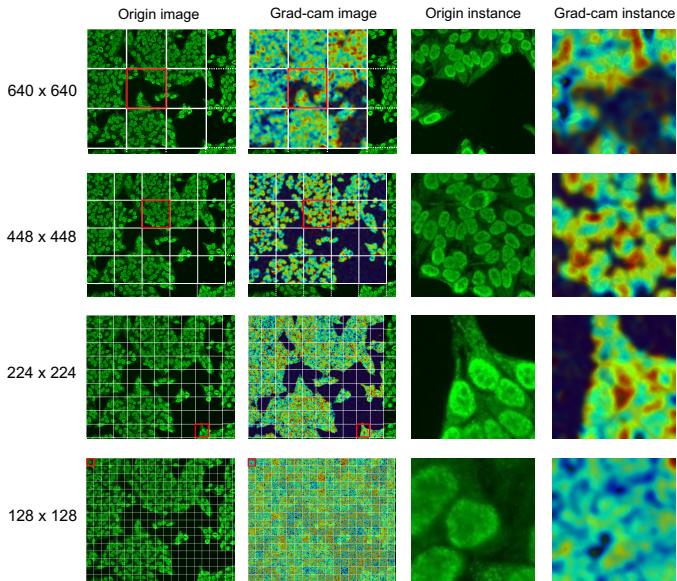


Fig. 5. Grad-CAM visualization for different patch sizes. Patch size of 448 × 448 demonstrates the most focused activation on foreground ANA patterns.

V. CONCLUSION

In this paper, we have proposed a new framework consisting of three novel modules to address the multi-instance multi-label challenge of ANA detection. Empirical results on the ANA dataset demonstrate that the proposed framework consistently outperforms existing state-of-the-art approaches across a range of metrics and settings. The comparison results on other MIML medical image domains verify the generalization ability of our proposed framework. Despite these promising results, the current framework relies on heuristic-based confidence estimation and fixed-size patch sampling, which may overlook nuanced patterns in more complex visual contexts. In the future, we plan to explore continuous attention mechanisms to better capture nuanced visual patterns and improve the understanding of how pseudo-labels are assigned. We will also extend our framework to more medical imaging tasks.

REFERENCES

- [1] S. Ghosh and V. Chaudhary, “Feature analysis for automatic classification of hep-2 fluorescence patterns : Computer-aided diagnosis of autoimmune diseases,” in *Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*, 2012, pp. 174–177.
- [2] G. Thibault and J. Angulo, “Efficient statistical/morphological cell texture characterization and classification,” in *Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*, 2012, pp. 2440–2443.
- [3] C. Vununu, S.-H. Lee, and K.-R. Kwon, “A deep feature extraction method for hep-2 cell image classification,” *Electronics*, vol. 8, no. 1, 2019. [Online]. Available: <https://www.mdpi.com/2079-9292/8/1/20>
- [4] H. Lei, T. Han, F. Zhou, Z. Yu, J. Qin, A. Elazab, and B. Lei, “A deeply supervised residual network for hep-2 cell classification via cross-modal transfer learning,” *Pattern Recognition*, vol. 79, pp. 290–302, 2018. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0031320318300608>
- [5] C. Vununu, S.-H. Lee, O.-J. Kwon, and K.-R. Kwon, “A dynamic learning method for the classification of the hep-2 cell images,” *Electronics*, vol. 8, no. 8, 2019. [Online]. Available: <https://www.mdpi.com/2079-9292/8/8/850>
- [6] X. Jia, L. Shen, X. Zhou, and S. Yu, “Deep convolutional neural network based hep-2 cell classification,” in *2016 23rd International Conference on Pattern Recognition (ICPR)*, 2016, pp. 77–80.
- [7] H. Xie, Y. He, H. Lei, T. Han, Z. Yu, and B. Lei, “Deeply supervised residual network for hep-2 cell classification,” in *2018 24th International Conference on Pattern Recognition (ICPR)*, 2018, pp. 699–703.
- [8] V. B. S. Prasath, Y. M. Kassim, Z. A. Oraibi, J.-B. Guiriec, A. Hafiane, G. Seetharaman, and K. Palaniappan, “Hep-2 cell classification and segmentation using motif texture patterns and spatial features with random forests,” in *2016 23rd International Conference on Pattern Recognition (ICPR)*, 2016, pp. 90–95.
- [9] J. Liu, B. Xu, L. Shen, J. Garibaldi, and G. Qiu, “Hep-2 cell classification based on a deep autoencoding-classification convolutional neural network,” in *2017 IEEE 14th International Symposium on Biomedical Imaging (ISBI 2017)*, 2017, pp. 1019–1023.
- [10] D. BS, K. Subramaniam, and N. HR, “Hep-2 cell classification using artificial neural network approach,” in *2016 23rd International Conference on Pattern Recognition (ICPR)*, 2016, pp. 84–89.
- [11] Z. Gao, L. Wang, L. Zhou, and J. Zhang, “Hep-2 cell image classification with deep convolutional neural networks,” *IEEE Journal of Biomedical and Health Informatics*, vol. 21, no. 2, pp. 416–428, 2017.
- [12] H. Li, H. Huang, W.-S. Zheng, X. Xie, and J. Zhang, “Hep-2 specimen classification via deep cnns and pattern histogram,” in *2016 23rd International Conference on Pattern Recognition (ICPR)*, 2016, pp. 2145–2149.
- [13] S. Manivannan, W. Li, S. Akbar, R. Wang, J. Zhang, and S. J. McKenna, “An automated pattern recognition system for classifying indirect immunofluorescence images of hep-2 cells and specimens,” *Pattern Recognition*, vol. 51, pp. 12–26, 2016. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0031320315003465>
- [14] C. Xiao et al., “Confusion-resistant federated learning via diffusion-based data harmonization on non-iid data,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 137495–137520, 2024.
- [15] C.-W. Ngo, Y.-G. Jiang, X.-Y. Wei, F. Wang, W. Zhao, H.-K. Tan, and X. Wu, “Experimenting viroeo-374: Bag-of-visual-words and visual-based ontology for semantic video indexing and search,” in *IEEE Computer Society*, 2007.
- [16] C.-W. Ngo, Y.-G. Jiang, X.-Y. Wei, W. Zhao, F. Wang, X. Wu, and H.-K. Tan, “Beyond semantic search: What you observe may not be what you think,” in *IEEE Computer Society*, 2008.
- [17] R. Tong, J. Liu, S. Liu, X. Hu, and L. Wang, “Renaissance of rnns in streaming clinical time series: Compact recurrence remains competitive with transformers,” *arXiv preprint arXiv:2510.16677*, 2025.
- [18] R. Tong, J. Liu, S. Liu, J. Xu, L. Wang, and T. Wang, “Does bigger mean better? comparative analysis of cnns and biomedical vision language modules in medical diagnosis,” *arXiv preprint arXiv:2510.00411*, 2025.
- [19] Y. Jiang, W. Zhang, X. Zhang, X.-Y. Wei, C. W. Chen, and Q. Li, “Prior knowledge integration via llm encoding and pseudo event regulation for video moment retrieval,” in *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024, pp. 7249–7258.
- [20] C. Xiao, C. Zhao, Z. Ke, and F. Shen, “Curiosity meets cooperation: A game-theoretic approach to long-tail multi-label learning,” *arXiv preprint arXiv:2510.17520*, 2025.
- [21] C. Xiao, L. Hou, L. Fu, and W. Chen, “Diffusion-based self-supervised imitation learning from imperfect visual servoing demonstrations for robotic glass installation,” in *2025 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2025, pp. 10401–10407.
- [22] J. Ge, J. Cao, X. Zhu, X. Zhang, C. Liu, K. Wang, and B. Liu, “Consistencies are all you need for semi-supervised vision-language tracking,” in *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024, pp. 1895–1904.
- [23] J. Ge, J. Cao, Y. Bao, B. Cao, and B. Liu, “Gal: combining global and local contexts for interpersonal relation extraction toward document-level chinese text,” *Neural Computing and Applications*, vol. 36, no. 11, pp. 5715–5731, 2024.

- [24] J. Ge, J. Cao, X. Chen, X. Zhu, W. Liu, C. Liu, K. Wang, and B. Liu, "Beyond visual cues: Synchronously exploring target-centric semantics for vision-language tracking," *ACM Transactions on Multimedia Computing, Communications and Applications*, vol. 21, no. 5, pp. 1–21, 2025.
- [25] J.-W. Ge, J.-X. Cao, Z.-X. Zhao, and B. Liu, "Fsd-gan: Generative adversarial training for face swap detection via the latent noise fingerprint," *Journal of Computer Science and Technology*, vol. 40, no. 2, pp. 397–412, 2025.
- [26] J. Ge, X. Zhang, J. Cao, X. Zhu, W. Liu, Q. Gao, B. Cao, K. Wang, C. Liu, B. Liu *et al.*, "Gen4track: A tuning-free data augmentation framework via self-correcting diffusion model for vision-language tracking," in *Proceedings of the 33rd ACM International Conference on Multimedia*, 2025, pp. 3037–3046.
- [27] N. Agmon-Levin, J. Damoiseaux, C. Kallenberg, U. Sack, T. Witte, M. Herold, X. Bossuyt, L. Musset, R. Cervera, A. Plaza-Lopez, C. Dias, M. J. Sousa, A. Radice, C. Eriksson, O. Hultgren, M. Viander, M. Khamashta, S. Regenass, L. E. C. Andrade, A. Wiik, A. Tincani, J. Rönnelid, D. B. Bloch, M. J. Fritzler, E. K. L. Chan, I. Garcia-De La Torre, K. N. Konstantinov, R. Lahita, M. Wilson, O. Vainio, N. Fabien, R. A. Sinico, P. Meroni, and Y. Shoenfeld, "International recommendations for the assessment of autoantibodies to cellular antigens referred to as anti-nuclear antibodies," *Annals of the Rheumatic Diseases*, vol. 73, no. 1, pp. 17–23, 2014. [Online]. Available: <https://ard.bmjjournals.org/content/73/1/17>
- [28] C. von Muhlen, I. Garcia-De La Torre, M. Infantino, J. Damoiseaux, L. Andrade, O. Carballo, K. Conrad, P. Francescantonio, M. Fritzler, M. Herold, W. Klötz, W. Cruvinel, T. Mimori, M. Satoh, L. Musset, and E. Chan, "How to report the antinuclear antibodies (anti-cell antibodies) test on hep-2 cells: guidelines from the icap initiative," *Immunologic Research*, vol. 69, no. 6, pp. 594–608, Dec. 2021.
- [29] Z.-L. Zhang and M.-L. Zhang, "Multi-instance multi-label learning with application to scene classification," in *Advances in Neural Information Processing Systems*, B. Schölkopf, J. Platt, and T. Hoffman, Eds., vol. 19. MIT Press, 2006.
- [30] J. Keeler, D. Rumelhart, and W. Leow, "Integrated segmentation and recognition of hand-printed numerals," in *Advances in Neural Information Processing Systems*, R. Lippmann, J. Moody, and D. Touretzky, Eds., vol. 3. Morgan-Kaufmann, 1990.
- [31] T. G. Dietterich, R. H. Lathrop, and T. Lozano-Pérez, "Solving the multiple instance problem with axis-parallel rectangles," *Artificial Intelligence*, vol. 89, no. 1, pp. 31–71, 1997. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0004370296000343>
- [32] M.-L. Zhang and Z.-H. Zhou, "A review on multi-label learning algorithms," *IEEE Transactions on Knowledge and Data Engineering*, vol. 26, no. 8, pp. 1819–1837, 2014.
- [33] X. Wang, H. Chen, C. Gan, H. Lin, Q. Dou, E. Tsougenis, Q. Huang, M. Cai, and P.-A. Heng, "Weakly supervised deep learning for whole slide lung cancer image analysis," *IEEE Transactions on Cybernetics*, vol. 50, no. 9, pp. 3950–3962, 2020.
- [34] N. Coudray, P. S. Ocampo, T. Sakellaropoulos, N. Narula, M. Snuderl, D. Fenyö, A. L. Moreira, N. Razavian, and A. Tsirigos, "Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning," *Nature medicine*, vol. 24, no. 10, pp. 1559–1567, 2018.
- [35] X.-Y. Wei, Z.-Q. Yang, X.-L. Zhang, G. Liao, A.-L. Sheng, S. K. Zhou, Y. Wu, and L. Du, "Deep collocative learning for immunofixation electrophoresis image analysis," *IEEE transactions on medical imaging*, vol. 40, no. 7, pp. 1898–1910, 2021.
- [36] P. Duygulu, K. Barnard, J. F. G. de Freitas, and D. A. Forsyth, "Object recognition as machine translation: Learning a lexicon for a fixed image vocabulary," in *Computer Vision — ECCV 2002*, A. Heyden, G. Sparr, M. Nielsen, and P. Johansen, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2002, pp. 97–112.
- [37] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y. Zheng, "Nus-wide: A real-world web image database from national university of singapore," in *Proceedings of the ACM International Conference on Image and Video Retrieval*, ser. CIVR '09. New York, NY, USA: Association for Computing Machinery, 2009. [Online]. Available: <https://doi.org/10.1145/1646396.1646452>
- [38] L. Fei-Fei and P. Perona, "A bayesian hierarchical model for learning natural scene categories," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, vol. 2, 2005, pp. 524–531 vol. 2.
- [39] J. Read, B. Pfahringer, G. Holmes, and E. Frank, "Classifier chains for multi-label classification," in *Machine Learning and Knowledge Discovery in Databases*, W. Buntine, M. Grobelnik, D. Mladenić, and J. Shawe-Taylor, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 254–269.
- [40] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision — ECCV 2014*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Cham: Springer International Publishing, 2014, pp. 740–755.
- [41] C. Han, X. Pan, L. Yan, H. Lin, B. Li, S. Yao, S. Lv, Z. Shi, J. Mai, J. Lin, B. Zhao, Z. Xu, Z. Wang, Y. Wang, Y. Zhang, H. Wang, C. Zhu, C.-Y. Lin, L. Mao, M. Wu, L. Duan, J. Zhu, D. Hu, Z. Fang, Y. Chen, Y. Zhang, Y. Li, Y. Zou, Y.-Z. Yu, X. Li, H. L. Li, Y. hai Cui, G. Han, Y. Xu, J. Xu, H. Yang, C. Li, Z. Liu, C. Lu, X. Chen, C. Liang, Q. Zhang, and Z. Liu, "Wssss4luad: Grand challenge on weakly-supervised tissue semantic segmentation for lung adenocarcinoma," *ArXiv*, vol. abs/2204.06455, 2022.
- [42] H. Li, H. Huang, W.-S. Zheng, X. Xie, and J. Zhang, "Hep-2 specimen classification via deep cnns and pattern histogram," in *2016 23rd International Conference on Pattern Recognition (ICPR)*. IEEE, 2016, pp. 2145–2149.
- [43] S. Manivannan, W. Li, S. Akbar, R. Wang, J. Zhang, and S. J. McKenna, "An automated pattern recognition system for classifying indirect immunofluorescence images of hep-2 cells and specimens," *Pattern Recognition*, vol. 51, pp. 12–26, 2016.
- [44] V. Snell, W. Christmas, and J. Kittler, "Hep-2 fluorescence pattern classification," *Pattern Recognition*, vol. 47, no. 7, pp. 2338–2347, 2014. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0031320313004226>
- [45] C. Cortes and V. Vapnik, "Support-vector networks," *Mach. Learn.*, vol. 20, no. 3, p. 273–297, sep 1995. [Online]. Available: <https://doi.org/10.1023/A:1022627411411>
- [46] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [47] W. Zhang, Q. Tian, Y. Cao, W. Fan, D. Jiang, Y. Wang, Q. Li, and X.-Y. Wei, "Graphate: advancing multilevel and multi-label anatomical therapeutic chemical classification via atom-level graph learning," *Briefings in Bioinformatics*, vol. 26, no. 2, p. bbaf194, 2025.
- [48] K. Al-Dulaimi, V. Chandran, K. Nguyen, J. Banks, and I. Tomeo-Reyes, "Benchmarking hep-2 specimen cells classification using linear discriminant analysis on higher order spectra features of cell shape," *Pattern Recognition Letters*, vol. 125, pp. 534–541, 2019. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167865518304811>
- [49] J. Feng and Z.-H. Zhou, "Deep mml network," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, ser. AAAI'17. AAAI Press, 2017, p. 1884–1890.
- [50] M. Kumar, B. Packer, and D. Koller, "Self-paced learning for latent variable models," in *Advances in Neural Information Processing Systems*, J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, Eds., vol. 23. Curran Associates, Inc., 2010.
- [51] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proceedings of the 26th Annual International Conference on Machine Learning*, ser. ICML '09. New York, NY, USA: Association for Computing Machinery, 2009, p. 41–48. [Online]. Available: <https://doi.org/10.1145/1553374.1553380>
- [52] F. Khan, B. Mutlu, and J. Zhu, "How do humans teach: On curriculum learning and teaching dimension," in *Advances in Neural Information Processing Systems*, J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Weinberger, Eds., vol. 24. Curran Associates, Inc., 2011.
- [53] D. Zhang, D. Meng, C. Li, L. Jiang, Q. Zhao, and J. Han, "A self-paced multiple-instance learning framework for co-saliency detection," in *2015 IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 594–602.
- [54] D. Zhang, D. Meng, and J. Han, "Co-saliency detection via a self-paced multiple-instance learning framework," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 5, pp. 865–878, 2017.
- [55] C. Li, F. Wei, J. Yan, X. Zhang, Q. Liu, and H. Zha, "A self-paced regularization framework for multilabel learning," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 6, pp. 2660–2666, 2018.
- [56] Y. Zhong, B. Du, and C. Xu, "Learning to reweight examples in multi-label classification," *Neural Networks*, vol. 142, pp. 428–436, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0893608021001106>
- [57] M. Ren, W. Zeng, B. Yang, and R. Urtasun, "Learning to reweight examples for robust deep learning," in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings

- of Machine Learning Research, J. Dy and A. Krause, Eds., vol. 80. PMLR, 10–15 Jul 2018, pp. 4334–4343. [Online]. Available: <https://proceedings.mlr.press/v80/ren18a.html>
- [58] J. Shu, Q. Xie, L. Yi, Q. Zhao, S. Zhou, Z. Xu, and D. Meng, *Meta-Weight-Net: Learning an Explicit Mapping for Sample Weighting*. Red Hook, NY, USA: Curran Associates Inc., 2019.
- [59] M. Amgad, L. A. Atteya, H. Hussein, K. H. Mohammed, E. Hafiz, M. A. Elsebaie, A. M. Alhusseini, M. A. AlMoslemany, A. M. Elmatboly, P. A. Pappalardo *et al.*, “Nucls: A scalable crowdsourcing, deep learning approach and dataset for nucleus classification, localization and segmentation,” *arXiv preprint arXiv:2102.09099*, 2021.
- [60] M. Amgad, H. Elfandy, H. Hussein, L. A. Atteya, M. A. Elsebaie, L. S. Abo Elnasr, R. A. Sakr, H. S. Salem, A. F. Ismail, A. M. Saad *et al.*, “Structured crowdsourcing enables convolutional segmentation of histology images,” *Bioinformatics*, vol. 35, no. 18, pp. 3461–3467, 2019.
- [61] J. Gamper, N. A. Koohbanani, K. Benes, S. Graham, M. Jahanifar, S. A. Khurram, A. Azam, K. Hewitt, and N. Rajpoot, “Pannuke dataset extension, insights and baselines,” *arXiv preprint arXiv:2003.10778*, 2020.
- [62] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems* 32, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- [63] D. Kingma and J. L. Ba, “3rd international conference on learning representations, iclr 2015-conference track proceedings,” in *International Conference on Learning Representations, ICLR) Adam: A method for stochastic optimization Go to reference in article*, 2015.
- [64] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [65] Q. Lai, J. Zhou, Y. Gan, C.-M. Vong, and C. P. Chen, “Single-stage broad multi-instance multi-label learning (bmiml) with diverse inter-correlations and its application to medical image classification,” *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 8, no. 1, pp. 828–839, 2024.
- [66] J. Ye, J. He, X. Peng, W. Wu, and Y. Qiao, “Attention-driven dynamic graph convolutional network for multi-label image recognition,” 2020. [Online]. Available: <https://arxiv.org/abs/2012.02994>
- [67] E. Ben-Baruch, T. Ridnik, N. Zamir, A. Noy, I. Friedman, M. Protter, and L. Zelnik-Manor, “Asymmetric loss for multi-label classification,” 2021. [Online]. Available: <https://arxiv.org/abs/2009.14119>
- [68] Z.-H. Zhou, M.-L. Zhang, S.-J. Huang, and Y.-F. Li, “Multi-instance multi-label learning,” *Artificial Intelligence*, vol. 176, no. 1, pp. 2291–2320, 2012.
- [69] M.-L. Zhang, “A k-nearest neighbor based multi-instance multi-label learning algorithm,” in *2010 22nd IEEE international conference on tools with artificial intelligence*, vol. 2. IEEE, 2010, pp. 207–212.
- [70] Z.-H. Zhou and M.-L. Zhang, “Multi-instance multi-label learning with application to scene classification,” in *Advances in neural information processing systems*, 2006, pp. 1609–1616.
- [71] S.-J. Huang, W. Gao, and Z.-H. Zhou, “Fast multi-instance multi-label learning,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 11, pp. 2614–2627, 2018.
- [72] S. Liu, L. Zhang, X. Yang, H. Su, and J. Zhu, “Query2label: A simple transformer way to multi-label classification,” 2021. [Online]. Available: <https://arxiv.org/abs/2107.10834>
- [73] T. Ridnik, G. Sharir, A. Ben-Cohen, E. Ben-Baruch, and A. Noy, “MI-decoder: Scalable and versatile classification head,” 2021. [Online]. Available: <https://arxiv.org/abs/2111.12933>
- [74] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, “Cbam: Convolutional block attention module,” 2018. [Online]. Available: <https://arxiv.org/abs/1807.06521>
- [75] R. Liu, J. Huang, T. H. Li, and G. Li, “Causality compensated attention for contextual biased visual recognition,” in *The Eleventh International Conference on Learning Representations*, 2023. [Online]. Available: <https://openreview.net/forum?id=8XqDnrmZQNF>
- [76] J. Zeng, X. Gao, L. Gao, Y. Yu, L. Shen, and X. Pan, “Recognition of rare antinuclear antibody patterns based on a novel attention-based enhancement framework,” *Briefings in Bioinformatics*, vol. 25, no. 2, p. bbad531, Jan. 2024.
- [77] Q. Xie, P. Chen, Z. Li, and R. Xie, “Automatic segmentation and classification for antinuclear antibody images based on deep learning,” *Computational Intelligence and Neuroscience*, vol. 2023, no. 1, p. 1353965, 2023.
- [78] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” 2015. [Online]. Available: <https://arxiv.org/abs/1409.1556>
- [79] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” 2015. [Online]. Available: <https://arxiv.org/abs/1512.00567>
- [80] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” 2015. [Online]. Available: <https://arxiv.org/abs/1512.03385>
- [81] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, “Cvt: Introducing convolutions to vision transformers,” 2021. [Online]. Available: <https://arxiv.org/abs/2103.15808>
- [82] T. Ridnik, H. Lawen, A. Noy, E. B. Baruch, G. Sharir, and I. Friedman, “Tresnet: High performance gpu-dedicated architecture,” 2020. [Online]. Available: <https://arxiv.org/abs/2003.13630>
- [83] A. Shrivastava, A. Gupta, and R. Girshick, “Training region-based object detectors with online hard example mining,” 2016. [Online]. Available: <https://arxiv.org/abs/1604.03540>
- [84] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-cam: Visual explanations from deep networks via gradient-based localization,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 618–626.