
Navigating the MIL Trade-Off: Flexible Pooling for Whole Slide Image Classification

Hossein Jafarinia¹, Danial Hamdi^{1*}, Amirhossein Alamdar^{1*}, Elahe Zahiri²
Soroush Vafaie Tabar¹, Alireza Alipanah¹, Nahal Mirzaie¹, Saeed Razavi¹
Amir Najafi¹, Mohammad Hossein Rohban¹

¹Computer Engineering Department, Sharif University of Technology

²Department of Mathematical Sciences, Sharif University of Technology

{jafarinia, amirhossein.alamdar, elahe.zahiri,
soroush.vafaie96, alireza.alipanah46, nahal.mirzaie,
saeed.razavi, amir.najafi, rohban}@sharif.edu
danial.hamdi@outlook.com

Abstract

Multiple Instance Learning (MIL) is a standard weakly supervised approach for Whole Slide Image (WSI) classification, where performance hinges on both feature representation and MIL pooling strategies. Recent research has predominantly focused on Transformer-based architectures adapted for WSIs. However, we argue that this trend faces a fundamental limitation: data scarcity. In typical settings, Transformer models yield only marginal gains without access to large-scale datasets—resources that are virtually inaccessible to all but a few well-funded research labs. Motivated by this, we revisit simple, non-attention MIL with unsupervised slide features and analyze temperature- β -controlled log-sum-exp (LSE) pooling. For slides partitioned into N patches, we theoretically show that LSE has a smooth transition at a critical $\beta_{\text{crit}} = \mathcal{O}(\log N)$ threshold, interpolating between mean-like aggregation (stable, better generalization but less sensitive) and max-like aggregation (more sensitive but looser generalization bounds). Grounded in this analysis, we introduce Maxsoft—a novel MIL pooling function that enables flexible control over this trade-off, allowing adaptation to specific tasks and datasets. To further tackle real-world deployment challenges such as specimen heterogeneity, we propose PerPatch augmentation—a simple yet effective technique that enhances model robustness. Empirically, Maxsoft achieves state-of-the-art performance in low-data regimes across four major benchmarks (CAMELYON16, CAMELYON17, TCGA-Lung, and SICAP-MIL), often matching or surpassing large-scale foundation models. When combined with PerPatch augmentation, this performance is further improved through increased robustness. Code is available at <https://github.com/jafarinia/maxsoft>

1 Introduction

Whole Slide Image (WSI) analysis using machine learning holds significant promise for supporting the complex and labor-intensive workflow of pathologists [1–3]. However, the extremely large and variable size of WSIs—typically on the order of $150,000 \times 150,000$ pixels—renders the direct application of standard computer vision models, such as Vision Transformers (ViTs) [4], infeasible. The prevailing solution is to divide WSIs into smaller patches and apply a weakly supervised framework known as Multiple Instance Learning (MIL). MIL enables joint modeling of slide-level and patch-level predictions by decomposing the task into a patch-level representation encoder followed by a pooling function that aggregates patch features for WSI-level classification [5, 6].

*Equal contribution.

Following the success of the attention mechanism [7] and the introduction of the Transformer architecture [8], attention-based methods such as ABMIL were adopted in the context of pathology MIL, leading to modest performance gains [5]. In parallel, the introduction of self-supervised learning (SSL)—beginning with SimCLR [9] and its application in DSMIL [10]—enabled training of encoders from scratch on domain-specific patches, resulting in substantial improvements. Since then, much of the research has centered on pairing powerful SSL-trained encoders with increasingly complex Transformer-based architectures to push state-of-the-art (SoTA). However, these works often lack thorough analysis of performance attribution and frequently credit improvements to architectural complexity without sufficient empirical justification [11–16].

In this paper, we demonstrate that Transformer-based MIL methods yield only marginal gains in low-data settings (Figure 3), and once a strong representation encoder is in place, Transformer architectures offer little to no added benefit (see Tables 1 and 3 when the Encoder is Prov-GigaPath [17]). Drawing on our comprehensive encoder experiments in Tables 1 and 3 and the data-quality analysis in Appendix N, we attribute this effect to their well-documented reliance on large-scale training data [4, 15, 18–31]. Unfortunately, most publicly available WSI datasets are relatively small [18, 19], comprising only a few hundred to a few tens of thousands of slides—insufficient for Transformer training (WSI-level, not patch-level). To our knowledge, the largest existing dataset contains approximately 200,000 WSIs [17], is not publicly accessible, and is still arguably too small for training Transformer models. Moreover, the computational resources required to process such datasets are concentrated in a few well-resourced laboratories, making this direction impractical for most groups working on WSI classification.

This motivates a reconsideration of classical (non-attention-based) MIL pooling strategies such as **max pooling** and **mean pooling**. In particular, we analyze **log-sum-exp (LSE) pooling**, a temperature-parameterized function:

$$\text{LSE}_\beta(q_1, \dots, q_N) \triangleq \frac{1}{\beta} \log \left(\frac{1}{N} \sum_{i=1}^N e^{\beta q_i} \right),$$

defined for real-valued inputs q_1, \dots, q_N with an adjustable parameter $\beta \geq 0$ (inverse temperature). The LSE function smoothly interpolates between mean and max behavior: when $\beta \ll 1$, it approximates the mean; when $\beta \gg 1$, it approaches the maximum. Through a combination of theoretical and empirical analyses, we observe smooth phase transitions in model behavior as β crosses a critical threshold of order $\mathcal{O}(\log N)$, where N denotes the number of patches. In the *small- β regime* (i.e., $\beta \ll \mathcal{O}(\log N)$), the model exhibits improved generalization but reduced sensitivity (see Theorem 1). In contrast, in the *large- β regime* (i.e., $\beta \gg \mathcal{O}(\log N)$), sensitivity increases at the cost of looser generalization bounds (see Theorem 2) and less stable training. This trade-off reveals a form of *Pareto optimality* in the temperature parameter, enabling task-specific tuning to balance generalization and sensitivity. Additionally, while max pooling aligns well with the MIL inductive bias, its non-differentiability and high gradient variance make it challenging to optimize, often resulting in unwanted test-time performance fluctuations (see Tables 1 and 3).

Motivated by these observations and insights from Backward Pass Differentiable Approximation (BPDA) [32], we propose **Maxsoft pooling**—a novel MIL pooling strategy that applies LSE with $\beta = \infty$ (i.e., max) during the forward pass, and the gradient of LSE with a moderate β during the backward pass. Conceptually, this hard-forward/soft-backward design is closely related to straight-through estimators [33] (e.g., Straight-Through Gumbel-Softmax) [34]. This design improves optimization stability, generalization, and inference-time sensitivity. Across four major pathology datasets and five standard MIL benchmarks, Maxsoft consistently demonstrates strong performance. On the challenging CAMELYON17 dataset [35], we achieve a WSI-level AUC of 1.0 and a patch-level AUC of 0.93—to our knowledge, the best reported results to date in WSI classification.

From a different perspective, real-world deployment presents challenges such as staining variability, artifacts introduced during slide preparation, and acquisition noise. To enhance robustness under these conditions, we explored data augmentation strategies that reflect such clinically relevant variations. We found that previously proposed WSI-specific augmentations offer negligible or no benefit, as they fail to introduce meaningful or diverse transformations. To address this, we also introduce **PerPatch augmentation**—a simple yet previously unexplored technique that leverages the inherent patch-level granularity of WSIs to increase training diversity in MIL. We observed that combining PerPatch with existing methodologies further improves results in a measurable and consistent manner, and also reduces performance fluctuations.

In summary, our main contributions are:

- A theoretical and empirical analysis of LSE pooling, motivated by its ability to interpolate between max and mean behaviors through a tunable temperature parameter. In particular, we establish several smooth and competing phase transitions in model behavior, revealing trade-offs between generalization and sensitivity.
- The introduction of **Maxsoft**, a novel MIL pooling method derived from our analysis, which demonstrates strong performance, especially in low-data regimes.
- The proposal of **PerPatch augmentation**, a simple yet effective data augmentation strategy that leverages the patch-level structure of WSIs to substantially increase data diversity and yield measurable performance improvements.

2 Related Work

2.1 Transformer-based Architectures

Motivated by the success of the attention mechanism [7], the Transformer architecture [8], and biological insights suggesting that spatial relationships between regions in a WSI influence diagnostic outcomes, recent work has predominantly focused on Transformer-based MIL pooling architectures to capture inter-patch dependencies. WSIs typically consist of a large number of patches (e.g., around 20,000 per slide), making the naive application of self-attention—with its $\mathcal{O}(N^2)$ complexity—ineffective due to prohibitive GPU memory requirements. Among these approaches, ABMIL [5] was the first to introduce attention mechanisms in this context, while DSMIL [10] and Snuffy [15] proposed *sparse* self-attention variants. Notably, Snuffy provided theoretical justification that their sparse attention mechanism can approximate full self-attention under certain conditions [15].

Beyond Transformers, there are graph-based methods, which are hampered by graph construction and nucleus segmentation, and many graph-ViT hybrids collapse to self-attention—effectively Transformer MIL (e.g., GTP) [13, 36–39] and prototype-based approaches, which are also emerging and are sometimes not strictly MIL (e.g., PANTHER) [40, 41].

2.2 Classical MIL Pooling Functions

In addition to max and mean, several other pooling functions have been used, such as noisy-or [42] (any single cause), noisy-and [43] (sufficient proportion of causes), Integrated Segmentation and Recognition (ISR) [44] (a smooth OR with evidence accumulation), and smooth approximations to max/mean such as LSE [45], generalized mean (GM), and smoothmax [46].

2.3 Pathology Foundation Models

Recent efforts to develop Foundation Models for pathology [17, 47–49] have largely focused on training representation encoders using massive numbers of patches extracted from as many WSIs as possible. For instance, UNI [47] is trained on 100 million patches, while Prov-GigaPath [17] uses 1.3 billion patches from 200,000 WSIs—both based on DINOv2 [50]. These approaches yield highly robust encoders that can boost performance across various MIL pooling strategies. However, building such models faces key challenges: limited public data, high computational cost, lack of flexibility compared to generative foundation models like LLMs [51, 52], and sensitivity to training errors such as poor hyperparameter choices. Consequently, progress remains confined to a few well-resourced labs. Complementing these efforts, the recently proposed R²T-MIL targets foundation model-level performance in low-data regimes via online feature re-embedding during MIL pooling training.

2.4 Augmentation Methods

WSI-level Augmentation modifies slide-level representations by selecting a fixed number of representative patches, inspired by Mixup [53] and Mask Augmentation [54]. Methods such as ReMix [55], RankMix [56], PseMix [19], MixupMIL [57], and Attention-Guided Mixup [58] apply various mixing strategies after aggressive patch downsampling (e.g., using centroids, ranked attention scores, or pseudo-bags). This leads to substantial information loss and degraded performance, often

requiring teacher-student distillation for limited gains (see Tables 2 and 4). Additionally, many such augmentations lack realism.

Patch-level Augmentation methods apply standard image transformations—such as random rotation, color jitter, and Hematoxylin-Eosin-DAB (HED) jitter—directly to patches prior to feature extraction. These augmentations are typically inspired by patch classification research [59, 60]. To increase efficiency, EMBAUGMENTER [61] and AugDiff [62] employ latent generative models to simulate augmented patch representations during MIL training. SSRDL [63] builds on this by introducing a DINO-based [64] framework that performs online sampling from a learned distribution of augmented features and reuses parts of the model during MIL pooling training. While these methods improve efficiency and robustness, they often suffer from limited diversity—e.g., applying the same augmentation uniformly across patches (EMBAUGMENTER, AugDiff)—and reduced interpretability, especially in SSRDL, where the specific form of augmentation is not directly observable.

3 Background: MIL Formulation

In a binary image classification setting, the dataset $D = \{(I_1, y_1), \dots, (I_n, y_n)\}$ consists of images (or equivalently *bags of patches*) I_i , where each bag $I_i = \{\mathbf{X}_1^{(i)}, \dots, \mathbf{X}_N^{(i)}\}$ contains N corresponding instances. Each bag label $y_i \in \{0, 1\}$ is determined by its individual instance labels $\{y_1^{(i)}, \dots, y_N^{(i)}\}$, where $y_j^{(i)} \in \{0, 1\}$ are unknown during training. Under the standard MIL assumption, we have

$$y_i = \begin{cases} 0, & \text{if } \sum_{j=1}^N y_j^{(i)} = 0, \\ 1, & \text{otherwise,} \end{cases} = \max_{j=1, \dots, N} \{y_j^{(i)}\} \quad \forall i = 1, \dots, n.$$

Thus, a bag I_i is labeled positive if at least one of its instances is labeled positive; otherwise, it is labeled negative. The MIL model is trained by optimizing the log-likelihood function:

$$P(y|I) = \phi(I)^y (1 - \phi(I))^{1-y},$$

where $\phi(I) = \phi(\mathbf{X}_1, \dots, \mathbf{X}_N) \in [0, 1]$ represents the predicted probability of ($y = 1$) given the bag I . Since MIL assumes no ordering or dependency among instances, $\phi(I)$ must be a permutation-invariant function. This is ensured by the *fundamental theorem of symmetric functions with monomials* [65] and a similar result by [66], which states that for any permutation-invariant function ϕ satisfying Hausdorff continuity, there exist functions ψ and $\Phi : \mathbb{R} \rightarrow \mathbb{R}$, and a permutation-invariant function $\pi : \mathbb{R}^N \rightarrow \mathbb{R}$ such that

$$\phi(I) = \Phi(\pi(\psi(\mathbf{X}_1), \dots, \psi(\mathbf{X}_N))). \quad (1)$$

Here, ψ and Φ are continuous transformations, and π is a permutation-invariant function (such as sum or max). MIL methods typically follow two primary approaches: i) **Instance-level approach**, where ψ is an instance classifier and Φ is the identity function. ii) **Embedding-level approach**, where ψ is a feature extractor, and Φ maps the extracted features to a bag classification score.

In Deep MIL, ψ typically uses pre-trained vision backbones to extract features from bag instances [10–12, 15, 18, 67, 68]. The aggregation function π ranges from non-parametric methods like max pooling to parametric ones like attention mechanisms, as detailed in Section 2.1. For multi-class and multi-task classification, Φ 's output dimension is adjusted to the number of classes [69–71].

4 Our Method

While SoTA MIL pooling architectures focus on embedding-level methods, this work emphasizes instance-level approaches. In clinical practice, pathologists prioritize patch-level predictions as they allow validation of the model's reasoning. As a result, instance-level methods offer full interpretability. Although max pooling and mean pooling can be applied at both embedding and instance levels, we focus exclusively on their instance-level variants throughout this work (i.e., Φ in (1) is identity).

We begin by following the standard MIL pooling procedure: each image in D is divided into N non-overlapping patches $\mathbf{X}_1, \dots, \mathbf{X}_N$. We model the function $\psi(\cdot)$ in (1) by passing each patch through a fixed, pre-trained ViT, followed by a trainable fully-connected MLP, and concluding with a sigmoid function in order to compute $q_i \triangleq \psi(\mathbf{X}_i) \in [0, 1]$ for $i = 1, \dots, N$ (see Section 5 for more details). The next subsection establishes a theoretical foundation for the trade-offs associated with the choice of aggregation function $\pi(\cdot)$.

4.1 Theoretical Analysis

Assume, instead of the standard max/mean pooling, we choose a soft LSE $_{\beta}$ function for $\pi(\cdot)$ in (1). We theoretically demonstrate that smaller values of β generally lead to better generalization and a more stable training process. Conversely, larger values of β are more suitable for reducing training error and increasing the model’s sensitivity in real-world scenarios, where only a few patches may contain malignant patterns, and the model must detect and respond to them effectively. A full version of this section is provided in Appendix A, with a summary of informal results here:

Theorem 1 ((Informal) Effect of β on Generalization Error). *Assuming mild conditions—such as Lipschitz continuity of the loss function, i.i.d. sampling of images in dataset D , and local statistical dependence among image patches—we show that the generalization error undergoes a smooth phase transition as β increases from below $\mathcal{O}(\log N)$ to above. In the small- β regime, the generalization gap can be bounded as $\mathcal{O}((nN)^{-1/2} + n^{-1})$, while in the large- β regime, it can be as large as $\Omega(n^{-1/2})$.*

The formal version of this theorem (Theorem 4), including precise bounds, threshold conditions, and constants, is provided in Appendix A, along with a complete proof. Next, we focus on sensitivity analysis. Notably, in contrast to generalization—which favors smaller values of β —high sensitivity in noisy regimes requires sufficiently large β .

Theorem 2 ((Informal) Effect of β on True Positive / False Negative Rate). *Assume a general statistical model for the distribution of $\psi(\mathbf{X}_i)$ s on each patch, and suppose the number of cancer-indicative patches (i.e., patches exhibiting malignant patterns) is, with high probability, substantially smaller than N . Then, for any $\alpha \in (0, 1)$, if $\beta \geq \mathcal{O}(\log(N/\alpha))$, the true positive rate of classifier LSE $_{\beta}(\psi(\mathbf{X}_1), \dots, \psi(\mathbf{X}_N))$ is at least $1 - \alpha$. Conversely, if β remains $\mathcal{O}(1)$ with respect to N , the false negative rate is at least $1/2$.*

Again, the formal version of this theorem (Theorem 6), along with a complete proof and further explanations, is provided in Appendix A. Specifically, we show that choosing $\beta \gg \mathcal{O}(\log N)$ is both necessary and sufficient to achieve a high true positive rate and to mitigate false negatives. The key conclusion of this section is that all values of $\beta \in (0, +\infty)$ are *Pareto-optimal*: depending on which aspect—generalization or sensitivity—is prioritized, and how the trade-off is weighted, the optimal choice of β will vary. This insight supports the central idea behind our method: using smaller values of β during training to promote generalization, while selecting larger values (possibly even ∞) at inference time to enhance sensitivity.

4.2 Maxsoft pooling

High gradient variance can lead to instability during training, slowing convergence and requiring more iterations to reach a desired accuracy [72–74]. Thus, reducing variance improves stability, enhances generalization, and promotes more reliable learning dynamics [75–77]. This underscores the need to regulate gradient variance to balance stability with effective optimization. max pooling and mean pooling, the two dominant MIL pooling strategies, represent opposite ends of this spectrum: max pooling excels in instance discrimination but suffers from instability and worse generalization due to its non-smooth, non-differentiable nature, while mean pooling is stable, fully differentiable, and robust for bag classification but weaker at distinguishing instances. We provide a unified theoretical explanation for these trade-offs in Section 4.1.

Leveraging this insight, we introduce Maxsoft pooling (see Figure 1), a novel strategy that combines the stability and smoothness of mean pooling with the discriminative power of max pooling. In the forward pass, we use LSE $_{\infty}(q_{1:N}) = \max_i q_i$, while in the backward pass we approximate the gradient of the max operator by $\nabla \text{LSE}_{\beta}(q_{1:N})$ for some moderate $\beta < \infty$. For more details on the learning procedure, see Algorithm 1. Maxsoft pooling retains the exact behavior of max pooling during the forward pass, but it employs a differentiable function in the backward pass by replacing the max operator with the LSE. This ensures smoother gradient computation, improved stability, and reduced generalization gap (see [32], and Theorem 1).

4.3 PerPatch Augmentation

Existing WSI-level augmentation methods, as outlined in Section 2.4, often yield marginal gains or even degrade performance [19, 55, 62, 63]. Therefore, we shift focus to patch-level augmentation,

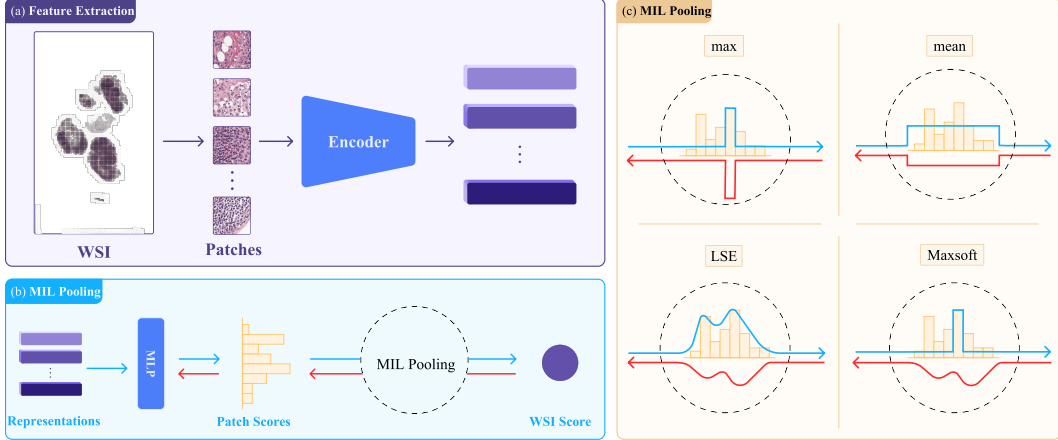


Figure 1: Overview of MIL (a) The WSIs are segmented into patches, followed by embedding extraction via a pre-trained encoder. (b) Embeddings are fed into MIL pooling for patch. (c) mean, max, LSE, and Maxsoft with their respective forward and backward behaviors.

which introduces more meaningful variability. Latent generative models [61, 62] add computational overhead without consistent benefits (Tables 2, 4). Instead, we precompute multiple augmented versions of each patch and employ the PerPatch sampling strategy, enabling a substantial increase in WSI diversity during MIL pooling training.

PerPatch augmentation differs from PerSlide by sampling each patch independently—either from its original version or one of its augmented variants—instead of applying the same augmentation transformation across all patches in a single WSI (see Figure 5). Specifically, given a WSI I , we generate multiple augmented variants $I^{(k)} = \{\mathbf{X}_1^{(k)}, \dots, \mathbf{X}_N^{(k)}\}$, $k = 0, 1, \dots, m$, where $I^{(0)} = I$ is the original WSI and m denotes the number of augmentations. We then construct an augmented WSI \hat{I} by independently sampling each patch from its available versions: $\hat{I} = \{\hat{\mathbf{X}}_1, \dots, \hat{\mathbf{X}}_N\}$ with $\hat{\mathbf{X}}_j = \mathbf{X}_j^{(k)}$, $k \sim \text{Uniform}(\{0, 1, 2, \dots, m\})$. Thus, each patch $\hat{\mathbf{X}}_j$ is independently sampled from either the original version ($k = 0$) or one of its m augmented variants ($1 \leq k \leq m$). For a visual and procedural comparison between PerSlide and PerPatch augmentations, refer to Figure 5 and Algorithms 2 and 3. We argue that PerPatch creates significantly more diverse augmentations during MIL pooling training compared to PerSlide, and this may explain why the method consistently improves results, as shown in Tables 2 and 4.

5 Experiments and Results

Our major experiments are on major WSI classification datasets. Analyses of classical MIL pooling functions are in Appendix O. Although out-of-distribution (OOD) generalization is beyond our scope, we include such experiments in Appendix Q. We also evaluate on five canonical MIL datasets, reported in Appendix R.

5.1 Datasets

We evaluate on four large-scale pathology WSI datasets—CAMELYON16, CAMELYON17, TCGA-Lung and SICAP-MIL [78, 35, 79, 80]. CAMELYON16 contains 270 training / 129 test slides; CAMELYON17 has 1,000 slides across four tumor-type labels (normal, macro, micro, ITC); TCGA-Lung spans 1,042 slides split between LUAD and LUSC; SICAP-MIL comprises 349 slides with normal/abnormal labels and primary/secondary Gleason grades. Further dataset statistics and split protocols are given in Appendix E.

5.2 Experimental Setup

All WSIs are tiled into 256×256 patches at $20\times$ magnification, discarding background. For CAMELYON16 we use the official split; CAMELYON17, TCGA-Lung and SICAP-MIL are each partitioned

into 60% train, 15% val, 25% test (with SICAP-MIL’s official split). Details on train/val/test allocations and tiling thresholds appear in Appendix F.

5.3 Evaluation Metrics

We report slide-level AUC and accuracy, patch-level AUC and F1 (omitting patch accuracy due to class imbalance), plus F1 and Expected Calibration Error (ECE) to assess confidence calibration in cancer detection. Formal definitions and details for ECE are in Appendix G.

5.4 Models and Implementation Details

Encoders: We compare domain-specific and natural-image backbones including DINO Domain (ViT-S/16) trained on around a few million in-domain patches, DINO Natural (ViT-S/16) [64] pre-trained on ImageNet-1K, UNI (ViT-L/16) [47] with 100 million pathology patches, Prov-GigaPath (ViT-G/14) [17] with 1.3 billion patches, ResNet-50 [81] pre-trained on ImageNet-1K, PLIP [82], a CLIP-based model fine-tuned on around 208 thousand text-annotated pathology patches, and SSRDL-trained ViTs [63] via its Representation Augmentation Module on domain patches.

MIL-Pooling: We evaluate max, mean, LSE, ABMIL [5], DSMIL [10], Snuffy [15], R²T-MIL [16] and our Maxsoft.

Miscellaneous: For completeness, we also evaluate GTP [37] and PANTHER [41] as non-MIL WSI classifiers.

Augmentations: We compare standard WSI-level baselines—with basic augmentations and simple per-slide transforms (Random Rotation, Gaussian Blur, and Color Jitter; see Appendix I)—and prior augmentation strategies (ReMix [55], RankMix [56], AugDiff [62], and SSRDL’s latent augmentations [63]) against PerPatch, our patch-level method. Hyperparameters, optimizer settings, learning rates, and batch sizes are provided in Appendix H.

5.5 Results and Analysis

Across all datasets, stronger encoders yield better performance: DINO Natural < DINO Domain < UNI < Prov-GigaPath. With Prov-GigaPath, every MIL pooling strategy attains high scores (Tables 1 and 3). For visualizations of the resulting representations, see Appendix K.

Max pooling shows high variability, whereas mean pooling offers lower variance and on average better generalization. Guided by our theoretical analysis, we select β via validation: $\beta = 5$ for CAMELYON16/17, $\beta = 10$ for TCGA-Lung, and $\beta = 3.5$ for SICAP-MIL. This choice reflects the datasets’ characteristics—sparser positives in CAMELYON demand higher sensitivity, while more frequent positives in TCGA-Lung and SICAP-MIL favor generalization.

LSE matches or outperforms both Transformer-based and non-Transformer-based methods (ABMIL, DSMIL, Snuffy, R²T-MIL, GTP, and PANTHER) despite its simplicity. Maxsoft attains the highest overall performance and substantially lower variance than max pooling, validating the use of $\beta = \infty$ in the forward pass with a moderate β for gradients. These results indicate that, at current WSI dataset sizes, encoder quality dominates architectural complexity and attention-based pooling yields limited gains.

On augmentation (Tables 2 and 4), PerSlide yields inconsistent or negative effects, whereas PerPatch uniformly reduces variance and often improves accuracy, indicating robustness. WSI-level methods like ReMix and RankMix underperform due to aggressive patch downsampling and because DINO [64] embeddings lack augmentation–interpolation consistency between the image and feature spaces. Among prior approaches, only AugDiff and, at times, SSRDL show improvements—likely because, despite the issues above, they still produce diverse, non-redundant feature augmentations. SSRDL [63] offers some gains on TCGA-Lung but at high cost, since each new augmentation requires retraining a self-supervised model from scratch. Overall, PerPatch delivers the largest gains through richer, more varied augmentations. For analysis of special cases, see Appendix L.2.

Notably, Maxsoft with DINO Domain surpasses Prov-GigaPath in accuracy while requiring $0.001\times$ less computation and having $0.0026\times$ fewer parameters. In terms of augmentation speed, AugDiff is

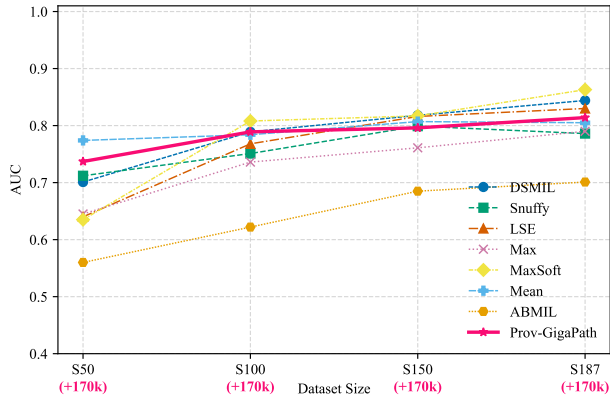


Figure 3: Performance of major Transformer-based and classical MIL pooling methods on the SICAP-MIL dataset across different training set sizes.

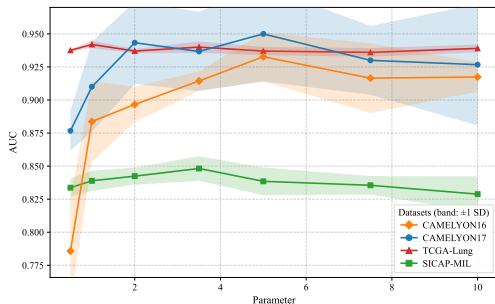


Figure 4: The impact of β in Maxsoft architecture on CAMELYON16, CAMELYON17, TCGA-Lung, and SICAP-MIL datasets.

Augmentation	Slide		Patch
	AUC	ECE	AUC
No Aug	0.983 _{.055}	0.121 _{.051}	0.839 _{.019}
Random Rotation	0.965 _{.021}	0.133 _{.060}	0.854 _{.020}
Random Elastic Deformation	0.890 _{.028}	0.193 _{.003}	0.861 _{.031}
Random Affine Transformation	0.900 _{.113}	0.193 _{.102}	0.822 _{.020}
Random Gaussian Blurring	0.995 _{.007}	0.081 _{.002}	0.882 _{.011}
Random Color Jitter	0.955 _{.064}	0.084 _{.001}	0.876 _{.045}
Random HED Jitter	0.955 _{.050}	0.156 _{.047}	0.874 _{.039}

Table 5: Result of each base augmentation on the CAMELYON17 dataset, reported as mean_{.std}.

Impact of Base Augmentations. To evaluate the augmentations from Section 5.4, we ablated them individually. We also examined the spatial transforms Random Elastic Deformation and Random Affine Transformation. As shown in Table 5, these spatial operations consistently reduced performance. We hypothesize that such transformations fail to capture clinical variability and primarily inject noise rather than useful diversity.

7 Conclusion and Discussion

This work highlights the value of instance-level representation learning under current data constraints in WSI diagnosis. In response, we present a theoretical and empirical analysis of classical MIL pooling and introduce Maxsoft—a simple, resource-efficient pooling function with strong, adaptable performance. A limitation is that, despite clinically acceptable visualizations, patch-level performance lags behind WSI-level accuracy. Another weakness is that, while our theory offers numerical guidance for selecting β in Maxsoft, a small dataset-specific search is still required; future work should make this selection tuning-free and integrated into training. Moreover, standard augmentations (e.g., Random Affine, Elastic Deformation) do not capture real distribution shifts, underscoring the need for more realistic techniques. Finally, we advocate studying cancer detection via anomaly detection [83–88], especially in data-scarce regimes, since it naturally targets rare or novel abnormalities without large labeled cancer sets [89–92].

Acknowledgments

We thank Mehrab Moradzadeh, Mohammad Mosayyebi, Arian Komaei Koma, Amir Hossein Saberi, Mohammad Azizmalayeri, Hossein Mirzaei, Hosein Hasani, and the anonymous reviewers for their helpful discussions and feedback on this work.

References

- [1] Andrew H Song, Guillaume Jaume, Drew FK Williamson, Ming Y Lu, Anurag Vaidya, Tiffany R Miller, and Faisal Mahmood. Artificial intelligence for digital and computational pathology. *Nature Reviews Bioengineering*, 1(12):930–949, 2023.
- [2] Jana Lipkova, Richard J Chen, Bowen Chen, Ming Y Lu, Matteo Barbieri, Daniel Shao, Anurag J Vaidya, Chengkuan Chen, Luoting Zhuang, Drew FK Williamson, et al. Artificial intelligence for multimodal data integration in oncology. *Cancer cell*, 40(10):1095–1110, 2022.
- [3] Artem Shmatko, Narmin Ghaffari Laleh, Moritz Gerstung, and Jakob Nikolas Kather. Artificial intelligence in histopathology: enhancing cancer research and clinical oncology. *Nature cancer*, 3(9):1026–1038, 2022.
- [4] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [5] Maximilian Ilse, Jakub Tomczak, and Max Welling. Attention-based deep multiple instance learning. In *International conference on machine learning*, pages 2127–2136. PMLR, 2018.
- [6] Jaeseok Jang and Hyuk-Yoon Kwon. Are multiple instance learning algorithms learnable for instances? *Advances in Neural Information Processing Systems*, 37:10575–10612, 2024.
- [7] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [8] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PmLR, 2020.
- [10] Bin Li, Yin Li, and Kevin W Eliceiri. Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14318–14328, 2021.
- [11] Linhao Qu, Xiaoyuan Luo, Shaolei Liu, Manning Wang, and Zhijian Song. Dgmil: Distribution guided multiple instance learning for whole slide image classification. In *International conference on medical image computing and computer-assisted intervention*, pages 24–34. Springer, 2022.
- [12] Yushan Zheng, Jun Li, Jun Shi, Fengying Xie, Jianguo Huai, Ming Cao, and Zhiguo Jiang. Kernel attention transformer for histopathology whole slide image analysis and assistant cancer diagnosis. *IEEE Transactions on Medical Imaging*, 42(9):2726–2739, 2023.
- [13] Gianpaolo Bontempo, Angelo Porrello, Federico Bolelli, Simone Calderara, and Elisa Ficarra. Das-mil: distilling across scales for mil classification of histological wsis. In *International conference on medical image computing and computer-assisted intervention*, pages 248–258. Springer, 2023.
- [14] Olga Fourkioti, Matt De Vries, Chen Jin, Daniel C Alexander, and Chris Bakal. Camil: Context-aware multiple instance learning for cancer detection and subtyping in whole slide images. *arXiv preprint arXiv:2305.05314*, 2023.
- [15] Hossein Jafarinia, Alireza Alipanah, Danial Hamdi, Saeed Razavi, Nahal Mirzaie, and Mohammad Hossein Rohban. Snuffy: Efficient whole slide image classifier. *arXiv e-prints*, pages arXiv–2408, 2024.
- [16] Wenhao Tang, Fengtao Zhou, Sheng Huang, Xiang Zhu, Yi Zhang, and Bo Liu. Feature re-embedding: Towards foundation model-level performance in computational pathology. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11343–11352, 2024.

- [17] Hanwen Xu, Naoto Usuyama, Jaspreet Bagga, Sheng Zhang, Rajesh Rao, Tristan Naumann, Cliff Wong, Zelalem Gero, Javier González, Yu Gu, et al. A whole-slide foundation model for digital pathology from real-world data. *Nature*, 630(8015):181–188, 2024.
- [18] Hongrun Zhang, Yanda Meng, Yitian Zhao, Yihong Qiao, Xiaoyun Yang, Sarah E Coupland, and Yalin Zheng. Dtf-d-mil: Double-tier feature distillation multiple instance learning for histopathology whole slide image classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18802–18812, 2022.
- [19] Pei Liu, Luping Ji, Xinyu Zhang, and Feng Ye. Pseudo-bag mixup augmentation for multiple instance learning-based whole slide image classification. *IEEE Transactions on Medical Imaging*, 43(5):1841–1852, 2024.
- [20] Lucas Beyer, Xiaohua Zhai, and Alexander Kolesnikov. Better plain vit baselines for imagenet-1k. *arXiv preprint arXiv:2205.01580*, 2022.
- [21] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in neural information processing systems*, 34:24261–24272, 2021.
- [22] Peng Li, Lu Huang, Jin Li, Haiyan Yan, and Dongjing Shan. Graph-based vision transformer with sparsity for training on small datasets from scratch. *Scientific Reports*, 15(1):24520, 2025.
- [23] Namuk Park and Songkuk Kim. How do vision transformers work? *arXiv preprint arXiv:2202.06709*, 2022.
- [24] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? *Advances in neural information processing systems*, 34:12116–12128, 2021.
- [25] Zhiying Lu, Hongtao Xie, Chuanbin Liu, and Yongdong Zhang. Bridging the gap between vision transformers and convolutional neural networks on small datasets. *Advances in Neural Information Processing Systems*, 35:14663–14677, 2022.
- [26] Yahui Liu, Enver Sangineto, Wei Bi, Nicu Sebe, Bruno Lepri, and Marco Nadai. Efficient training of visual transformers with small datasets. *Advances in Neural Information Processing Systems*, 34:23818–23830, 2021.
- [27] Ali Hassani, Steven Walton, Nikhil Shah, Abulikemu Abuduweili, Jiachen Li, and Humphrey Shi. Escaping the big data paradigm with compact transformers. *arXiv preprint arXiv:2104.05704*, 2021.
- [28] Seung Hoon Lee, Seunghyun Lee, and Byung Cheol Song. Vision transformer for small-size datasets. *arXiv preprint arXiv:2112.13492*, 2021.
- [29] Ibrahim Batuhan Akkaya, Senthilkumar S Kathiresan, Elahe Arani, and Bahram Zonooz. Enhancing performance of vision transformers on small datasets through local inductive bias incorporation. *Pattern Recognition*, 153:110510, 2024.
- [30] Francisco Carrillo-Perez, Marija Pizurica, Michael G Ozawa, Hannes Vogel, Robert B West, Christina S Kong, Luis Javier Herrera, Jeanne Shen, and Olivier Gevaert. Synthetic whole-slide image tile generation with gene expression profile-infused deep generative models. *Cell Reports Methods*, 3(8), 2023.
- [31] Jun Li, Junyu Chen, Yucheng Tang, Ce Wang, Bennett A Landman, and S Kevin Zhou. Transforming medical imaging with transformers? a comparative review of key properties, current progresses, and future perspectives. *Medical image analysis*, 85:102762, 2023.
- [32] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International conference on machine learning*, pages 274–283. PMLR, 2018.
- [33] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- [34] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- [35] Peter Bandi, Oscar Geessink, Quirine Manson, Marcory Van Dijk, Maschenka Balkenhol, Meyke Hermsen, Babak Ehteshami Bejnordi, Byungjae Lee, Kyunghyun Paeng, Aoxiao Zhong, et al. From detection of individual metastases to classification of lymph node status at the patient level: the camelyon17 challenge. *IEEE transactions on medical imaging*, 38(2):550–560, 2018.

- [36] Siemen Brussee, Giorgio Buzzanca, Anne MR Schrader, and Jesper Kers. Graph neural networks in histopathology: Emerging trends and future directions. *Medical Image Analysis*, page 103444, 2025.
- [37] Yi Zheng, Rushin H Gindra, Emily J Green, Eric J Burks, Margrit Betke, Jennifer E Beane, and Vijaya B Kolachalama. A graph-transformer for whole slide image classification. *IEEE transactions on medical imaging*, 41(11):3003–3015, 2022.
- [38] Chaitanya K Joshi. Transformers are graph neural networks. *arXiv preprint arXiv:2506.22084*, 2025.
- [39] Chen Cai, Truong Son Hy, Rose Yu, and Yusu Wang. On the connection between mpnn and graph transformer. In *International conference on machine learning*, pages 3408–3430. PMLR, 2023.
- [40] Jin-Gang Yu, Zihao Wu, Yu Ming, Shule Deng, Yuanqing Li, Caifeng Ou, Chunjiang He, Baiye Wang, Pusheng Zhang, and Yu Wang. Prototypical multiple instance learning for predicting lymph node metastasis of breast cancer from whole-slide pathological images. *Medical Image Analysis*, 85:102748, 2023.
- [41] Andrew H Song, Richard J Chen, Tong Ding, Drew FK Williamson, Guillaume Jaume, and Faisal Mahmood. Morphological prototyping for unsupervised slide representation learning in computational pathology. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11566–11578, 2024.
- [42] Cha Zhang, John Platt, and Paul Viola. Multiple instance boosting for object detection. *Advances in neural information processing systems*, 18, 2005.
- [43] Oren Z Kraus, Jimmy Lei Ba, and Brendan J Frey. Classifying and segmenting microscopy images with deep multiple instance learning. *Bioinformatics*, 32(12):i52–i59, 2016.
- [44] James Keeler, David Rumelhart, and Wee Leow. Integrated segmentation and recognition of hand-printed numerals. *Advances in neural information processing systems*, 3, 1990.
- [45] Jan Ramon and Luc De Raedt. Multi instance neural networks. In *Proceedings of the ICML-2000 workshop on attribute-value and relational learning*, pages 53–60, 2000.
- [46] Florentin Bieder, Robin Sandkühler, and Philippe C Cattin. Comparison of methods generalizing max-and average-pooling. *arXiv preprint arXiv:2103.01746*, 2021.
- [47] Richard J Chen, Tong Ding, Ming Y Lu, Drew FK Williamson, Guillaume Jaume, Andrew H Song, Bowen Chen, Andrew Zhang, Daniel Shao, Muhammad Shaban, et al. Towards a general-purpose foundation model for computational pathology. *Nature Medicine*, 30(3):850–862, 2024.
- [48] Xiyue Wang, Junhan Zhao, Eliana Marostica, Wei Yuan, Jietian Jin, Jiayu Zhang, Ruijiang Li, Hongping Tang, Kanran Wang, Yu Li, et al. A pathology foundation model for cancer diagnosis and prognosis prediction. *Nature*, 634(8035):970–978, 2024.
- [49] Nanne Aben, Edwin D de Jong, Ioannis Gatopoulos, Nicolas Känzig, Mikhail Karasikov, Axel Lagré, Roman Moser, Joost van Doorn, Fei Tang, et al. Towards large-scale training of pathology foundation models. *arXiv preprint arXiv:2404.15217*, 2024.
- [50] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=a68SUt6zFt>. Featured Certification.
- [51] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [52] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [53] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.

- [54] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with dropout. *arXiv preprint arXiv:1708.04552*, 2017.
- [55] Jiawei Yang, Hanbo Chen, Yu Zhao, Fan Yang, Yao Zhang, Lei He, and Jianhua Yao. Remix: A general and efficient framework for multiple instance learning based whole slide image classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 35–45. Springer, 2022.
- [56] Yuan-Chih Chen and Chun-Shien Lu. Rankmix: Data augmentation for weakly supervised learning of classifying whole slide images with diverse sizes and imbalanced categories. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23936–23945, 2023.
- [57] Michael Gadermayr, Lukas Koller, Maximilian Tschuchnig, Lea Maria Stangassinger, Christina Kreutzer, Sebastien Couillard-Despres, Gertie Janneke Oostingh, and Anton Hittmair. Mixup-mil: Novel data augmentation for multiple instance learning and a study on thyroid cancer diagnosis. In *International conference on medical image computing and computer-assisted intervention*, pages 477–486. Springer, 2023.
- [58] Farchan Hakim Raswa, Chun-Shien Lu, and Jia-Ching Wang. Attention-guided prototype mixing: Diversifying minority context on imbalanced whole slide images classification learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 7624–7633, 2024.
- [59] David Tellez, Geert Litjens, Péter Bándi, Wouter Bulten, John-Melle Bokhorst, Francesco Ciompi, and Jeroen Van Der Laak. Quantifying the effects of data augmentation and stain color normalization in convolutional neural networks for computational pathology. *Medical image analysis*, 58:101544, 2019.
- [60] Khrystyna Faryna, Jeroen Van der Laak, and Geert Litjens. Tailoring automated data augmentation to h&e-stained histopathology. In *Medical imaging with deep learning*, 2021.
- [61] Imaad Zaffar, Guillaume Jaume, Nasir Rajpoot, and Faisal Mahmood. Embedding space augmentation for weakly supervised learning in whole-slide images. In *2023 IEEE 20th International Symposium on Biomedical Imaging (ISBI)*, pages 1–4. IEEE, 2023.
- [62] Liuxi Dai, Yifeng Wang, Haoqian Wang, Yongbing Zhang, et al. Augdiff: Diffusion based feature augmentation for multiple instance learning in whole slide image. *IEEE Transactions on Artificial Intelligence*, 2024.
- [63] Kunming Tang, Zhiguo Jiang, Kun Wu, Jun Shi, Fengying Xie, Wei Wang, Haibo Wu, and Yushan Zheng. Self-supervised representation distribution learning for reliable data augmentation in histopathology wsi classification. *IEEE Transactions on Medical Imaging*, 2024.
- [64] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.
- [65] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. *Advances in neural information processing systems*, 30, 2017.
- [66] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017.
- [67] Ming Y Lu, Drew FK Williamson, Tiffany Y Chen, Richard J Chen, Matteo Barbieri, and Faisal Mahmood. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature biomedical engineering*, 5(6):555–570, 2021.
- [68] Zhuchen Shao, Hao Bian, Yang Chen, Yifeng Wang, Jian Zhang, Xiangyang Ji, et al. Transmil: Transformer based correlated multiple instance learning for whole slide image classification. *Advances in neural information processing systems*, 34:2136–2147, 2021.
- [69] Ladislav Lenc and Pavel Král. Two-level neural network for multi-label document classification. In *Artificial Neural Networks and Machine Learning–ICANN 2017: 26th International Conference on Artificial Neural Networks, Alghero, Italy, September 11–14, 2017, Proceedings, Part II 26*, pages 368–375. Springer, 2017.
- [70] Ladislav Lenc and Pavel Král. Combination of neural networks for multi-label document classification. In *International Conference on Applications of Natural Language to Information Systems*, pages 278–282. Springer, 2017.

- [71] Andrew Maxwell, Runzhi Li, Bei Yang, Heng Weng, Aihua Ou, Huixiao Hong, Zhaoxian Zhou, Ping Gong, and Chaoyang Zhang. Deep learning architectures for multi-label classification of intelligent health risk prediction. *BMC bioinformatics*, 18:121–131, 2017.
- [72] Florian Hübler, Ilyas Fatkhullin, and Niao He. From gradient clipping to normalization for heavy tailed sgd, 2024. URL <https://arxiv.org/abs/2410.13849>.
- [73] Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning, 2018. URL <https://arxiv.org/abs/1606.04838>.
- [74] Fartash Faghri, David Duvenaud, David J. Fleet, and Jimmy Ba. A study of gradient variance in deep learning, 2020. URL <https://arxiv.org/abs/2007.04532>.
- [75] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima, 2017. URL <https://arxiv.org/abs/1609.04836>.
- [76] Dong Yin, Ashwin Pananjady, Max Lam, Dimitris Papailiopoulos, Kannan Ramchandran, and Peter Bartlett. Gradient diversity: a key ingredient for scalable distributed learning, 2018. URL <https://arxiv.org/abs/1706.05699>.
- [77] Dominic Masters and Carlo Luschi. Revisiting small batch training for deep neural networks, 2018. URL <https://arxiv.org/abs/1804.07612>.
- [78] Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes van Diest, Bram van Ginneken, Nico Karssemeijer, Geert J. S. Litjens, Jeroen A. van der Laak, Meyke Hermesen, Quirine F. Manson, Maschenka C. A. Balkenhol, Oscar G. F. Geessink, Nikolaos Stathonikos, Marcory Crf van Dijk, Peter Bult, Francisco Beca, Andrew H. Beck, Dayong Wang, Aditya Khosla, Rishab Gargeya, Humayun Irshad, Aoxiao Zhong, Qi Dou, Quanzheng Li, Hao Chen, Huangjing Lin, Pheng-Ann Heng, Christian Hass, Elia Bruni, Quincy Kwan-Sut Wong, Ugur Halici, Mustafa Ümit Öner, Rengul Cetin-Atalay, Matt Berseth, Vitali Khvatkov, A. Vylegzhanin, Oren Z. Kraus, Muhammad Shaban, Nasir M. Rajpoot, Ruqayya Awan, Korsuk Sirinukunwattana, Talha Qaiser, Yee-Wah Tsang, David Tellez, Jonas Annuscheit, Peter Hufnagl, Mira Valkonen, Kimmo Kartasalo, Leena Latonen, Pekka Ruusuvuori, Kaisa Liimatainen, Shadi Albarqouni, Bharti Mungal, Amitha Anna George, Stefanie Demirci, Nassir Navab, Seiryu Watanabe, Shigeto Seno, Yoichi Takenaka, Hideo Matsuda, Hady Ahmady Phoulady, Vassili A. Kovalev, Alexander Kalinovsky, Vitali Liauchuk, Gloria Bueno, Maria del Milagro Fernández-Carrobles, Ismael Serrano, Oscar Deniz, Daniel Racoceanu, and Rui Venâncio. Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *JAMA*, 318:2199–2210, 2017. URL <https://api.semanticscholar.org/CorpusID:205086555>.
- [79] Lee AD Cooper, Elizabeth G Demicco, Joel H Saltz, Reid T Powell, Arvind Rao, and Alexander J Lazar. Pancancer insights from the cancer genome atlas: the pathologist’s perspective. *The Journal of pathology*, 244(5):512–524, 2018.
- [80] Julio Silva-Rodríguez, Arne Schmidt, María A Sales, Rafael Molina, and Valery Naranjo. Proportion constrained weakly supervised histopathology image classification. *Computers in Biology and Medicine*, 147:105714, 2022.
- [81] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [82] Zhi Huang, Federico Bianchi, Mert Yuksekgonul, Thomas J Montine, and James Zou. A visual–language foundation model for pathology image analysis using medical twitter. *Nature medicine*, 29(9):2307–2316, 2023.
- [83] Milda Pocevičiūtė, Gabriel Eilertsen, and Claes Lundström. Unsupervised anomaly detection in digital pathology using gans. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 1878–1882. IEEE, 2021.
- [84] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In *International conference on machine learning*, pages 4393–4402. PMLR, 2018.
- [85] Chun-Liang Li, Kihyuk Sohn, Jinsung Yoon, and Tomas Pfister. Cutpaste: Self-supervised learning for anomaly detection and localization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9664–9674, 2021.

- [86] Vitjan Zavrtanik, Matej Kristan, and Danijel Skočaj. Draem-a discriminatively trained reconstruction embedding for surface anomaly detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8330–8339, 2021.
- [87] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in neural information processing systems*, 31, 2018.
- [88] Mojtaba Nafez, Amirhossein Koochakian, Arad Maleki, Jafar Habibi, and Mohammad Hossein Rohban. Patchguard: Adversarially robust anomaly detection and localization through vision transformers and pseudo anomalies. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 20383–20394, 2025.
- [89] J Dippel, N Preniřl, J Hense, et al. Ai-based anomaly detection for clinical-grade histopathological diagnostics. *nejm ai* 1: Aioa2400468, 2024.
- [90] Cosmin I Bercea, Benedikt Wiestler, Daniel Rueckert, and Julia A Schnabel. Generalizing unsupervised anomaly detection: towards unbiased pathology screening. In *Medical Imaging with Deep Learning*, 2023.
- [91] Yu Cai, Weiwen Zhang, Hao Chen, and Kwang-Ting Cheng. Medianomaly: A comparative study of anomaly detection in medical images. *Medical Image Analysis*, page 103500, 2025.
- [92] Ioannis Lagogiannis, Felix Meissen, Georgios Kaissis, and Daniel Rueckert. Unsupervised pathology detection: a deep dive into the state of the art. *IEEE transactions on medical imaging*, 43(1):241–252, 2023.
- [93] John Duchi and Hongseok Namkoong. Variance-based regularization with convex objectives. *Journal of Machine Learning Research*, 20(68):1–55, 2019.
- [94] Robert G Gallager. *Stochastic processes: theory for applications*. Cambridge University Press, 2013.
- [95] Michel Ledoux. Measure concentration, transportation cost, and functional inequalities. *Summer School on Singular Phenomena and Scaling in Mathematical Models, Bonn*, pages 10–13, 2003.
- [96] Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- [97] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library, 2019.
- [98] Paul Tourniaire, Marius Ilie, Paul Hofman, Nicholas Ayache, and Hervé Delingette. Ms-clam: Mixed supervision for the classification and localization of tumors in whole slide images. *Medical Image Analysis*, 85:102763, 2023.
- [99] Tong Ding, Sophia J Wagner, Andrew H Song, Richard J Chen, Ming Y Lu, Andrew Zhang, Anurag J Vaidya, Guillaume Jaume, Muhammad Shaban, Ahrong Kim, et al. Multimodal whole slide foundation model for pathology. *arXiv preprint arXiv:2411.19666*, 2024.
- [100] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017.
- [101] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2017. URL <https://api.semanticscholar.org/CorpusID:53592270>.
- [102] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [103] Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo. Reliable fidelity and diversity metrics for generative models. In *International conference on machine learning*, pages 7176–7185. PMLR, 2020.
- [104] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? In *International conference on machine learning*, pages 5389–5400. PMLR, 2019.

- [105] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. Measuring robustness to natural distribution shifts in image classification. *Advances in Neural Information Processing Systems*, 33:18583–18599, 2020.
- [106] John P Miller, Rohan Taori, Aditi Raghunathan, Shiori Sagawa, Pang Wei Koh, Vaishaal Shankar, Percy Liang, Yair Carmon, and Ludwig Schmidt. Accuracy on the line: on the strong correlation between out-of-distribution and in-distribution generalization. In *International conference on machine learning*, pages 7721–7735. PMLR, 2021.
- [107] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8340–8349, 2021.
- [108] Thomas G Dietterich, Richard H Lathrop, and Tomás Lozano-Pérez. Solving the multiple instance problem with axis-parallel rectangles. *Artificial intelligence*, 89(1-2):31–71, 1997.
- [109] Stuart Andrews, Ioannis Tsochantaridis, and Thomas Hofmann. Support vector machines for multiple-instance learning. *Advances in neural information processing systems*, 15, 2002.

A Theoretical Analysis

In this section, we first provide a mathematical formulation of the multiple instance learning (MIL) setting considered in this work. We begin by establishing the notation and then describe the main problem setup, along with the data generation model. Specifically, we introduce a statistical model that captures the generation process of images in the dataset, as well as the set (or bag) of patches associated with each image. Our modeling approach is intentionally kept general, incorporating only standard and widely accepted assumptions such as the local dependence of image patches and the i.i.d. nature of images in the training set. Additionally, although our experiments use the cross-entropy loss, the theoretical analysis accommodates any proper, continuous, Lipschitz loss function.

We then proceed to the theoretical analysis, starting with uniform convergence bounds that provide generalization guarantees for our method. In particular, we study the role of the hyperparameter β in the log-sum-exp formulation and its influence on generalization. Our main result in this part is stated in Theorem 4.

Subsequently, we shift focus to the sensitivity analysis in Section A.4. Notably, in contrast to generalization—which favors smaller values of β —high sensitivity in noisy regimes requires sufficiently large β . In Theorem 6, we show that choosing $\beta \gg \mathcal{O}(\log N)$ is both necessary and sufficient to achieve a high true positive rate and to mitigate false negatives.

The key conclusion of this section is that all values of $\beta \in (0, +\infty)$ are *Pareto-optimal*: depending on which aspect—generalization or sensitivity—is prioritized and how the trade-off is weighted, the optimal choice of β will vary. This insight supports the central idea behind our method: using smaller values of β during training to promote generalization, while selecting larger values (possibly even ∞) at inference time to enhance sensitivity.

A.1 Data Generation Model

Let $n, N \in \mathbb{N}$. Assume we observe n i.i.d. images $I_1, \dots, I_n \stackrel{i.i.d.}{\sim} P_0$, where $P_0 \in \mathcal{M}(\mathbb{R}^{D_1 \times D_2})$ is an unknown distribution over images of size $D_1 \times D_2$, with $D_1, D_2 \in \mathbb{N}$. Assume each image I_j for $j \in [n]$ is decomposed into N non-overlapping (or possibly overlapping) patches:

$$\text{Patch}(I_j) \triangleq (\mathbf{X}_1^{(j)}, \dots, \mathbf{X}_N^{(j)}).$$

Each image patch $\mathbf{X}_i^{(j)}$ (for $j \in [n], i \in [N]$) is passed through a fixed vision transformer model network ViT(\cdot) to produce a feature vector $\mathbf{x}_i^{(j)}$:

$$\mathbf{x}_i^{(j)} \triangleq \text{ViT}(\mathbf{X}_i^{(j)}) \in \mathcal{X}, \quad (2)$$

where $\mathcal{X} \subseteq \mathbb{R}^d$ is the feature space, and d denotes the feature dimension. Define the *bag of features* for image I_j as

$$\text{Bag}(I_j) \triangleq (\mathbf{x}_1^{(j)}, \dots, \mathbf{x}_N^{(j)}), \quad j \in [n].$$

Assumption 1 (Local Dependence of Image Patches). *Let $I \sim P_0$, and consider the N feature vectors $(\mathbf{x}_1, \dots, \mathbf{x}_N)$ in $\text{Bag}(I)$. Then, we assume they form a locally dependent stochastic process such that:*

$$\mathbf{x}_i \perp \mathbf{x}_j \quad \forall i, j \in [N] \text{ with } |i - j| > B, \quad (3)$$

where \perp indicates statistical independence according to P_0 , and $1 \leq B \ll N$ is a fixed parameter independent of N .

This assumption has been widely validated empirically across a broad range of natural and medical image datasets. Indeed, the statistical correlations among different patches of such images are predominantly local: patches that are spatially distant—i.e., whose patch indices differ by more than a certain *bandwidth* $B = \mathcal{O}(1)$ —are typically almost independent, even though they might be identically or non-identically distributed.

A.2 Model Training via SGD

Consider the class of Multi-Layer Perceptrons (MLPs) defined as

$$\mathcal{F} \triangleq \{f_{\theta} : \mathcal{X} \rightarrow \mathbb{R} \mid \theta \in \Theta\},$$

where Θ denotes the parameter space, i.e., the set of all possible weight configurations of the neural networks. For each $\theta \in \Theta$, the function f_{θ} represents an MLP that maps input feature vectors from \mathcal{X} to real-valued outputs.

Definition 3 (LSE Loss). Fix $\beta \in (0, +\infty)$ and a Lipschitz loss function $\mathcal{L}(y||y') : [0, 1]^2 \rightarrow \mathbb{R}_+$ for (soft) labels $y, y' \in [0, 1]$. For $\mathbf{x} \in \mathcal{X}$, define the *log-sum-exp aggregation function*:

$$A(\mathbf{x}_1, \dots, \mathbf{x}_N) \triangleq \frac{1}{\beta} \log \left(\frac{1}{N} \sum_{i=1}^N \exp(\beta \cdot \sigma \circ f_{\theta}(\mathbf{x}_i)) \right),$$

where $\sigma \circ f_{\theta}$ is the composition of a learned function $f_{\theta} : \mathcal{X} \rightarrow \mathbb{R}$ and an activation $\sigma : \mathbb{R} \rightarrow [0, 1]$. The LSE-based loss function for parameters $\theta \in \Theta$ for the empirical dataset $\{(I_j, y_j) \mid j = 1, \dots, n\}$ is then given by:

$$\mathbb{L}_{\text{LSE}}(\theta) \triangleq \frac{1}{n} \sum_{j=1}^n \mathcal{L} \left(y_j \parallel A(\mathbf{x}_1^{(j)}, \dots, \mathbf{x}_N^{(j)}) \right). \quad (4)$$

Note that we have

$$\begin{aligned} \lim_{\beta \rightarrow 0} A(\mathbf{x}_1, \dots, \mathbf{x}_N) &= \frac{1}{N} \sum_{i=1}^N \sigma \circ f_{\theta}(\mathbf{x}_i), \\ \text{while } \lim_{\beta \rightarrow \infty} A(\mathbf{x}_1, \dots, \mathbf{x}_N) &= \max_{i \in [N]} \sigma \circ f_{\theta}(\mathbf{x}_i). \end{aligned} \quad (5)$$

Consider the loss function defined in Definition 3, and suppose that an empirical risk minimization (ERM) procedure yields the estimator $\hat{\theta}_{\text{ERM}}$ defined as

$$\hat{\theta}_{\text{ERM}} \triangleq \underset{\theta \in \Theta}{\text{argmin}} \mathbb{L}_{\text{LSE}}(\theta). \quad (6)$$

The minimization problem is assumed to be solved using a standard Stochastic Gradient Descent (SGD) algorithm with an arbitrary batch size. It is important to note that such optimization problems are generally non-convex and, therefore, are not expected to be solved to global optimality. In practice, due to early stopping and the inherent complexity of the loss landscape, the algorithm typically converges to a local minimum or a stationary point. However, this does not impact the validity of our subsequent analysis. Specifically, our generalization guarantees are stated as high-probability uniform convergence bounds that hold for all $\theta \in \Theta$, regardless of the particular solution returned by the optimization algorithm.

A.3 Generalization Bounds

Define the generalization error as

$$\text{GE}(\hat{\theta}_{\text{ERM}}) \triangleq \mathbb{E}_{P_0} \left[\mathbb{L}_{\text{LSE}}(\hat{\theta}_{\text{ERM}}) \right] - \mathbb{L}_{\text{LSE}}(\hat{\theta}_{\text{ERM}}).$$

Theorem 4 (Effect of β on Generalization Error). *Let Assumption 1 hold for some B , and define $q_{i,j} \triangleq \sigma \circ f_{\hat{\theta}_{\text{ERM}}}(\mathbf{x}_i^j)$ for all images $j \in [n]$ and patch ids $i \in [N]$. Also let $U \triangleq \frac{1}{n} \sum_{j=1}^n \max_i q_{i,j}$ and $M \triangleq \frac{1}{n} \sum_{j=1}^n \text{median}_i q_{i,j}$. Then, having $\beta \ll \frac{\log(N/2)}{U-M}$, with high probability results into the following bound:*

$$\text{GE}(\hat{\theta}_{\text{ERM}}) \leq \mathcal{O} \left((nN)^{-1/2} + n^{-1} \right).$$

On the other hand, if $\beta \gg \frac{\log(N/2)}{U-M}$, then, depending on the tail behavior of the $q_{i,j}$'s, the generalization error may (again with high probability) satisfy

$$\text{GE}(\hat{\theta}_{\text{ERM}}) \asymp \mathcal{O}(n^{-1/2}).$$

Proof of Theorem 4. We first show that the LSE behaves similarly to an exponentially weighted average, exhibiting a smooth phase transition as β increases from below to above the threshold

$$\frac{\log \frac{N}{2}}{U - M}.$$

In the small- β regime, the LSE behaves roughly like a simple average, and thus the variance of the loss decreases with both n and N . In contrast, in the large- β regime, the LSE approximates the maximum operator, which may exhibit weaker concentration properties depending on the tail behavior of the variables $\sigma \circ f_{\theta}(\mathbf{x}_i^{(j)})$.

Lemma 1. For $\beta \geq 0$, the log-sum-exp (LSE) and the exponentially weighted average of N values q_1, \dots, q_N bound each other as follows:

$$\left| \frac{1}{\beta} \log \left(\frac{1}{N} \sum_{i=1}^N e^{\beta q_i} \right) - \sum_{i=1}^N q_i \cdot \frac{e^{\beta q_i}}{\sum_{j=1}^N e^{\beta q_j}} \right| \leq C \cdot \min \left\{ \beta, \frac{\log N}{\beta} \right\}, \quad (7)$$

where C is a constant independent of β , but dependent on the empirical distribution of the q_i s.

Proof. For $\beta \ll 1$, we use the Taylor expansion of the exponential function: $e^x = 1 + x + \frac{x^2}{2} + \mathcal{O}(x^3)$. Applying this to $e^{\beta q_i}$, we get

$$\begin{aligned} \frac{1}{\beta} \log \left(\frac{1}{N} \sum_{i=1}^N e^{\beta q_i} \right) &= \frac{1}{\beta} \log \left(\frac{1}{N} \sum_{i=1}^N \left(1 + \beta q_i + \frac{\beta^2 q_i^2}{2} + \mathcal{O}(\beta^3) \right) \right) \\ &= \frac{1}{\beta} \log \left(1 + \beta \bar{q} + \frac{\beta^2 \bar{q}^2}{2} + \mathcal{O}(\beta^3) \right) \\ &= \bar{q} + \frac{\beta}{2} (\bar{q}^2 - \bar{q}^2) + \mathcal{O}(\beta^2), \end{aligned} \quad (8)$$

where $\bar{q} = \frac{1}{N} \sum_{i=1}^N q_i$ and $\bar{q}^2 = \frac{1}{N} \sum_{i=1}^N q_i^2$. For the exponentially weighted average, we expand the weights similarly:

$$\begin{aligned} \sum_{i=1}^N q_i \cdot \frac{e^{\beta q_i}}{\sum_{j=1}^N e^{\beta q_j}} &= \sum_{i=1}^N q_i \cdot \frac{1 + \beta q_i + \mathcal{O}(\beta^2)}{N (1 + \beta \bar{q} + \mathcal{O}(\beta^2))} \\ &= \bar{q} + \beta (\bar{q}^2 - \bar{q}^2) + \mathcal{O}(\beta^2). \end{aligned} \quad (9)$$

Thus, the absolute difference between the two expressions is

$$|A - B| \leq \frac{\beta}{2} \cdot \text{Var}(q_1, \dots, q_N) + \mathcal{O}(\beta^2).$$

For $\beta \gg 1$, denote $q_{\max} = \max_i q_i$ and let i^* be the index where this maximum is achieved. We write

$$\begin{aligned} \frac{1}{\beta} \log \left(\frac{1}{N} \sum_{i=1}^N e^{\beta q_i} \right) &= \frac{1}{\beta} \log \left(\frac{1}{N} e^{\beta q_{\max}} \left[1 + \sum_{i \neq i^*} e^{-\beta (q_{\max} - q_i)} \right] \right) \\ &= q_{\max} - \frac{\log N}{\beta} + \frac{1}{\beta} \log \left(1 + \sum_{i \neq i^*} e^{-\beta (q_{\max} - q_i)} \right) \\ &= q_{\max} - \frac{\log N}{\beta} + \mathcal{O}(\beta^{-2}). \end{aligned} \quad (10)$$

For the exponentially weighted average,

$$\begin{aligned} \sum_{i=1}^N q_i \cdot \frac{e^{\beta q_i}}{\sum_{j=1}^N e^{\beta q_j}} &= \frac{q_{\max} + \sum_{i \neq i^*} q_i e^{-\beta (q_{\max} - q_i)}}{1 + \sum_{i \neq i^*} e^{-\beta (q_{\max} - q_i)}} \\ &= q_{\max} + \mathcal{O}(\beta^{-2}). \end{aligned} \quad (11)$$

Therefore, the difference in this regime satisfies $|A - B| \leq \frac{\log N}{\beta} + \mathcal{O}(\beta^{-2})$. Combining both cases, we obtain the claimed bound:

$$\left| \frac{1}{\beta} \log \left(\frac{1}{N} \sum_{i=1}^N e^{\beta q_i} \right) - \sum_{i=1}^N q_i \cdot \frac{e^{\beta q_i}}{\sum_{j=1}^N e^{\beta q_j}} \right| \leq C \cdot \min \left\{ \beta, \frac{\log N}{\beta} \right\}.$$

□

Hence, the two formulations become asymptotically equivalent as $\beta \rightarrow 0^+$ and as $\beta \rightarrow \infty$. In these limiting regimes—which are of primary interest in this analysis—we may treat the two expressions interchangeably.

Lemma 2. *Assume real values satisfying $q_1 \leq q_2 \leq \dots \leq q_N$ and let $\beta \geq 0$. Define:*

$$\Delta \triangleq q_N - q_1, \quad \delta \triangleq q_{\lceil N/2 \rceil} - q_1.$$

Then, for all $j = 1, \dots, N$, the exponentially weighted probabilities satisfy:

$$\frac{1}{N} \cdot \frac{2}{e^{\beta\delta} + e^{\beta\Delta}} \leq p_j \triangleq e^{\beta q_j} \left(\sum_{i=1}^N e^{\beta q_i} \right)^{-1} \leq \frac{1}{N} \cdot \frac{2e^{\beta\Delta}}{1 + e^{\beta\delta}}. \quad (12)$$

Proof. To prove the upper-bound, we can write

$$\begin{aligned} e^{\beta q_j} \left(\sum_{i=1}^N e^{\beta q_i} \right)^{-1} &\leq e^{\beta q_{\max}} \left(\sum_{i=1}^N e^{\beta q_i} \right)^{-1} \\ &= e^{\beta\Delta} \left(\sum_{i=1}^N e^{\beta(q_i - q_1)} \right)^{-1} \\ &\leq e^{\beta\Delta} \left(\frac{N}{2} + \frac{N}{2} e^{\beta\delta} \right)^{-1} \\ &= \frac{1}{N} \cdot \frac{2e^{\beta\Delta}}{1 + e^{\beta\delta}}. \end{aligned} \quad (13)$$

The lower-bound can be achieved in a similar fashion:

$$\begin{aligned} e^{\beta q_j} \left(\sum_{i=1}^N e^{\beta q_i} \right)^{-1} &\geq e^{\beta q_1} \left(\sum_{i=1}^N e^{\beta q_i} \right)^{-1} \\ &= \left(\sum_{i=1}^N e^{\beta(q_i - q_1)} \right)^{-1} \\ &\geq \left(\frac{N}{2} e^{\beta\Delta} + \frac{N}{2} e^{\beta\delta} \right)^{-1} \\ &= \frac{1}{N} \cdot \frac{2}{e^{\beta\delta} + e^{\beta\Delta}}. \end{aligned} \quad (14)$$

This completes the proof. □

Corollary 1. *Under the setting of Lemma 2, for*

$$\beta \ll \frac{\log(N/2)}{\Delta - \delta},$$

we have $p_j \leq c(\beta)/N$ for all j , where $c(\beta)$ is a constant independent of N . Conversely, when

$$\beta \gg \frac{\log(N/2)}{\Delta - \delta},$$

there exists an index $j^ \in [N]$ such that $p_{j^*} \geq 1 - c'(\beta)$ and $p_j \leq c'(\beta)$ for all $j \neq j^*$, with $c'(\beta)$ also independent of N . Moreover,*

$$\lim_{\beta \rightarrow 0^+} c(\beta) = 1, \quad \lim_{\beta \rightarrow \infty} c'(\beta) = 0.$$

Proof. The proof is straightforward. We only need to show that the quantity

$$\beta^* \triangleq \frac{\log(N/2)}{\Delta - \delta}$$

serves as the critical threshold such that, for all $\beta < \beta^*$, the upper bound in Lemma 2 drops below 1. Specifically, we verify that

$$\frac{1}{N} \cdot \frac{2e^{\beta^* \Delta}}{1 + e^{\beta^* \delta}} \leq 1 \implies \beta^* \simeq \frac{\log \frac{N}{2}}{\Delta - \delta}.$$

For $\beta \ll \beta^*$, the resulting distribution over the p_j values is nearly uniform. In contrast, when $\beta \gg \beta^*$, the probability mass concentrates sharply on the index j corresponding to $q_j = q_{\max}$, with $p_j \approx 1$ and all other p_j values approaching zero. \square

Lemma 3 (Uniform Convergence Bound in Eq. (2) of [93]). *Let $\mathcal{H} = \{h_{\boldsymbol{\theta}} : \mathcal{X} \rightarrow \mathbb{R} \mid \boldsymbol{\theta} \in \Theta\}$ be a learnable class of functions, and let P_0 be a probability measure over \mathcal{X} . Suppose $\ell : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$ is a loss function satisfying mild regularity conditions. Then, with high probability and uniformly over all $\boldsymbol{\theta} \in \Theta$, the following generalization bound holds:*

$$\mathbb{E}_{P_0} [\ell(\mathbf{X}; \boldsymbol{\theta})] \leq \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{X}_i; \boldsymbol{\theta}) + C_1 \sqrt{\frac{\text{Var}(\ell(\mathbf{X}; \boldsymbol{\theta}))}{n}} + \frac{C_2}{n}, \quad (15)$$

where $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{i.i.d.}{\sim} P_0$, and C_1, C_2 are constants depending on the model and confidence parameters.

The proof can be found in [93], as well as in several other related works cited therein. To bound the generalization gap, we must bound the variance $\text{Var}(\ell(I, y; \boldsymbol{\theta}))$, where

$$\ell(I, y; \boldsymbol{\theta}) = \mathcal{L} \left(y \left\| \frac{1}{\beta} \log \left(\frac{1}{N} \sum_{i=1}^N e^{\beta \sigma \circ f_{\boldsymbol{\theta}}(\mathbf{x}_i)} \right) \right. \right),$$

and (I, y) is a sample drawn from P_0 , with $\text{Bag}(I) = (\mathbf{x}_1, \dots, \mathbf{x}_N)$.

Lemma 4. *Assume that the function $h \mapsto \mathcal{L}(y|h)$ is L -Lipschitz for some $L \geq 0$ and all $y \in [0, 1]$. Then, by the Efron-Stein inequality,*

$$\text{Var}(\ell(I, y; \boldsymbol{\theta})) \leq L^2 \cdot \text{Var} \left(\frac{1}{\beta} \log \left(\frac{1}{N} \sum_{i=1}^N e^{\beta \sigma \circ f_{\boldsymbol{\theta}}(\mathbf{x}_i)} \right) \right).$$

The proof follows directly from the Lipschitz continuity of \mathcal{L} . Recall that due to the statement of the theorem we assume $\mathcal{L} : [0, 1]^2 \rightarrow \mathbb{R}$ is smooth and Lipschitz. Define

$$U \triangleq \frac{1}{n} \sum_{j=1}^n \max_{i=1, \dots, N} \sigma \left(f_{\boldsymbol{\theta}}(\mathbf{x}_i^{(j)}) \right), \quad M \triangleq \frac{1}{n} \sum_{j=1}^n \text{median}_{i=1, \dots, N} \sigma \left(f_{\boldsymbol{\theta}}(\mathbf{x}_i^{(j)}) \right).$$

Analysis for small β : Assume that

$$\beta \ll \frac{\log(N/2)}{U - M}.$$

Then, using Lemma 2 and related results, we obtain:

$$\begin{aligned} \text{Var}(\ell(I, y; \boldsymbol{\theta})) &\leq L^2 c^2(\beta) \cdot \text{Var} \left(\frac{1}{N} \sum_{i=1}^N \sigma \circ f_{\boldsymbol{\theta}}(\mathbf{x}_i) \right) \\ &= \frac{L^2 c^2(\beta)}{N^2} \sum_{i,j=1}^N \text{Cov}(\sigma \circ f_{\boldsymbol{\theta}}(\mathbf{x}_i), \sigma \circ f_{\boldsymbol{\theta}}(\mathbf{x}_j)) \\ &\stackrel{(i)}{=} \frac{L^2 c^2(\beta)}{N^2} \sum_{\substack{i,j=1 \\ |i-j| \leq B}}^N \text{Cov}(\sigma \circ f_{\boldsymbol{\theta}}(\mathbf{x}_i), \sigma \circ f_{\boldsymbol{\theta}}(\mathbf{x}_j)) \\ &\stackrel{(ii)}{\leq} \frac{BL^2 c^2(\beta)}{4N}, \end{aligned} \quad (16)$$

where (i) follows from a mixing or local-dependence assumption, and (ii) uses the fact that any variable bounded in $[0, 1]$ has variance at most $1/4$.

Therefore, the generalization gap is bounded by:

$$C_1 \sqrt{\frac{\text{Var}(\ell(\mathbf{X}; \boldsymbol{\theta}))}{n}} + \frac{C_2}{n} \leq \mathcal{O}\left(\frac{1}{\sqrt{Nn}} + \frac{1}{n}\right),$$

where the constants are $\mathcal{O}(1)$ in the limit as $n, N \rightarrow \infty$, assuming β is sufficiently small. These constants may also depend on other complexity measures (e.g., Rademacher complexity or VC-dimension in binary classification settings).

Analysis for large β : Assume that

$$\beta \gg \frac{\log(N/2)}{U - M}.$$

Then, based on Corollary 1, we have

$$\text{Var}(\ell(I, y; \boldsymbol{\theta})) \asymp L^2(1 - c'(\beta)) \cdot \text{Var}\left(\max_{i=1, \dots, N} \sigma \circ f_{\boldsymbol{\theta}}(\mathbf{x}_i)\right).$$

The term above—i.e., the variance of the maximum of N weakly-dependent random variables—does not, in general, admit a closed-form analytical expression. Moreover, depending on the tail behavior of the variables $\sigma \circ f_{\boldsymbol{\theta}}(\mathbf{x}_i)$, the variance of their maximum may not even decrease with N .

Lemma 5 (Var($\max_i Z_i$) for exponentials). *Assume $Z_1, \dots, Z_N \stackrel{i.i.d.}{\sim} \lambda e^{-\lambda x}$ for $x \geq 0$ and some $\lambda > 0$, and define $M_N \triangleq \max_{i=1, \dots, N} Z_i$. Then:*

$$\text{Var}(M_N) = \frac{1}{\lambda^2} \sum_{k=1}^N \frac{1}{k^2},$$

which converges to $\pi^2/(6\lambda^2)$ as $N \rightarrow \infty$, and hence does not vanish as N grows.

The proof can be found in Chapter 8 of [94]. In general, unless one is dealing with degenerate cases with sharply truncated support (e.g., a uniform distribution on $[0, 1]$), the variance of the *mean* decays faster than that of the *maximum* of $\lceil N/B \rceil$ random variables [94, 95]. This completes the proof. \square

A.4 Selection Sensitivity Analysis

In this section, we aim to theoretically analyze the effect of the parameter β on the *sensitivity* of the inference-time performance of the model $\text{LSE}_{\beta}(\sigma \circ f_{\boldsymbol{\theta}}(\mathbf{x}_{1:N}))$, where $\text{LSE} * \beta$ is defined as:

$$\text{LSE}_{\beta}(q_1, \dots, q_N) \triangleq \frac{1}{\beta} \log\left(\frac{1}{N} \sum_{i=1}^N e^{\beta q_i}\right). \quad (17)$$

A commonly accepted assumption in the analysis of whole-slide histopathology images is that often only a small fraction of image patches may contain cancerous patterns—yet this is sufficient to label the entire slide as cancerous. The remaining patches may contain entirely normal tissue.

Another widely accepted assumption is that the values of $\sigma \circ f_{\boldsymbol{\theta}}(\mathbf{x}_i)$ across different patches of an image I —i.e., for $\text{Bag}(I) = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ and some learned parameters $\boldsymbol{\theta} \in \Theta$ —exhibit statistical variation. For instance, for "normal tissue" patches $\mathbf{x}_i \mid i \in \mathcal{N}$, the outputs are distributed around a mean value μ_0 , whereas for "abnormal" or "critical tissue" patches $\mathbf{x}_i \mid i \in \mathcal{C}$, they are distributed around a higher mean μ_c , with $\mu_c > \mu_0$. Here, $\mathcal{N}, \mathcal{C} \subseteq [N]$ form a bipartition of the patch indices corresponding to normal and abnormal regions, respectively (so that $|\mathcal{N}| = N - |\mathcal{C}|$). In real-world scenarios, it is typically the case that $|\mathcal{C}| \ll N$.

The variances of the aforementioned class-conditional distributions, as well as the mean gap $\mu_c - \mu_0$, depend on how well the parameters $\boldsymbol{\theta}$ have been learned during training.

In what follows, we argue that—independent of the training quality (as captured by the class-conditional variances and the mean difference $\mu_c - \mu_0$)—choosing a larger value of β leads to

improved accuracy bias and higher sensitivity, especially in the regime where $|\mathcal{C}| \ll N$. Conversely, when $|\mathcal{C}| \geq |\mathcal{N}|$, selecting a smaller value of β is preferable. However, since the former condition ($|\mathcal{C}| \ll N$) is the one that arises in practice, our analysis supports choosing larger values of β for better sensitivity.

Assumption 2. Let $I \sim P_0$ denote a random image, and let $\text{Bag}(I) = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ denote its random patches. Let Q_1 be a distribution supported over $[N] \cup \{0\}$. Then, suppose that $N_c \sim Q_1$ denotes the number of abnormal patches in the image, which can be zero upon being a “normal” image. Then, there exist N_c unknown patches containing cancerous patterns, indexed by $\mathcal{C} \subseteq [N]$ if $N_c > 0$, and the remaining patches are normal, indexed by $\mathcal{N} = [N] \setminus \mathcal{C}$. For these patches, we assume:

$$\sigma \circ f_{\theta}(\mathbf{x}_i) \sim Q_n \quad \text{for } i \in \mathcal{N}, \quad (18)$$

$$\sigma \circ f_{\theta}(\mathbf{x}_i) \sim Q_c \quad \text{for } i \in \mathcal{C}, \quad (19)$$

where Q_n and Q_c are unknown distributions satisfying the following:

- $\mathbb{E}_{Q_c}(X) - \mathbb{E}_{Q_n}(X) \geq \Delta$ for some known $\Delta > 0$,
- $\text{Var}(Q_n), \text{Var}(Q_c) \leq V$ for some known $V \geq 0$.

No further assumptions are made about the distributions Q_1 , Q_n , or Q_c , since they all depend on the learned parameter $\theta \in \Theta$.

For simplicity, we threshold the value of $\text{LSE}_{\beta}(\sigma \circ f_{\theta}(\mathbf{x}_{1:N}))$ for an image I at

$$T \triangleq \frac{\mu_c + \mu_n}{2},$$

where values above this threshold are considered indicative of cancer, and values below it indicate a normal image. While this threshold can be optimized in more complex settings, such refinements do not affect the core analysis presented in this section. Also, consider the following standard definition for performance measure:

Definition 5 (TP/FP/TN/FN Rates of θ). For a random image $I \sim P_0$ and corresponding bag of patch features $\mathbf{x}_1, \dots, \mathbf{x}_N$, and under Assumption 2, we define the true/false positive (TP/FP) and true/false negative error rates of a given parameter $\theta \in \Theta$ as follows:

- True Positive rate is the probability of $\text{LSE}_{\beta}(\sigma \circ f_{\theta}(\mathbf{x}_{1:N})) \geq T$ given $N_c \geq 1$,
- False Positive rate is the probability of $\text{LSE}_{\beta}(\sigma \circ f_{\theta}(\mathbf{x}_{1:N})) \geq T$ given $N_c = 0$,
- True Negative rate is the probability of $\text{LSE}_{\beta}(\sigma \circ f_{\theta}(\mathbf{x}_{1:N})) < T$ given $N_c = 0$,
- False Negative rate is the probability of $\text{LSE}_{\beta}(\sigma \circ f_{\theta}(\mathbf{x}_{1:N})) < T$ given $N_c \geq 1$.

Theorem 6 (Effect of β on Sensitivity). Let Assumption 2 hold for some unknown sub-Gaussian distributions Q_0, Q_1 , and Q_2 . For $\alpha \in (0, 1)$, assume

$$\beta \geq \frac{2C_{\alpha} \mathbb{E} \left[\log \frac{N}{N_c} \mid N_c \geq 1 \right]}{\Delta - 2C'_{\alpha} \sqrt{N} N_c^{-\gamma}} = \mathcal{O} \left(\frac{1}{\Delta} \mathbb{E} \left[\log \frac{N}{N_c} \mid N_c \geq 1 \right] \right),$$

where C_{α} and C'_{α} have polylogarithmic dependence on α^{-1} , and $\gamma \in (0, 1)$ is a constant depending on the tail properties of Q_1 and Q_2 . Then, we have $\text{TP} \geq 1 - \alpha$. Conversely, conditioned on $1 \leq N_c < N/2$, assume that

$$\beta \leq \mathcal{O} \left(\Delta \left(1 - \frac{2N_c}{N} \right) \right),$$

where the constants only depend on Q_1 and Q_2 . Then, we have $\text{FN} \geq 1/2$ conditionally.

Proof of Theorem 6. According to Assumption 2 and based on the definition of log-sum-exp (LSE), we have

$$\begin{aligned} \text{LSE}_{\beta}(\sigma \circ f_{\theta}(\mathbf{x}_{1:N})) &= \frac{1}{\beta} \log \left(\frac{1}{N} \sum_{i=1}^N \exp(\beta \sigma \circ f_{\theta}(\mathbf{x}_i)) \right) \\ &= \frac{1}{\beta} \log \left(\frac{N_c}{N} e^{\beta \Gamma_c} + \frac{N - N_c}{N} e^{\beta \Gamma_n} \right), \end{aligned} \quad (20)$$

where

$$\Gamma_c \triangleq \frac{1}{\beta} \log \left(\frac{1}{N_c} \sum_{i \in \mathcal{C}} e^{\beta q_i} \right), \quad \Gamma_n \triangleq \frac{1}{\beta} \log \left(\frac{1}{N - N_c} \sum_{i \in \mathcal{N}} e^{\beta q'_i} \right), \quad (21)$$

with $q_1, \dots, q_{N_c} \stackrel{i.i.d.}{\sim} Q_c$ and $q'_1, \dots, q'_{N-N_c} \stackrel{i.i.d.}{\sim} Q_n$. Moreover,

$$\begin{aligned} \Gamma_c &= \mu_c + \frac{1}{\beta} \log \left(\frac{1}{N_c} \sum_{i=1}^{N_c} e^{\beta(q_i - \mu_c)} \right) \triangleq \mu_c + \Delta q_c, \\ \Gamma_n &= \mu_n + \frac{1}{\beta} \log \left(\frac{1}{N - N_c} \sum_{i=1}^{N-N_c} e^{\beta(q'_i - \mu_n)} \right) \triangleq \mu_n + \Delta q_n, \end{aligned} \quad (22)$$

where Δq_c and Δq_n are the β -LSE of N_c and $N - N_c$ zero-mean sub-Gaussian random variables, respectively. It is known that as β ranges from 0 to $+\infty$, the expected values $\mathbb{E}[\Delta q_c]$, $\mathbb{E}[\Delta q_n]$ grow from 0 to $\mathcal{O}(\sqrt{V} \log N_c)$ and $\mathcal{O}(\sqrt{V} \log(N - N_c))$, respectively [96]. Additionally, their variances vanish at the rate $\{N_c, N - N_c\}^{-\gamma}$ for some $\gamma \in (1/2, 1)$ depending on the tail behavior of Q_c and Q_n . Therefore,

$$\mathbb{P} \left(\Delta q_c \leq C_\alpha \sqrt{V} (\log N_c + N_c^{-\gamma}), \quad \Delta q_n \leq C_\alpha \sqrt{V} (\log(N - N_c) + (N - N_c)^{-\gamma}) \right) \geq 1 - \alpha/2,$$

and

$$\mathbb{P} \left(\Delta q_c \geq -C_\alpha \sqrt{V} N_c^{-\gamma}, \quad \Delta q_n \geq -C_\alpha \sqrt{V} (N - N_c)^{-\gamma} \right) \geq 1 - \alpha/2,$$

for all $\alpha \in (0, 1)$, where C_α has polylogarithmic dependence on α^{-1} .

Guarantee on TP rate: With probability at least $1 - \alpha$ given N_c , we have

$$\begin{aligned} \text{LSE}_\beta(\sigma \circ f_\theta(\mathbf{x}_{1:N})) &= \mu_n + \Delta - \frac{\log(N/N_c)}{\beta} + \frac{1}{\beta} \log \left(e^{\beta \Delta q_c} + \left(\frac{N - N_c}{N_c} \right) e^{-\beta \Delta} e^{\beta \Delta q_n} \right) \\ &\geq \mu_n + \Delta - \frac{\log(N/N_c)}{\beta} + \Delta q_c \\ &\geq \mu_n + \Delta - \frac{\log(N/N_c)}{\beta} - C_\alpha \sqrt{V} N_c^{-\gamma}. \end{aligned} \quad (23)$$

Thus, to ensure $\text{LSE}_\beta(\sigma \circ f_\theta(\mathbf{x}_{1:N})) \geq T$, it suffices that

$$\beta \geq \frac{\log(N/N_c)}{\frac{\Delta}{2} - C_\alpha \sqrt{V} N_c^{-\gamma}}.$$

Taking expectation over N_c , it suffices that

$$\beta \geq \frac{2C'_\alpha \mathbb{E}[\log(N/N_c) | N_c \geq 1]}{\Delta - 2C_\alpha \sqrt{V} N_c^{-\gamma}},$$

which implies $\text{TP} \geq 1 - \alpha$. This proves the first part.

Guarantee on FN rate: Next, we show that small values of β can sharply degrade the TP rate, leading to an increase in FN rate. When β is small, the distributions of Δq_c and Δq_n are nearly symmetric. This directly implies that β -LSE value is symmetrically distributed around

$$\frac{1}{\beta} \log \left(\frac{N_c}{N} e^{\beta \mu_c} + \frac{N - N_c}{N} e^{\beta \mu_n} \right).$$

Therefore, conditioned on $N_c \geq 1$, we consider the case where

$$\frac{1}{\beta} \log \left(\frac{N_c}{N} e^{\beta \mu_c} + \frac{N - N_c}{N} e^{\beta \mu_n} \right) \leq \frac{\mu_c + \mu_n}{2},$$

which implies $\text{FN} \geq 1/2$. Using Taylor expansion for small β :

$$\begin{aligned} \frac{\mu_c + \mu_n}{2} &\geq \frac{1}{\beta} \log \left(\frac{N_c}{N} e^{\beta \mu_c} + \frac{N - N_c}{N} e^{\beta \mu_n} \right) \\ &= \frac{N_c}{N} \mu_c + \frac{N_c}{2N} \beta \mu_c^2 + \frac{N - N_c}{N} \mu_n + \frac{N - N_c}{2N} \beta \mu_n^2 + \mathcal{O}(\beta^2). \end{aligned} \quad (24)$$

Rearranging yields the sufficient condition:

$$\beta \leq \frac{\Delta \left(1 - \frac{2N_c}{N}\right)}{\mu_n^2 + \frac{N_c}{N} (\Delta^2 + 2\mu_n \Delta)},$$

which guarantees a conditional FN rate of at least $1/2$. This completes the proof. \square

Theorem 6 shows that selecting a sufficiently large value of β ensures an arbitrarily high true positive (TP) rate. Specifically, setting $\beta \gg \mathcal{O}\left(\frac{1}{\Delta} \log N\right)$ —with constants depending only polylogarithmically on the target confidence level—is sufficient to achieve any desired TP rate. Conversely, in the presence of noise, choosing a small β leads to a substantial false negative (FN) rate, significantly degrading performance. Therefore, the theorem recommends using large values of β to ensure robust sensitivity.

B Algorithms

B.1 Maxsoft

To avoid using complex notation, we explain our algorithm on a binary classification problem. With simple modifications, our method can be extended to other MIL problems such as multiclass classification or regression.

As defined, LSE_β is an aggregation function in form of

$$\text{LSE}_\beta(q_1, \dots, q_N) \triangleq \frac{1}{\beta} \log \left(\frac{1}{N} \sum_{i=1}^N e^{\beta q_i} \right), \quad (25)$$

represents a scalar instance value in the range $[0, 1]$. For simplicity, we denote $\text{LSE}_\beta(q_{1:N})$ as $\text{LSE}_\beta(q_1, \dots, q_N)$. We can see that

$$\frac{\partial \text{LSE}_\beta}{\partial q_i}(q_{1:N}) = \frac{e^{\beta q_i}}{\sum_{j=1}^N e^{\beta q_j}}.$$

which is simply the Softmax function when applied to q_1, \dots, q_N .

In multiple instance learning of WSI we aim to train a model to reduce $\text{L}_{\text{LSE}}(\theta)$ with respect to θ (see Definition 3). Let us introduce likelihood probability with

$$p_j = P(y_j = 1 | I_j) = P(y_j = 1 | q_1^{(j)}, \dots, q_N^{(j)}), \quad (26)$$

rephrasing $\text{L}_{\text{LSE}}(\theta)$ as $\frac{1}{n} \sum_{j=1}^n \mathcal{L}(y_j || p_j)$. We specifically use the cross entropy loss function for \mathcal{L} .

In Maxsoft, we perform the forward pass with $p_j = \max_i(q_1^{(j)}, \dots, q_N^{(j)})$. Afterward, in the backward step, we approximate $\text{L}_{\text{LSE}}(\theta)$ by replacing the max aggregation operator with LSE_β . More precisely, we approximate $\nabla_{\theta} \text{L}_{\text{LSE}}$ as:

$$\nabla_{\theta} \text{L}_{\text{LSE}} \approx \frac{1}{n} \sum_{j=1}^n \frac{\partial \mathcal{L}}{\partial p_j} \left(\max(q_{1:N}^j) \right) \cdot \text{Softmax}(q_{1:N}^j) \cdot J_{q_{1:N}^j}(\theta). \quad (27)$$

where J is the Jacobian matrix. Notice that in case we calculated the original gradients, the Softmax term would be replaced with the max gradient (which due to non-differentiability has instable behavior) and in case we initially used LSE in the forward pass, $\frac{\partial \mathcal{L}}{\partial p_j} \left(\max(q_{1:N}^j) \right)$ would be replaced with $\frac{\partial \mathcal{L}}{\partial p_j} \left(\text{LSE}_\beta(q_{1:N}^j) \right)$.

Algorithm 1 Multiple Instance Learning of WSI with Maxsoft

- 1: **Input:** WSI Dataset $D = \{(I_1, y_1), \dots, (I_n, y_n)\}$
 - 2: **Input:** β : Hyperparameter for LSE_β
 - 3: **Input:** ViT: Frozen Vision Transformer feature extractor
 - 4: **Input:** f_θ : Trainable MIL classification head
 - 5: **Input:** α : Step size (learning rate)
 - 6: **Output:** f_{θ^*} : Trained model

 - 7: **for** each training epoch **do**
 - 8: **for** $j = 1$ to n **do**
 - 9: Extract patches $\{\mathbf{X}_1^{(j)}, \dots, \mathbf{X}_N^{(j)}\}$ from image I_j
 - 10: Obtain features $\{\mathbf{x}_i^{(j)} = \text{ViT}(\mathbf{X}_i^{(j)})\}_{i=1}^N$
 - 11: Compute instance-level scores $q_i^{(j)} = \sigma(f_\theta(\mathbf{x}_i^{(j)}))$ for $i = 1, \dots, N$
 - 12: Compute bag-level prediction: $p_j = \max(q_{1:N}^{(j)})$ (which is also test time prediction)
 - 13: Compute loss $\mathcal{L}(y_j \| p_j)$
 - 14: **end for**
 - 15: Compute gradient approximation:
$$\nabla_{\theta} \text{L}_{\text{LSE}} \approx \frac{1}{n} \sum_{j=1}^n \frac{\partial \mathcal{L}}{\partial p_j} (Y_j \| \max(q_{1:N}^{(j)})) \cdot \text{Softmax}(q_{1:N}^{(j)}) \cdot J_{q_{1:N}^{(j)}}(\theta)$$
 - 16: Update model parameters: $\theta \leftarrow \theta - \alpha \cdot \nabla_{\theta} \text{L}_{\text{LSE}}$
 - 17: **end for**
 - 18: Return trained model: $\theta^* \leftarrow \theta$
 - 19: **return** f_{θ^*}
-

Algorithm 2 PerSlide Augmentation for WSI Tasks

- 1: **Input:** WSI Dataset $D = \{(I_1, y_1), \dots, (I_n, y_n)\}$, Augmentation functions $\mathcal{T} = \{\tau_1, \dots, \tau_m\}$, Epochs E , Number of instances per data point N
 - 2: **Output:** Trained model

 - 3: **for** epoch $t = 1$ to E **do**
 - 4: Initialize augmented dataset $D_t^{\text{aug}} \leftarrow \emptyset$
 - 5: **for** each $I_i \in D$ **do**
 - 6: Sample augmentation $\tau \sim \text{Uniform}(\mathcal{T})$
 - 7: Apply τ to the whole image I_i to get \hat{I}_i
 - 8: Add (\hat{I}_i, y_i) to D_t^{aug}
 - 9: **end for**
 - 10: TrainOneEpoch(D_t^{aug}, t)
 - 11: **end for**
-

Algorithm 3 PerPatch Augmentation

```
1: Input: WSI Dataset  $D = \{(I_1, y_1), \dots, (I_n, y_n)\}$ , Augmentation functions  $\mathcal{T} = \{\tau_1, \dots, \tau_m\}$ ,  
   Epochs  $E$ , Number of instances per data point  $N$ , ViT Frozen Vision Transformer feature  
   extractor  
2: Output: Trained model  
  
3: Precompute Stage:  
4: for each image  $I_i \in D$  do  
5:   Set first augmented image variant as the original  $I_i^{(0)} \leftarrow I_i$   
6:   for each augmentation  $\tau_k \in \mathcal{T}$ , where  $k = 1, \dots, m$  do  
7:     Initialize augmented image variant  $I_i^{(k)} \leftarrow \emptyset$   
8:     for each patch  $X_j \in I_i$  where  $j = 1, 2, \dots, N$  do  
9:       Apply augmentation:  $X_j^{(k)} = \tau_k(X_j^i)$   
10:      Compute embedding for further optimization:  $x_j^{(k)} = \text{ViT}(X_j^{(k)})$   
11:      Add  $x_j^{(k)}$  to  $I_i^{(k)}$   
12:    end for  
13:  end for  
14: end for  
  
15: Training Stage:  
16: for epoch  $t = 1$  to  $E$  do  
17:   Initialize augmented dataset  $D_t^{\text{aug}} \leftarrow \emptyset$   
18:   for each image  $I_i \in D$  do  
19:     Initialize empty set  $\hat{I}_i \leftarrow \emptyset$   
20:     for each patch  $X_j \in I_i$  where  $j = 1, 2, \dots, N$  do  
21:       Sample augmentation index  $k \sim \text{Uniform}(\{0, \dots, m\})$   
22:       Retrieve embedding  $x_j^{(k)}$  from  $I_i^{(k)}$   
23:       Add  $x_j^{(k)}$  to  $\hat{I}_i$   
24:     end for  
25:     Add  $(\hat{I}_i, y_i)$  to  $D_t^{\text{aug}}$   
26:   end for  
27:   TrainOneEpoch( $D_t^{\text{aug}}, t$ )  
28: end for
```

For cases with multiple classes we simply modify MLP’s linear head to have multiple heads for each class (or in the worst case for architectures which differ from MLP, multiple classification head and Maxsoft per class), extending equation 26 for cases with $y = k$. To further facilitate the usage of our Maxsoft aggregation function, we provide a simple Python implementation of it in the paper. It can be easily integrated into any code like other Pytorch [97] layers.

C PyTorch Implementation of Maxsoft Pooling

Listing 1 presents the PyTorch-style implementation of Maxsoft pooling.

Listing 1: Pytorch-style implementation for Maxsoft

```
class MaxSoftmaxSTE(torch.autograd.Function):
    @staticmethod
    def forward(ctx, input, beta):
        # Save the input tensor for backward pass
        ctx.save_for_backward(input)
        ctx.beta = beta

        # Perform the forward pass (select maximum value)
        max_val, _ = torch.max(input, dim=0)

        return max_val

    @staticmethod
    def backward(ctx, grad_output):
        # Retrieve the saved input tensor
        input, = ctx.saved_tensors
        beta = ctx.beta

        # Compute the Softmax over the input for gradient
        Softmax_grad = torch.Softmax(input * beta, dim=0)

        # Multiply the incoming gradient (grad_output) by
        # the Softmax weights
        grad = grad_output * Softmax_grad

        return grad, None
```

D PerSlide vs. PerPatch Figure

Figure 5 contains a visual comparison of the difference between PerSlide and PerPatch augmentations.

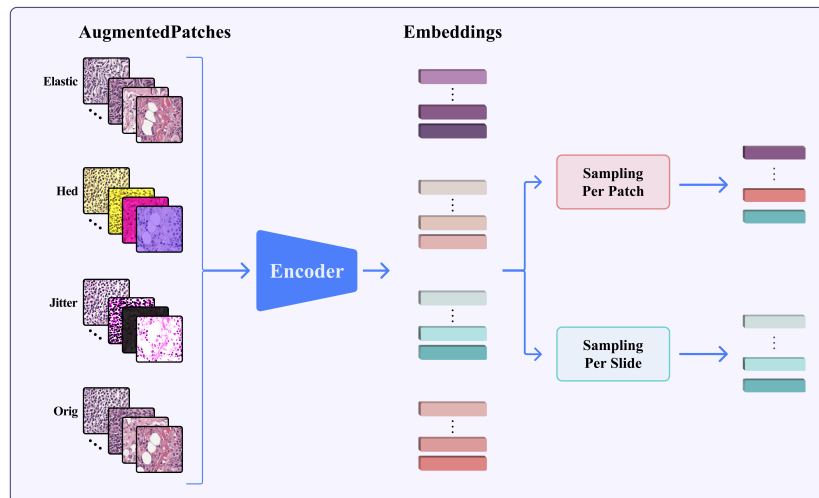


Figure 5: Comparison of PerSlide and PerPatch augmentations: PerSlide applies a single augmentation uniformly across all patches in a WSI, whereas PerPatch independently selects an augmentation for each patch from multiple variants, resulting in substantially higher diversity during MIL pooling training.

E Datasets

CAMELYON16 and CAMELYON17 are large-scale WSI datasets introduced for benchmarking automated methods in detecting metastatic breast cancer in lymph nodes. Developed for challenging settings, they serve as standards for evaluating tumor detection algorithms in histopathology. Automating this task can reduce pathologist workload, improve diagnostic consistency, and mitigate subjectivity [78, 35]. CAMELYON16 includes two classes—normal and tumor—while CAMELYON17 includes four: normal, macro, micro, and ITC. The ITC class (isolated tumor cells) is especially challenging due to its sparse, small tumor clusters (<0.2 mm or <200 cells). Both datasets are difficult because tumor regions occupy only a small area in positive WSIs. CAMELYON16 provides an official labeled test split and pixel-level annotations for all WSIs, whereas CAMELYON17 lacks test labels and includes pixel-level annotations for only 100 of its 500 training WSIs.

TCGA-Lung is a subset of The Cancer Genome Atlas (TCGA) comprising WSIs from two lung cancer subtypes: Lung Adenocarcinoma (LUAD) and Lung Squamous Cell Carcinoma (LUSC). After filtering out low-quality slides, the dataset includes 1,042 WSIs—530 LUAD and 512 LUSC. A key characteristic is that tumor regions occupy the majority of each slide. Additionally, most patients contribute multiple WSIs [79]. Pixel-level annotations are not available.

SICAP-MIL is a publicly available dataset designed to benchmark MIL-based approaches for prostate cancer grading in WSIs. It includes biopsy slides from 271 patients, scanned at $40\times$ and tiled into overlapping 512×512 patches at $10\times$ resolution. Each slide is globally labeled by expert pathologists with primary and secondary Gleason grades, reflecting the dominant tumor patterns. The dataset also introduces proportional constraints that represent the relative occurrence of each grade, supporting the development of constrained MIL methods that can rival fully supervised models. Slides are labeled as normal or abnormal and further annotated with Gleason grades (GG3, GG4, GG5). In abnormal WSIs, tumor-associated regions are present in a significant portion, though not the majority of the slide. Exact pixel-level annotations are not provided [80].

F Additional Experimental Setup

CAMELYON16 and CAMELYON17: WSIs are tiled into 256×256 non-overlapping patches at $20\times$ magnification, with background regions excluded following [67]. For CAMELYON16, we use the official test split. As CAMELYON17 lacks an official labeled test set and includes one low-quality slide, we discard that slide and randomly split the remaining data into approximately 60% training, 15% validation, and 25% testing, using the balanced splitting protocol from [78, 35, 98]. Specifically, we ensure each split has a roughly equal number of Normal, Macro, Micro, and ITC samples.

For CAMELYON16, we train on 50 WSIs and evaluate on the full test set. For CAMELYON17, we use the 99 high-quality annotated WSIs, as region-level annotation is clinically most relevant. This limited-WSI setup emulates extreme low-data regimes typical of rare cancers (≈ 70 WSIs/type [99]). The CAMELYON16 patching yields ≈ 1.5 M patches, averaging $\approx 7,900$ per WSI. For CAMELYON17, we obtain ≈ 800 k patches, about 8,000 per WSI on average. Patches overlapping annotated tumor regions are labeled tumor; all others are labeled normal. Each model is trained from scratch five times, and we report means and standard deviations for all metrics. We binarize CAMELYON17 following [14], treating ITC-labeled WSIs as tumorous. This is the most challenging setting, as some ITC slides contain only one or two tumor patches. For experiments using the complete versions of these two datasets, see Appendix P.

TCGA-Lung: WSIs are tiled into non-overlapping 256×256 patches at $20\times$ magnification, excluding background regions. This results in approximately 3.5 million patches, averaging around 3,500 patches per WSI. The dataset is roughly split into training (60%), validation (15%), and test (25%) sets and enforce patient-level grouping, whereby all slides from a patient are assigned to a single split (no cross-split leakage) [79].

SICAP-MIL: Each 512×512 patch is divided into four 256×256 patches, yielding approximately 34,000 patches in total, with each WSI containing around 100 patches on average. Each model is trained from scratch five times independently, and we report the mean and standard deviation for all performance metrics [80].

G Calibration Metric ECE

Let B_m denote the set of indices for samples with prediction confidence in the interval $I_m = (\frac{m-1}{M}, \frac{m}{M}]$. The accuracy within bin B_m is defined as:

$$\text{acc}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} 1(\hat{y}_i = y_i), \quad (28)$$

where \hat{y}_i is the predicted label and y_i is the ground-truth label for sample i . The average confidence in bin B_m is:

$$\text{conf}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \hat{p}_i, \quad (29)$$

where \hat{p}_i denotes the predicted confidence for sample i .

The Expected Calibration Error (ECE) is computed as:

$$ECE = \sum_{m=1}^M \frac{|B_m|}{n} |\text{acc}(B_m) - \text{conf}(B_m)|, \quad (30)$$

where n is the number of samples.

A perfectly calibrated model satisfies $\text{acc}(B_m) = \text{conf}(B_m)$ for all m , resulting in an ECE of 0 [100]. For instance, both $\hat{p}_i = 1$ with $\hat{y}_i = y_i$ and $\hat{p}_i = 0$ with $\hat{y}_i \neq y_i$ contribute to lower ECE [15].

H Implementation Details

H.1 MIL Models

For DINO Domain [64] on CAMELYON16, CAMELYON17, and TCGA-Lung [78, 35, 79], we train a ViT-S/16 from scratch on all patches from training WSIs using the default hyperparameters from the official DINO repository [64]. For SICAP-MIL [80], we apply the same default settings for ViT-S/16.

For all MIL methods, including Maxsoft pooling, we tune learning rate, weight decay, and weight initialization using the validation set. The best configuration is selected based on validation performance. No early stopping is applied. All models are trained for 500 epochs using the AdamW optimizer [101] with default parameters unless otherwise specified.

CAMELYON16. DINO Natural (LR 0.002, WD 0.05, Xavier-uniform); DINO Domain (LR 0.1, WD 0.05, truncated-normal init); UNI (LR 0.02, WD 0.05, Xavier-uniform); Prov-GigaPath (LR 0.02, WD 0.05, Xavier-uniform).

CAMELYON17. DINO Natural (LR 0.02, WD 0.005, Xavier-uniform); DINO Domain (LR 0.1, WD 0.05, truncated-normal); UNI (LR 0.02, WD 0.005, Xavier-uniform); Prov-GigaPath (LR 0.02, WD 0.05, Xavier-uniform).

TCGA-Lung. DINO Natural (LR 0.02, WD 0.05, Xavier-uniform); DINO Domain (LR 0.002, WD 0.005, truncated-normal); UNI (LR 0.002, WD 0.05, Xavier-uniform); Prov-GigaPath (LR 0.1, WD 0.05, orthogonal).

SICAP-MIL. DINO Natural (LR 0.02, WD 0.05, orthogonal); DINO Domain (LR 0.002, WD 0.005, Xavier-uniform); UNI (LR 0.002, WD 0.05, truncated-normal); Prov-GigaPath (LR 0.1, WD 0.05, truncated-normal).

H.2 Augmentations

Base PerPatch augmentations comprise Random Rotation, Random Gaussian Blur, and Random Color Jitter, selected per our analyses in Appendix J and Table 4; hyperparameters are dataset-specific: CAMELYON16—LR 0.1, WD 0.05, Xavier-uniform; CAMELYON17—LR 0.1, WD 0.05,

orthogonal; TCGA-Lung—LR 0.1, WD 0.05, orthogonal; SICAP-MIL—LR 0.1, WD 0.05, truncated-normal. All MIL models are trained with bag-level labels only. Experiments use PyTorch 2.1 and scikit-learn on an RTX 4090 [97].

H.3 Augmentations and Architectures

We further extend our range of experiments to assess the effect of different augmentation methods on various MIL pooling architectures on the CAMELYON17 and SICAP-MIL datasets [80]. As shown in Table 6, our PerPatch augmentation method demonstrates AUC improvements in most setups. Among the previous methods, AugDiff [62] shows performance stability across multiple architectures, whereas MixUp methods [55, 56] exhibit poor performance in some settings.

Table 6: The effect of various augmentation methods across MIL architectures on the Camelyon17 and SICAP datasets.

Augmentation	CAMELYON17													
	max pooling		mean pooling		ABMIL		DSMIL		Snuffy		LSE pooling		Maxsoft pooling	
	AUC	ECE	AUC	ECE	AUC	ECE	AUC	ECE	AUC	ECE	AUC	ECE	AUC	ECE
ReMix (MixUp)	0.743 ₃₈₅	0.245 ₁₆₉	0.823 ₀₀₆	0.118 ₀₃₀	0.833 ₀₃₈	0.193 ₀₃₈	0.723 ₃₂₄	0.162 ₀₂₆	0.688 ₁₇₆	0.179 ₀₄₄	0.777 ₁₈₆	0.212 ₀₄₆	0.840 ₀₂₀	0.242 ₀₁₇
RankMix (MixUp)	0.812 ₂₅₃	0.144 ₁₅₇	0.842 ₀₀₄	0.107 ₀₂₆	0.823 ₀₄₅	0.192 ₀₄₈	0.764 ₁₂₃	0.092 ₀₁₈	0.751 ₁₈₂	0.182 ₀₃₄	0.833 ₁₇₀	0.202 ₀₄₀	0.855 ₁₉₁	0.197 ₁₅₆
AugDiff	0.842 ₁₈₀	0.163 ₀₉₃	0.863 ₀₀₅	0.147 ₀₅₃	0.855 ₀₃₅	0.189 ₀₃₃	0.801 ₁₀₂	0.067₀₂₈	0.886 ₀₄₁	0.189 ₀₀₄	0.882 ₀₆₉	0.222 ₀₅₈	0.894 ₀₆₉	0.161 ₁₁₀
SSRDL	0.833 ₀₄₀	0.289 ₀₉₅	0.767 ₀₄₀	0.239 ₀₂₇	0.793 ₀₁₅	0.263 ₀₁₅	0.807 ₁₃₅	0.085 ₀₀₃	0.818 ₀₂₈	0.176 ₀₀₂	0.817 ₀₁₅	0.237 ₀₄₈	0.859 ₀₄₀	0.061 ₀₃₀
PerPatch	0.923₀₇₆	0.131₀₂₁	0.870 ₀₀₁	0.126 ₀₀₀	0.893₁₁₀	0.183 ₀₅₅	0.920 ₁₀₇	0.259 ₀₁₃	0.980₀₁₄	0.174₀₂₄	0.965₀₂₄	0.232 ₀₁₂	0.980₀₁₄	0.052₀₁₃

Augmentation	SICAP-MIL													
	max pooling		mean pooling		ABMIL		DSMIL		Snuffy		LSE pooling		Maxsoft pooling	
	AUC	ECE	AUC	ECE	AUC	ECE	AUC	ECE	AUC	ECE	AUC	ECE	AUC	ECE
ReMix (MixUp)	0.850 ₀₀₅	0.162 ₀₀₄	0.787 ₀₀₂	0.047 ₀₀₂	0.648 ₀₀₁	0.388 ₀₀₁	0.852 ₀₀₅	0.119 ₀₁₂	0.845 ₀₀₃	0.221 ₀₂₈	0.855 ₀₀₅	0.201 ₀₀₆	0.849₀₀₈	0.180 ₀₁₇
RankMix (MixUp)	0.846 ₀₁₃	0.157 ₀₀₉	0.808 ₀₀₂	0.042 ₀₀₂	0.683 ₀₀₁	0.381 ₀₀₁	0.849 ₀₀₄	0.133 ₀₂₂	0.855 ₀₀₂	0.193 ₀₁₁	0.836 ₀₀₂	0.220 ₀₂₂	0.845 ₀₀₈	0.280 ₀₁₇
AugDiff	0.861 ₀₂₉	0.155 ₀₀₅	0.809 ₀₀₈	0.041₀₀₃	0.724 ₀₀₅	0.260 ₀₀₉	0.856 ₀₀₅	0.129 ₀₀₈	0.860 ₀₀₇	0.191 ₀₀₄	0.863 ₀₀₂	0.188 ₀₂₁	0.818 ₀₁₉	0.180 ₀₀₅
SSRDL	0.806 ₀₀₈	0.160 ₀₂₂	0.803 ₀₀₇	0.224 ₀₂₃	0.657 ₀₀₂	0.406 ₀₀₄	0.840 ₀₀₄	0.125 ₀₁₂	0.806 ₀₁₉	0.259 ₀₀₈	0.812 ₀₀₅	0.182 ₀₀₈	0.838 ₀₀₉	0.161₀₀₅
PerPatch	0.870₀₀₁	0.155 ₀₀₃	0.810₀₀₁	0.142 ₀₀₂	0.831₀₁₀	0.219 ₀₂₁	0.862 ₀₀₆	0.118₀₀₈	0.866₀₀₃	0.186₀₀₃	0.869₀₀₀	0.176 ₀₀₂	0.827 ₀₀₂	0.223 ₀₁₈

I Augmentations Descriptions and Samples

Figure 6 presents examples of four CAMELYON17 [35] patches augmented with Random Rotation, Random Elastic Deformation, Random Affine Transformation, Random Gaussian Blurring, Random Color Jitter, and Random Hematoxylin-Eosin-DAB (HED) Jitter [59, 62]. The augmentations are defined as follows: **Random Rotation**: Rotates the image by a random angle within a predefined range. **Random Elastic Deformation**: Applies spatially varying smooth deformations to simulate elastic distortions. **Random Affine Transformation**: Combines translation, rotation, scaling, and shearing. **Random Gaussian Blurring**: Convolves the image with a Gaussian kernel to reduce high-frequency noise. **Random Color Jitter**: Randomly modifies brightness, contrast, saturation, and hue. **Random HED Jitter**: Perturbs the HED color space representation to simulate staining variations [59, 62].

J Augmentation Quality and Complexity

In this section, we evaluate the effectiveness of Base augmentations using four key metrics: FID, Density, and Coverage.

FID (Fréchet Inception Distance): FID is a widely used metric that quantitatively compares the distribution of generated images against real images in a deep feature space. In the context of image augmentation, FID helps assess how well the augmented images capture the underlying data distribution of the real dataset. A lower FID score indicates that the augmented images are closer in distribution to the original images, which implies that the generative model produces realistic and coherent augmentations. A good FID score suggests that the additional images maintain the essential visual characteristics of the real-world data, thereby potentially improving the downstream performance of learning algorithms [102].

Density: The density metric measures how densely the generated images occupy the feature space relative to the real images. In image augmentation, high density implies that the synthetic images generated by the model are not only realistic but also well-aligned with the clusters of real images in the feature space. This alignment is critical, as it suggests that the augmented images reinforce the intrinsic patterns found in the data rather than creating spurious or outlier representations. Evaluating density helps to understand whether the augmentation process is introducing variations that are

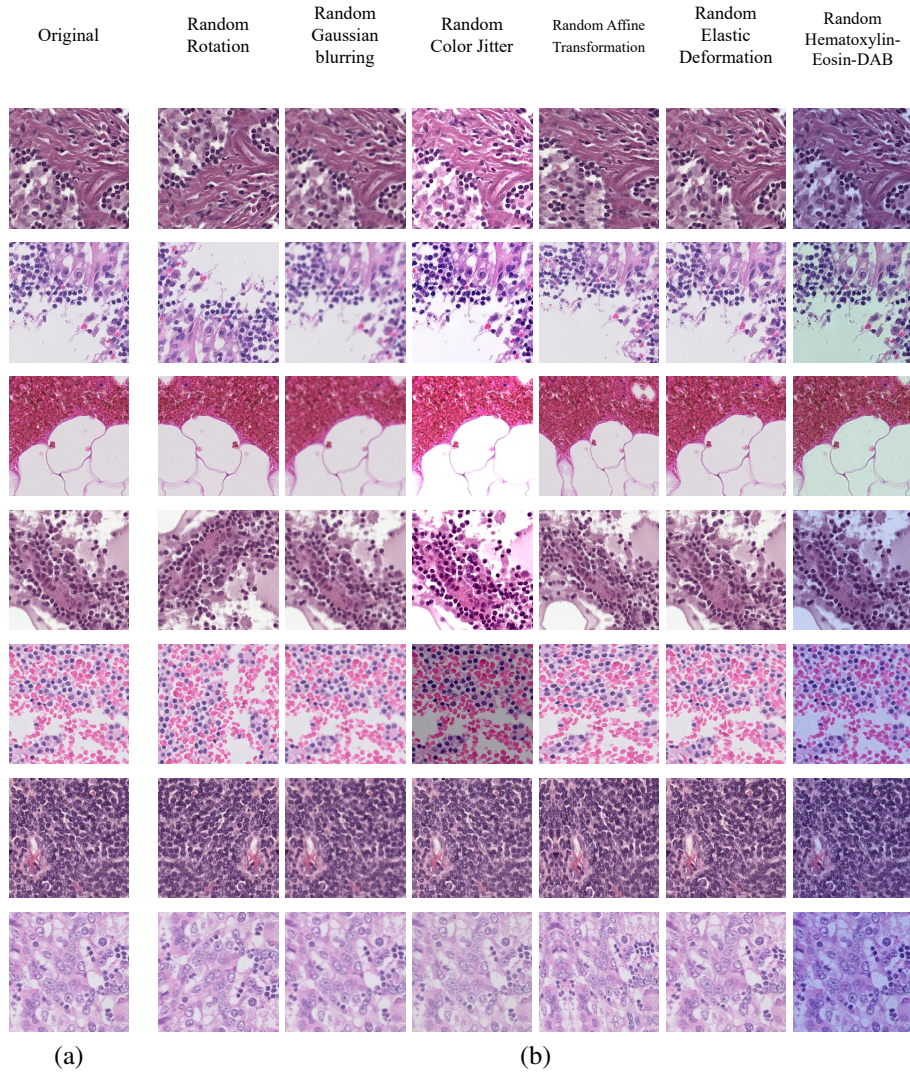


Figure 6: Augmentation schemes (a) sample CAMELYON17 patches (b) augmented versions of the same patches.

plausible and beneficial for training robust models, ensuring that the augmentation enriches the dataset with high-quality, representative samples [103].

Coverage: Coverage evaluates the diversity of the generated images by determining the proportion of the real image distribution that is covered by the synthetic samples. In the realm of image augmentation, high coverage indicates that the method is capable of producing a wide range of variations that collectively span the real data manifold. This is particularly important when the goal is to enhance a dataset by introducing new variations that help prevent overfitting and improve generalization in downstream tasks. A model with good coverage ensures that the augmented dataset is not biased toward a narrow subset of the data distribution, thereby providing a more comprehensive training set that captures the full spectrum of variability present in real-world images [103].

Table 7 shows that Random Hematoxylin–Eosin–DAB (HED) Jitter exhibits the lowest Density and Coverage and the highest FID. Despite its pathology-specific design, this augmentation neither reduces slide-level AUC variance nor improves ECE, indicating that the transformations it introduces do not reflect realistic imaging variability. In general, augmentations with higher Density and Coverage and lower FID perform better, as illustrated by the strong results of Random Rotation, Random Gaussian Blur, and Random Color Jitter, and the poor performance of Random Elastic

Deformation and Random Affine Transformation (see Table 5). We attribute the latter two failures to their limited relevance to real-world slide variations. Collectively, these findings indicate that the most effective strategy remains PerPatch combined with Random Rotation, Random Gaussian Blur, and Random Color Jitter.

Table 7: Quality Metrics for different augmentations on the CAMELYON17 dataset. The results are reported in the form of $mean.std$.

Augmentation	FID	Density	Coverage
Random Rotation	73.801 _{46.799}	0.785 _{0.066}	0.999 _{0.001}
Random Elastic Deformation	274.585 _{139.951}	0.353 _{0.162}	0.636 _{0.185}
Random Affine Transformation	135.422 _{70.692}	0.434 _{0.124}	0.858 _{0.108}
Random Gaussian Blurring	6.842 _{1.197}	0.986 _{0.030}	0.999 _{0.000}
Random Color Jitter	30.085 _{7.697}	0.672 _{0.053}	0.997 _{0.004}
Random Hematoxylin-Eosin-DAB (HED) Jitter	2395.841 _{286.359}	0.000 _{0.001}	0.001 _{0.001}

K UMAP of Patch Embeddings Based on Encoder

Figure 7 presents UMAP visualizations of patch representations obtained from DINO Natural, DINO Domain, UNI, and Prov-GigaPath on four test samples from the CAMELYON17 dataset. In general, as the strength of the representation encoder increases, the embeddings exhibit improved clustering by both label and patient, reflecting higher feature discriminability and domain alignment.

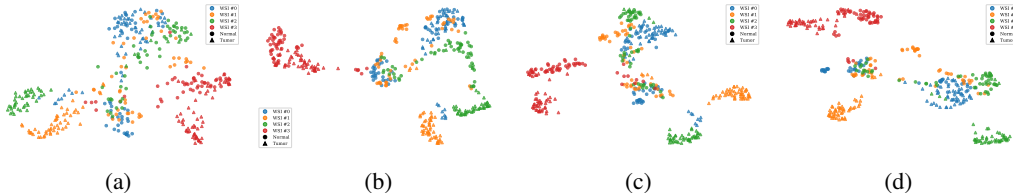


Figure 7: UMAP of representations. (a) with Dino Natural (b) Dino Domain (c) UNI (d) Prov-GigaPath.

L Error Analysis

L.1 Error Analysis on UNI for CAMELYON17

As we can see in Table 1, the results of UNI on the CAMELYON17 dataset are low. This is an exception to our finding that the better the representation encoder (in particular the more data it has been trained on) the better the performance of the MIL pooling. Through an error analysis we found that the problem with UNI comes from the fact that it cannot generalize well in WSIs with extremely low tumor regions (mostly ITC and then Micro subgroups). This only happens in CAMELYON17 since only this dataset has such small tumor regions and probably comes from the fact that UNI’s autopsy WSIs do not provide such samples. The results can be found at Table 8.

L.2 Error Analysis for Augmentations and ECE on CAMELYON

When *Per-Patch* augmentation is applied to the sparsely annotated CAMELYON17 WSIs, the network is repeatedly exposed to heavily transformed, minute tumorous regions. This increased morphological diversity strengthens its ability to recognise genuinely positive tissue and, as a result, *raises* accuracy on slides that contain tumour (Table 9). The same shift, however, slightly *erodes* performance on the overwhelmingly abundant tumour-free slides: features that were previously dismissed as benign are now more readily interpreted as malignant, leading to a higher false-positive rate.

From a calibration standpoint, the effect is likewise inverted relative to our original claim. The model becomes better calibrated on the rare, hard positives—confidence scores now align closely with their improved correctness—but it grows over-confident on negatives, which dominate the data distribution

Table 8: Accuracy results for each subgroup in CAMELYON17 on UNI representaton encoder.

Method	Slide			
	ITC ACC	Micro ACC	Macro ACC	Negative ACC
max pooling	0.531 _{.183}	0.604 _{.433}	0.535 _{.399}	0.787 _{.217}
mean pooling	0.484 _{.080}	0.667 _{.337}	0.514 _{.168}	0.758 _{.162}
ABMIL	0.417 _{.358}	0.583 _{.448}	0.542 _{.405}	0.681 _{.350}
DSMIL	0.521 _{.177}	0.604 _{.422}	0.535 _{.393}	0.821 _{.149}
Maxsoft pooling	0.713 _{.184}	0.308 _{.335}	0.753 _{.343}	0.816 _{.129}

yet see a fall in accuracy. Consequently, the overall ECE still increases, though the underlying driver is the miscalibration of negative patches rather than of positives.

In short, *Per-Patch* augmentation sharpens decision boundaries around the scarce tumour class, boosting sensitivity and reducing inter-run variance for positive findings, while sacrificing a portion of specificity on normal tissue. Post-hoc calibration targeted at the negative majority—e.g. temperature scaling on a validation set enriched for benign slides—offers a principled way to retain the newfound robustness to tumour heterogeneity without compromising probabilistic reliability.

Table 9: Accuracy results for each subgroup in CAMELYON17 with DINO Domain representation encoder with and without Augmentation.

Method	Slide			
	ITC ACC	Micro ACC	Macro ACC	Negative ACC
Maxsoft pooling No Aug	0.713 _{.184}	0.308 _{.335}	0.753 _{.343}	0.816 _{.129}
Maxsoft pooling PerPatch	1.000 _{.000}	0.500 _{.235}	0.500 _{.235}	0.750 _{.070}

M Additional ROI Detection Images

Figure 8 shows additional examples of patch-level classification on the CAMELYON17 dataset [78, 79]. Consistent with previous results, max pooling and Maxsoft pooling yield the most accurate ROI detections, while mean pooling performs noticeably worse.

N Sensitivity to Data Quality

Although our datasets are high quality and processed under stringent protocols, a simple visual inspection can verify their appearance; to test whether the underperformance of Transformer-based models in current data regimes stems from data quality, we reran all CAMELYON17 experiments under varied JPEG compression levels (50; 75—the common setting used in Tables 1; 100—the highest) and with Gaussian blur as a quality corruption. The results in Table 10 replicate prior trends: neither decreasing nor increasing quality alters Transformer behavior, rejecting data quality as the cause of Transformer-based MIL deficiencies in current data regimes.

O Experiment on Classical MIL Pooling Functions

We evaluate major classical MIL pooling functions on CAMELYON17 [35] to assess their effectiveness for WSI classification. As shown in Table 11, Maxsoft achieves the best slide- and patch-level metrics (except patch-level on DINO Natural); noisy-or performs poorly [43], ISR excels on patches, and smoothmax lies between LSE and Maxsoft. smoothmax remains without theoretical analysis; intuitively, its weaker performance may arise because its gradient weights instances by deviation from the current smoothmax value rather than Softmax-based importance, potentially hindering credit assignment relative to Maxsoft’s direct Softmax path and motivating future theoretical study.

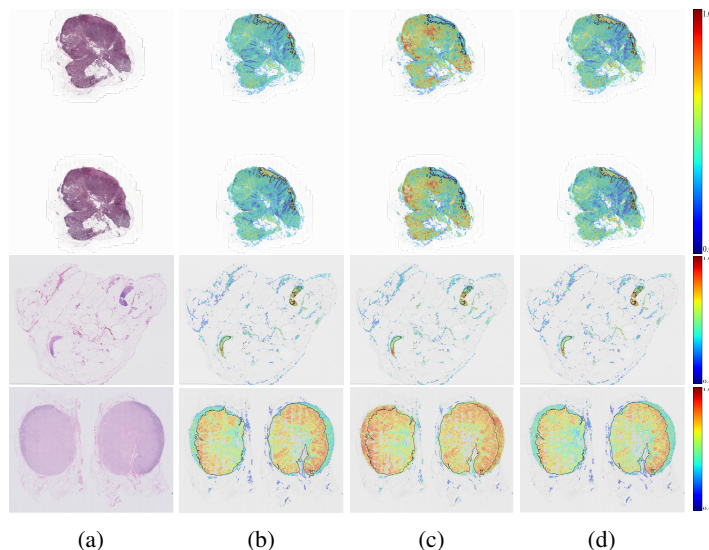


Figure 8: Overview of additional ROIs identified on representative WSIs from the CAMELYON17 dataset [35], using (a) max pooling, (b) mean pooling, and (c) Maxsoft pooling.

Table 10: MIL pooling performance under varied JPEG compression levels and Gaussian blur corruptions.

Method	JPEG 50		JPEG 75		JPEG 100		Gaussian Blur	
	AUC	ACC	AUC	ACC	AUC	ACC	AUC	ACC
max pooling	0.737	0.617	0.767	0.800	0.792	0.810	0.763	0.717
mean pooling	0.790	0.696	0.897	0.700	0.890	0.700	0.853	0.667
ABMIL	0.697	0.600	0.750	0.683	0.805	0.750	0.737	0.650
DSMIL	0.803	0.683	0.833	0.800	0.891	0.825	0.807	0.750
Snuffy	0.745	0.570	0.755	0.650	0.852	0.830	0.769	0.665
LSE pooling	0.857	0.783	0.850	0.800	0.923	0.817	0.863	0.800
Maxsoft pooling	0.970	0.883	0.983	0.867	0.987	0.950	0.963	0.883

P Full-Data Experiments on CAMELYON

For completeness, we report results from training on the full CAMELYON16 (80% train, 20% validation) and CAMELYON17 (60% train, 20% validation, 20% test) with 5-fold cross-validation, rather than the subsets in Appendix F, across MIL pooling architectures in Table 12 and augmentations in Table 3. Even in this setting, while some gains appear, the overall trend holds, and Maxsoft continues to surpass prior methods. Together with Tables 1 and 3, these results indicate that Maxsoft improves generalization regardless of data availability, consistent with its inductive bias.

Q Out-of-Distribution Generalization (CAMELYON17 \rightarrow CAMELYON16)

While our primary focus is in-distribution generalization, motivated by growing interest in out-of-distribution (OOD) generalization, we report OOD generalization results by training on CAMELYON17 [35] and evaluating on CAMELYON16 [78] with no exposure to the latter during training. As shown in Table 13, Maxsoft again overall surpasses prior methods, consistent with the established observation that stronger in-distribution accuracy often correlates with improved OOD performance [104–107]. A notable exception is DINO Domain, where mean pooling outperforms all methods—and even its counterparts with other encoders. We hypothesize that because CAMELYON17 contains positive slides with only a single malignant patch, many models (including Maxsoft) become overly

Table 11: Performance of major MIL pooling functions on CAMELYON17.

Encoder	Method	CAMELYON17					
		Slide				Patch	
		AUC	ACC	F1	ECE	AUC	F1
DINO Natural	LSE pooling	0.683 _{.093}	0.600 _{.087}	0.654 _{.070}	0.361 _{.075}	0.168 _{.291}	0.040 _{.069}
	GM pooling	0.670 _{.017}	0.700_{.000}	0.700 _{.000}	0.166 _{.023}	0.643 _{.002}	0.034 _{.000}
	ISR pooling	0.620 _{.030}	0.600 _{.000}	0.533 _{.000}	0.213 _{.022}	0.453 _{.046}	0.035 _{.060}
	noisy-and pooling	0.667 _{.005}	0.583 _{.028}	0.625 _{.000}	0.095_{.003}	0.655_{.002}	0.049_{.000}
	noisy-or pooling	0.500 _{.000}	0.500 _{.000}	0.000 _{.000}	0.500 _{.000}	0.568 _{.124}	0.021 _{.008}
	smoothmax pooling	0.643 _{.040}	0.583 _{.104}	0.714_{.031}	0.248 _{.055}	0.505 _{.015}	0.000 _{.001}
	Maxsoft pooling	0.710_{.010}	0.650 _{.050}	0.658 _{.083}	0.345 _{.040}	0.312 _{.024}	0.000 _{.000}
DINO Domain	LSE pooling	0.850 _{.070}	0.800 _{.050}	0.833 _{.021}	0.185 _{.066}	0.836 _{.026}	0.499_{.113}
	GM pooling	0.867 _{.023}	0.717 _{.029}	0.822 _{.003}	0.136 _{.106}	0.800 _{.012}	0.303 _{.021}
	ISR pooling	0.863 _{.031}	0.683 _{.029}	0.812 _{.041}	0.198 _{.018}	0.796 _{.011}	0.206 _{.022}
	noisy-and pooling	0.873 _{.006}	0.650 _{.000}	0.842 _{.037}	0.169 _{.028}	0.799 _{.014}	0.055 _{.003}
	noisy-or pooling	0.500 _{.000}	0.500 _{.000}	0.000 _{.000}	0.500 _{.000}	0.646 _{.121}	0.029 _{.011}
	smoothmax pooling	0.803 _{.035}	0.583 _{.076}	0.703 _{.118}	0.302 _{.066}	0.774 _{.028}	0.292 _{.058}
	Maxsoft pooling	0.983_{.055}	0.867_{.076}	0.935_{.072}	0.121_{.051}	0.839_{.019}	0.386 _{.237}
UNI	LSE pooling	0.603 _{.351}	0.617 _{.247}	0.628 _{.400}	0.386 _{.230}	0.709 _{.213}	0.210 _{.364}
	GM pooling	0.730 _{.121}	0.667 _{.144}	0.744 _{.097}	0.114_{.014}	0.700 _{.168}	0.030 _{.012}
	ISR pooling	0.607 _{.235}	0.583 _{.189}	0.712 _{.058}	0.397 _{.176}	0.830_{.130}	0.186 _{.322}
	noisy-and pooling	0.727 _{.138}	0.633 _{.126}	0.769 _{.069}	0.191 _{.033}	0.616 _{.154}	0.036 _{.014}
	noisy-or pooling	0.500 _{.000}	0.500 _{.000}	0.000 _{.000}	0.500 _{.000}	0.642 _{.121}	0.028 _{.020}
	smoothmax pooling	0.537 _{.319}	0.533 _{.104}	0.573 _{.343}	0.346 _{.143}	0.650 _{.263}	0.204 _{.352}
	Maxsoft pooling	0.753_{.071}	0.750_{.205}	0.779_{.040}	0.238 _{.172}	0.786 _{.255}	0.476_{.429}
Prov-GigaPath	LSE pooling	0.933 _{.031}	0.900 _{.050}	0.932 _{.027}	0.095 _{.048}	0.943 _{.004}	0.741 _{.108}
	GM pooling	0.900 _{.050}	0.767 _{.058}	0.825 _{.107}	0.135 _{.029}	0.915 _{.011}	0.374 _{.037}
	ISR pooling	0.930 _{.056}	0.867 _{.029}	0.893 _{.053}	0.144 _{.035}	0.916 _{.026}	0.658 _{.165}
	noisy-and pooling	0.893 _{.006}	0.783 _{.029}	0.881 _{.033}	0.150 _{.070}	0.836 _{.016}	0.138 _{.008}
	noisy-or pooling	0.500 _{.000}	0.500 _{.000}	0.000 _{.000}	0.500 _{.000}	0.757 _{.130}	0.301 _{.021}
	smoothmax pooling	0.917 _{.025}	0.700 _{.050}	0.889 _{.000}	0.245 _{.021}	0.930 _{.004}	0.704 _{.035}
	Maxsoft pooling	1.000_{.000}	0.933_{.029}	1.000_{.000}	0.062_{.022}	0.948_{.009}	0.744_{.034}

sensitive, whereas mean pooling is less affected. This explains mean pooling’s advantage within DINO Domain; its superiority over all encoders, however, warrants further analysis.

R MIL Datasets Out of Pathology Context

To evaluate the broader applicability of our approach beyond the pathology domain, we test our proposed method—alongside prior pathology MIL methods—on classical MIL benchmark datasets. These include MUSK1 and MUSK2, which model molecular binding: each molecule is represented by multiple conformations and is labeled positive if at least one conformation binds to a target protein, though the binding instance is not identified.

Animal-based datasets follow a similar MIL assumption. The Elephant dataset comprises 200 bags (100 positive, 100 negative), where each bag contains instances derived from segmented image features. A bag is labeled positive if it contains at least one elephant instance. The Tiger and Fox datasets follow the same structure, with bags labeled positive if they contain at least one instance of a tiger or fox, respectively. Instance-level labels are unavailable in all cases—only bag-level supervision is provided [108, 109].

Following standard protocol [108, 109], we conduct 10-fold cross-validation with five runs per fold and report the mean and standard deviation for each metric.

Table 12: MIL pooling results on the full CAMELYON16 and CAMELYON17 datasets [78, 35].

Encoder	Method	CAMELYON16				CAMELYON17					
		Slide		Patch		Slide					
		AUC	ACC	F1	ECE	AUC	F1	AUC	ACC	F1	ECE
DINO Natural	max pooling	0.691 _{.259}	0.726 _{.152}	0.533 _{.363}	0.079 _{.067}	0.751 _{.281}	0.248 _{.215}	0.728 _{.082}	0.713 _{.064}	0.631 _{.036}	0.212 _{.016}
	mean pooling	0.619 _{.002}	0.713 _{.000}	0.413 _{.000}	0.092 _{.001}	0.939.000	0.328 _{.001}	0.739 _{.000}	0.720 _{.000}	0.669 _{.005}	0.060.000
	ABMIL	0.857 _{.009}	0.858 _{.012}	0.789 _{.039}	0.143 _{.014}	0.810 _{.051}	0.383 _{.054}	0.795 _{.076}	0.777 _{.075}	0.723 _{.086}	0.255 _{.106}
	DSMIL	0.768 _{.123}	0.762 _{.063}	0.670 _{.147}	0.081 _{.011}	0.545 _{.143}	0.136 _{.046}	0.783 _{.117}	0.767 _{.042}	0.659 _{.131}	0.127 _{.063}
	Snuffy	0.792 _{.088}	0.767 _{.041}	0.706 _{.082}	0.083 _{.020}	0.712 _{.351}	0.216 _{.182}	0.798 _{.058}	0.803 _{.029}	0.707 _{.061}	0.111 _{.039}
	LSE pooling	0.884 _{.004}	0.837 _{.008}	0.826 _{.014}	0.049.006	0.894 _{.004}	0.432.007	0.816 _{.046}	0.793 _{.025}	0.707 _{.013}	0.135 _{.043}
	Maxsoft pooling	0.911.018	0.886.004	0.849.013	0.110 _{.002}	0.899 _{.004}	0.363 _{.018}	0.883.011	0.827.015	0.749.008	0.207 _{.029}
DINO Domain	max pooling	0.986 _{.001}	0.953 _{.004}	0.952 _{.006}	0.038 _{.003}	0.952 _{.001}	0.650 _{.021}	0.918 _{.015}	0.905 _{.021}	0.887 _{.007}	0.053 _{.012}
	mean pooling	0.621 _{.001}	0.667 _{.000}	0.413 _{.000}	0.023.005	0.939 _{.000}	0.294 _{.000}	0.741 _{.001}	0.680 _{.000}	0.657 _{.001}	0.102 _{.004}
	ABMIL	0.958 _{.018}	0.954 _{.008}	0.939 _{.017}	0.046 _{.009}	0.940 _{.008}	0.487 _{.215}	0.889 _{.016}	0.870 _{.010}	0.853 _{.007}	0.111 _{.017}
	DSMIL	0.963 _{.005}	0.948 _{.012}	0.951 _{.005}	0.037 _{.003}	0.465 _{.046}	0.114 _{.014}	0.920 _{.012}	0.883 _{.015}	0.862 _{.019}	0.068 _{.017}
	Snuffy	0.827 _{.023}	0.687 _{.038}	0.759 _{.018}	0.125 _{.041}	0.931 _{.006}	0.407 _{.005}	0.903 _{.008}	0.756 _{.015}	0.813 _{.013}	0.125 _{.019}
	LSE pooling	0.978 _{.000}	0.952 _{.000}	0.951 _{.000}	0.038 _{.004}	0.948 _{.000}	0.731.000	0.931 _{.022}	0.897 _{.011}	0.882 _{.009}	0.034 _{.001}
	Maxsoft pooling	0.993.002	0.961.000	0.958.000	0.037 _{.002}	0.953.001	0.664 _{.030}	0.954.013	0.927.011	0.904.009	0.027.006
UNI	max pooling	0.986 _{.005}	0.974 _{.012}	0.983 _{.012}	0.027 _{.013}	0.966 _{.007}	0.662.087	0.827 _{.038}	0.887 _{.049}	0.864 _{.047}	0.150 _{.039}
	mean pooling	0.583 _{.005}	0.594 _{.004}	0.474 _{.037}	0.370 _{.015}	0.821 _{.003}	0.232 _{.003}	0.839 _{.023}	0.800 _{.026}	0.726 _{.054}	0.338 _{.044}
	ABMIL	0.970 _{.015}	0.948 _{.039}	0.961 _{.016}	0.051 _{.041}	0.928 _{.002}	0.397 _{.115}	0.804 _{.224}	0.837 _{.170}	0.791 _{.234}	0.143 _{.206}
	DSMIL	0.971 _{.012}	0.961 _{.016}	0.956 _{.030}	0.057 _{.046}	0.717 _{.164}	0.177 _{.090}	0.874 _{.145}	0.907 _{.110}	0.840 _{.202}	0.139 _{.120}
	Snuffy	0.931 _{.023}	0.819 _{.049}	0.832 _{.046}	0.050 _{.035}	0.972 _{.001}	0.597 _{.010}	0.932 _{.026}	0.906 _{.025}	0.875 _{.040}	0.095 _{.009}
	LSE pooling	0.983 _{.005}	0.969 _{.008}	0.965 _{.022}	0.032 _{.009}	0.967 _{.003}	0.558 _{.105}	0.846 _{.196}	0.840 _{.176}	0.784 _{.235}	0.183 _{.156}
	Maxsoft pooling	0.998.001	0.974.009	0.990.000	0.026.007	0.973.001	0.595 _{.008}	0.982.006	0.957.006	0.943.002	0.038.0012
Prov-CigatPath	max pooling	0.980 _{.006}	0.969 _{.008}	0.976 _{.017}	0.027 _{.005}	0.964 _{.005}	0.556 _{.032}	0.958 _{.003}	0.943 _{.012}	0.921 _{.021}	0.057 _{.018}
	mean pooling	0.540 _{.030}	0.594 _{.027}	0.379 _{.074}	0.335 _{.024}	0.855 _{.015}	0.231 _{.008}	0.803 _{.002}	0.770 _{.008}	0.676 _{.000}	0.140 _{.004}
	ABMIL	0.982 _{.004}	0.964 _{.024}	0.976 _{.006}	0.034 _{.022}	0.963 _{.002}	0.688.134	0.963 _{.018}	0.950 _{.010}	0.937 _{.017}	0.050 _{.010}
	DSMIL	0.977 _{.009}	0.954 _{.021}	0.945 _{.011}	0.111 _{.047}	0.500 _{.000}	0.125 _{.000}	0.969 _{.014}	0.947 _{.012}	0.943 _{.027}	0.041 _{.011}
	Snuffy	0.962 _{.005}	0.913 _{.021}	0.902 _{.010}	0.121 _{.001}	0.889 _{.032}	0.555 _{.016}	0.942 _{.011}	0.890 _{.043}	0.870 _{.023}	0.076 _{.032}
	LSE pooling	0.980 _{.011}	0.972 _{.012}	0.976 _{.012}	0.028 _{.011}	0.952 _{.020}	0.602 _{.071}	0.968 _{.007}	0.956 _{.013}	0.866 _{.025}	0.271 _{.006}
	Maxsoft pooling	0.988.003	0.974.004	0.979.010	0.025.004	0.965.003	0.556 _{.016}	0.996.007	0.975.007	0.971.000	0.030.002
IN	R ² T-MIL	0.913 _{.007}	0.869 _{.006}	0.852 _{.020}	0.117 _{.013}	/	/	0.792 _{.045}	0.800 _{.014}	0.733 _{.013}	0.202 _{.006}
MLP	R ² T-MIL	0.946 _{.011}	0.891 _{.016}	0.879 _{.015}	0.095 _{.020}	/	/	0.896 _{.006}	0.890 _{.000}	0.887 _{.007}	0.108 _{.004}
UN	PANTHER	0.836 _{.001}	0.806 _{.000}	0.776 _{.002}	0.140 _{.025}	/	/	0.880 _{.004}	0.845 _{.007}	0.825 _{.009}	0.083 _{.021}

As shown in Table 14, Maxsoft pooling outperforms all baselines on nearly every dataset, achieving an AUC of 1.0 on MUSK1. The observed gap between AUC and accuracy is likely due to using a fixed decision threshold rather than tuning per dataset. These results demonstrate the versatility and strong generalization of Maxsoft pooling as a broadly applicable MIL aggregation method.

Table 13: OOD generalization results from models trained on CAMELYON17 and evaluated on CAMELYON16.

Encoder	Method	CAMELYON16					
		Slide				Patch	
		AUC	ACC	F1 Score	ECE	AUC	F1
DINO Natural	max	0.541 _{.085}	0.587 _{.065}	0.458 _{.047}	0.158 _{.083}	0.379 _{.083}	0.000 _{.000}
	mean	0.504 _{.002}	0.442 _{.000}	0.409 _{.000}	0.189 _{.001}	0.388 _{.002}	0.093 _{.001}
	DSMIL	0.527 _{.101}	0.550 _{.047}	0.391 _{.241}	0.235 _{.089}	0.577_{.127}	0.154_{.050}
	LSE	0.579 _{.018}	0.558 _{.013}	0.507_{.053}	0.149_{.033}	0.297 _{.028}	0.006 _{.003}
	Maxsoft	0.594_{.004}	0.589_{.035}	0.502 _{.070}	0.187 _{.115}	0.327 _{.053}	0.007 _{.002}
DINO Domain	max	0.477 _{.053}	0.390 _{.009}	0.413 _{.204}	0.397_{.061}	0.003 _{.004}	0.321 _{.129}
	mean	0.702_{.000}	0.628_{.000}	0.624_{.000}	0.574 _{.000}	0.148_{.000}	0.034 _{.000}
	DSMIL	0.580 _{.080}	0.491 _{.100}	0.505 _{.077}	0.464 _{.127}	0.102 _{.046}	0.130 _{.129}
	LSE	0.474 _{.015}	0.411 _{.021}	0.530 _{.018}	0.578 _{.012}	0.044 _{.013}	0.269 _{.048}
	Maxsoft	0.530 _{.074}	0.437 _{.056}	0.499 _{.036}	0.407 _{.048}	0.001 _{.001}	0.441_{.097}
UNI	max	0.653 _{.204}	0.548 _{.199}	0.532 _{.202}	0.165 _{.122}	0.682 _{.236}	0.253_{.388}
	mean	0.553 _{.008}	0.496 _{.008}	0.558 _{.006}	0.392 _{.008}	0.582 _{.010}	0.153 _{.004}
	DSMIL	0.692 _{.136}	0.636 _{.113}	0.650 _{.113}	0.163 _{.070}	0.500 _{.000}	0.125 _{.000}
	LSE	0.709 _{.174}	0.643 _{.182}	0.617 _{.174}	0.122 _{.000}	0.733 _{.192}	0.218 _{.279}
	Maxsoft	0.778_{.180}	0.765_{.126}	0.700_{.203}	0.120_{.000}	0.869_{.119}	0.157 _{.147}
Prov-GigaPath	max	0.931 _{.050}	0.827 _{.126}	0.878 _{.075}	0.085 _{.105}	0.958_{.003}	0.495 _{.012}
	mean	0.506 _{.004}	0.491 _{.009}	0.537 _{.003}	0.390 _{.009}	0.550 _{.003}	0.141 _{.001}
	DSMIL	0.897 _{.022}	0.716 _{.031}	0.793 _{.017}	0.161 _{.064}	0.500 _{.000}	0.125 _{.000}
	LSE	0.969 _{.009}	0.938 _{.048}	0.959 _{.027}	0.110 _{.052}	0.944 _{.024}	0.561_{.051}
	Maxsoft	0.979_{.001}	0.951_{.027}	0.963_{.015}	0.050_{.050}	0.938 _{.014}	0.526 _{.082}

Table 14: Results of MIL pooling on MUSK1, MUSK2, ELEPHANT, TIGER, and FOX datasets.

Method	MUSK1		MUSK2		ELEPHANT		TIGER		FOX	
	AUC	ACC	AUC	ACC	AUC	ACC	AUC	ACC	AUC	ACC
max pooling	0.881 _{.041}	0.778 _{.000}	0.861 _{.048}	0.767 _{.058}	0.753 _{.049}	0.700 _{.050}	0.879 _{.027}	0.850 _{.050}	0.675 _{.089}	0.567 _{.058}
mean pooling	0.762 _{.083}	0.704 _{.128}	0.750 _{.000}	0.700 _{.000}	0.980_{.000}	0.917 _{.029}	0.896 _{.012}	0.817 _{.058}	0.645 _{.033}	0.650 _{.050}
LSE pooling	0.976 _{.041}	0.815 _{.128}	0.931 _{.024}	0.750 _{.116}	0.957 _{.059}	0.900 _{.050}	0.909 _{.027}	0.883 _{.058}	0.720 _{.093}	0.625 _{.106}
ABMIL	0.976 _{.041}	0.870 _{.111}	0.833 _{.072}	0.800_{.173}	0.973 _{.021}	0.933 _{.029}	0.892 _{.067}	0.850 _{.100}	0.698 _{.104}	0.617 _{.029}
DSMIL	0.929 _{.000}	0.889 _{.000}	0.875 _{.042}	0.733 _{.058}	0.947 _{.015}	0.925 _{.007}	0.919 _{.030}	0.833 _{.029}	0.730 _{.079}	0.617 _{.029}
Snuffy	0.893 _{.051}	0.833 _{.079}	0.917 _{.000}	0.800 _{.141}	0.955 _{.007}	0.825 _{.035}	0.899 _{.057}	0.850 _{.071}	0.756 _{.059}	0.625 _{.035}
Maxsoft pooling	1.000_{.000}	0.889_{.000}	0.954_{.024}	0.767 _{.058}	0.958 _{.036}	0.950_{.050}	0.950_{.011}	0.900 _{.029}	0.760_{.033}	0.650_{.029}

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The scope of the effectiveness and main claims are clearly demonstrated in the abstract and introduction

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Yes, we provided formal assumptions of our theoretical analysis in the Appendix A and made clear reasoning on why we consider them true and how general our bounds are. We also discussed the limitations of our work in Section 7.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: Refer to Appendix A.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We talk about all models and datasets in Section 5, our method in the Section 4, extra details on hyperparameters, splits and training details are provided in Appendix F.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Codes and information of datasets that are constructed or reused in the paper are included in the GitHub repository.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The paper clearly describes the datasets, models evaluated, and evaluation metrics in Section 5.2 and Appendices E, F, and H.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: All experimental results are calculated over 5 different random initializations and their mean and std are reported.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Yes. We specified our compute resources in the supplementary materials. Additionally we compared our compute time relative to previous methods in the main tables.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: All authors read and confirm that the research in the paper conform with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: There is no societal impact of the work performed.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper has no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Every asset that we utilized for our implementations have been appropriately referenced, both within the paper itself and in the code (if needed). Although we did not specify the names of their respective licenses, you can find these details on the webpages we've cited.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: This paper makes the code and our self-pretrained weights and embeddings available through the external link.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This study uses only publicly available, de-identified datasets: CAMELYON16, CAMELYON17, TCGA-Lung, and SICAP-MIL. We performed no new data collection or interaction with human participants and conducted no crowdsourcing. Data were used under the datasets' terms; to our knowledge, original providers obtained the necessary approvals and/or consent. We made no attempt to re-identify individuals and followed dataset usage policies.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The study uses only publicly available, de-identified datasets (CAMELYON16, CAMELYON17, TCGA-Lung, SICAP-MIL). We conducted no new data collection or interaction with human participants; thus, this work does not constitute human-subjects research and IRB review was not required. We followed each dataset’s usage terms and made no attempt to re-identify individuals.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The paper only uses LLMs for text and code editing.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.