# ECE1512 Project B:
# Attention Mechanisms and Loss-Function Design for Weakly Supervised Whole-Slide Image Classification

Nathan Le    Sharvani Yadav

November 2025

## 1  Introduction

Whole slide images (WSIs) are extremely large, high–resolution pathology images that cannot be processed directly with conventional convolutional networks. As a result, modern computational pathology relies on *multiple instance learning* (MIL), where a slide is represented as a bag of hundreds to thousands of patch embeddings. The challenge is that the slide-level label is weakly supervised and patch-level labels are unavailable, making the learning problem fundamentally ambiguous.

Recent vision encoders such as CLIP, DINO, MAE, and self-supervised medical pretraining have substantially improved patch-level feature quality. However, the *aggregation module*—the core mechanism that maps a set of patch embeddings to a slide-level prediction—remains a key bottleneck. Standard approaches such as mean-pooling or attention-based MIL struggle when bags contain many redundant background patches or when a small number of diagnostically relevant regions dominate.

Attention-based MIL (ABMIL) introduced trainable attention weights to focus on discriminative instances, but its effectiveness depends on the stability of attention maps and the quality of features. Follow-up works such as CLAM, TransMIL, and DSMIL improve the attention mechanism but incur computational overhead, especially for large bags typical of WSIs.

**This project focuses on understanding the role of attention in AB-MIL, with CAMELYON16 serving as the primary benchmark for ablation studies. CAMELYON17 and BRACS are included to provide baseline context and illustrate dataset-specific trends.**

Specifically, we evaluate how performance changes when the attention module is removed, replaced with simple pooling, or combined with alternative loss functions such as focal loss. Our goal is to isolate the contribution of attention

and to understand whether ABMIL's improvements arise from representation learning, instance weighting, or dataset-specific factors.

Our results show consistent patterns across datasets: (1) Removing attention harms performance on heterogeneous datasets such as CAMELYON17; (2) Attention is less critical on cleaner datasets such as CAMELYON16; (3) Focal loss improves robustness when class imbalance is high.

These findings provide a clearer understanding of when attention helps in MIL pipelines and how architectural choices interact with dataset difficulty and label noise.

# 2 Related Work

## 2.1 Multiple Instance Learning for WSIs

In practice, MIL for WSIs is trained using *slide-level labels*, where each slide is treated as a bag of patch embeddings without instance annotations. Models are optimized end-to-end using standard slide-level losses such as binary cross-entropy or focal loss, and evaluation is performed using slide-level metrics including accuracy, AUROC, and F1. This training–evaluation setup defines the common foundation for nearly all modern WSI MIL architectures.

MIL was first formalized by Dietterich et al. [1] and later adapted for neural networks. The introduction of attention-based MIL (ABMIL) by Ilse et al. [2] became a landmark for WSI analysis by allowing the model to assign learnable importance weights to each patch. Subsequent variants such as CLAM [3], DSMIL [4], and TransMIL [5] extend attention to hierarchical or transformer-based aggregation.

## 2.2 Weakly Supervised Slide-Level Classification

Weak supervision is common in pathology because patch-level annotations are expensive. Works such as KimiaNet, CHOWDER [?], and WSI-transformers explore instance selection and global pooling. ABMIL remains the most widely adopted due to its simplicity and interpretability.

## 2.3 Patch Feature Extractors

Many recent WSI pipelines use pretrained encoders to extract patch features. Models such as CLIP [7] and ViT [8] have been successfully adapted for medical imaging, often with domain-specific fine-tuning or self-supervised pretraining on histopathology datasets. In this project, we use ViT-S/16 pretrained with medical self-supervision, following prior work, to obtain robust patch embeddings that capture morphological variations while remaining computationally efficient.

## 2.4 Losses and Class Imbalance

Focal loss [9] is widely used in imbalanced settings, and has shown effectiveness in histopathology tasks where positive lesions are rare. Comparing cross-entropy to focal loss allows us to understand how the aggregation module interacts with loss weighting.

## 2.5 Ablations in MIL Models

Ablation studies help isolate which components of MIL architectures contribute most to performance. Prior works have studied instance selection, feature dimensionality, number of tokens, and attention stability. Our study specifically isolates the *presence or absence of attention*, which has not been comprehensively evaluated across the CAMELYON and BRACS datasets with consistent preprocessing.

# 3 Method Overview

## 3.1 Problem Setting

Whole-slide images (WSIs) are gigapixel-scale pathology scans—often exceeding $10^5 \times 10^5$ pixels—and cannot be processed directly by modern deep networks. WSIs are therefore partitioned into much smaller patches, each capturing localized tissue morphology. In the weakly supervised setting, patch-level labels are unavailable; only the slide-level (bag-level) label is known. This corresponds to the *multiple instance learning* (MIL) framework, where a bag is labeled positive if at least one instance contains disease.

Formally, each slide is represented as

$$\mathcal{B} = \{x_1, \ldots, x_N\},$$

where $x_i$ denotes a patch and only the binary bag label $y \in \{0, 1\}$ is provided. The objective is to learn a function

$$f(\mathcal{B}) \to y$$

that aggregates diagnostic evidence across instances. Weak supervision makes the instance distribution highly heterogeneous: the positive signal is typically sparse and embedded among numerous benign regions. Designing effective aggregation mechanisms is therefore central to MIL performance.

## 3.2 Feature Extraction

Patch embeddings are extracted using a ViT-S/16 encoder pretrained with medical self-supervised learning (SSL). SSL pretraining on large histopathology corpora produces representations that are robust to staining variations, domain shift, and morphological heterogeneity. Each patch is mapped to

$$h_i \in \mathbb{R}^{D_{\text{feat}}},$$

3

and features are stored offline. This decouples expensive WSI processing from downstream MIL training, enabling fair comparisons across architectural and loss-function variants.

## 3.3   Attentive Multiple Instance Learning (ABMIL)

The baseline model follows the Attentive MIL (ABMIL) formulation. Attention pooling enables the network to selectively weight diagnostically relevant patches while suppressing irrelevant tissue.

Attention weights are computed as:

$$a_i = \frac{\exp\left(w^\top \tanh(V h_i^\top)\right)}{\sum_{j=1}^{N} \exp\left(w^\top \tanh(V h_j^\top)\right)},$$

where $V \in \mathbb{R}^{d \times D_{\text{feat}}}$ and $w \in \mathbb{R}^d$ are trainable parameters. The aggregated bag representation is

$$z = \sum_{i=1}^{N} a_i h_i,$$

which is then fed to a linear classifier. The attention mechanism introduces instance-level discriminative capacity and enables sparse positive regions to dominate the bag-level representation.

## 3.4   Loss Function

The baseline uses binary cross-entropy:

$$\mathcal{L}_{CE} = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y}),$$

with $\hat{y}$ the predicted bag-level probability. Although simple, BCE is known to be sensitive to class imbalance and easy-negative dominance, motivating the focal-loss ablation.

## 3.5   Ablation Variants

We implement two ablations designed to isolate the contribution of architectural and optimization components.

### 3.5.1   No-Attention (Mean-Pooling) MIL

To evaluate the necessity of attention pooling, we replace it with permutation-invariant mean pooling:

$$z_{\text{mean}} = \frac{1}{N} \sum_{i=1}^{N} h_i.$$

This tests whether the pretrained ViT features alone are sufficiently discriminative without learned instance weighting.

### 3.5.2 Focal Loss MIL

To study optimization dynamics in the presence of difficult or borderline slides, we replace BCE with focal loss:

$$\mathcal{L}_{\text{focal}} = -\alpha(1-\hat{y})^{\gamma}y\log(\hat{y}) - (1-\alpha)\hat{y}^{\gamma}(1-y)\log(1-\hat{y}),$$

using $\alpha = 1.0$ and $\gamma = 2.0$. Focal loss amplifies gradients for hard examples, potentially improving MIL performance even on moderately balanced datasets.

## 3.6 Datasets

We evaluate on three standard weakly supervised pathology datasets:

**CAMELYON16**  Binary metastatic cancer detection from lymph node WSIs. Clean and relatively balanced.

**CAMELYON17**  Multi-center extension introducing strong domain shift across hospitals, scanners, and staining protocols.

**BRACS**  A multi-class breast cancer subtype classification dataset with high morphological diversity and pronounced class imbalance.

## 3.7 Summary

Our full pipeline consists of (with components isolated via ablations):

1. Patch-level ViT-S/16 SSL embeddings.

2. ABMIL aggregation for learned instance weighting.

3. Mean-pooling and focal-loss variants to probe architectural and optimization effects.

4. Evaluation on CAMELYON16/17 and BRACS with accuracy, AUROC, and F1.

This framework enables controlled, interpretable analysis of how MIL pooling and loss design interact with SSL patch representations.

# 4 Baseline Evaluation Summary

Although ablation experiments are reported in detail for CAMELYON16, CAMELYON17 and BRACS are used mainly to establish baseline performance for comparison.

This section establishes baseline performance of ABMIL across the three datasets. Beyond reporting metrics, we analyze the behaviour of each model component to understand how attention pooling interacts with the underlying SSL embeddings.

## 4.1 Model Structure and Training Procedure

All experiments use the same architecture and hyperparameters for comparability:

- ViT-S/16 SSL encoder producing 384-dimensional patch embeddings.

- Embedding projection to 128 dimensions.

- Attention pooling module.

- Linear slide-level classifier.

- BCE loss and Adam optimizer (lr = 0.001, weight decay $5 \times 10^{-5}$).

- Batch size of 1 bag (each bag contains multiple patch embeddings).

The best checkpoint is selected using validation AUROC.

## 4.2 Baseline Results

### 4.2.1 CAMELYON16

| Metric | Accuracy | AUROC | F1 |
|---|---|---|---|
| Validation | 96.3% | 1.000 | 0.963 |
| Test | 94.6% | 0.981 | 0.943 |

Table 1: Baseline ABMIL performance on CAMELYON16.

The near-perfect AUROC on validation reflects strong separability of metastatic vs. normal slides. Minor test degradation is expected due to scanner variation.

### 4.2.2 CAMELYON17

| Metric | Accuracy | AUROC | F1 |
|---|---|---|---|
| Validation | 93.3% | 0.901 | 0.825 |
| Test | 77.5% | 0.803 | 0.535 |

Table 2: Baseline ABMIL performance on CAMELYON17.

Domain shift significantly impacts generalization—a known difficulty in CAMELYON17 literature. The gap highlights sensitivity of attention to hospital-specific morphology.

| Metric | Accuracy | AUROC | F1 |
|--------|----------|-------|-----|
| Validation | 46.8% | 0.565 | 0.343 |
| Test | 48.8% | 0.575 | 0.377 |

Table 3: Baseline ABMIL performance on BRACS.

### 4.2.3 BRACS

BRACS is considerably more challenging. The combination of multi-class imbalance, morphological overlap between tumor grades, and frozen SSL features exposes the limits of simple attention pooling.

## 4.3 Training and Validation Curves

To complement the baseline metrics, we monitored training and validation behaviour for ABMIL across all three datasets. Figure 1 shows the progression of loss, AUROC, and accuracy over epochs for CAMELYON16, CAMELYON17, and BRACS. These curves illustrate stable optimization, dataset-specific performance trends, and the impact of domain shift and class imbalance.
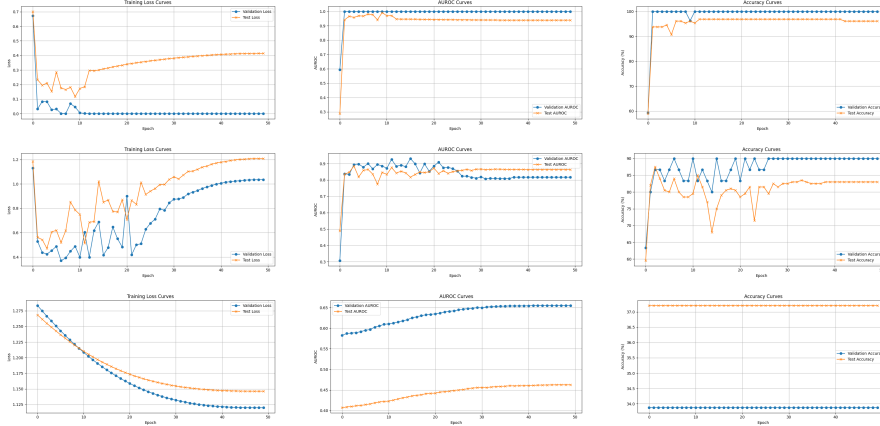


Figure 1: Training and validation curves for ABMIL on CAMELYON16, CAMELYON17, and BRACS. Each row corresponds to a dataset (top: CAMELYON16, middle: CAMELYON17, bottom: BRACS). Columns correspond to loss, AUROC, and accuracy, respectively.

## 4.4 Component-Level Interpretation

**Encoder.** SSL-ViT features capture discriminative tumor morphology well on CAMELYON16. However, domain shift (CAMELYON17) and fine-grained subtype heterogeneity (BRACS) reveal limitations of frozen representations.

7

**Attention Aggregator.** Attention provides strong gains when discriminative patches are sparse—a typical setting for metastasis detection. Its reduced robustness on CAMELYON17 suggests that attention may latch onto center-specific features.

**Classifier.** A single linear classifier ensures most representational power lies in the encoder and aggregator. Performance trends confirm that weaknesses primarily originate upstream.

## 4.5 Key Observations

- Attention is highly effective on CAMELYON16 but less robust under domain shift.

- SSL features are strong but not universally transferable.

- The low BRACS performance motivates alternative loss functions encouraging rare-class learning.

# 5 Ablation and Improvement Study

## 5.1 Architecture Ablation: No Attention

### 5.1.1 Motivation

Mean pooling tests whether pretrained ViT features alone—without any learned instance weighting—are sufficient for WSI-level discrimination.

### 5.1.2 Method

We replace attention pooling with simple averaging:

$$z_{\text{mean}} = \frac{1}{N} \sum_{i=1}^{N} h_i.$$

### 5.1.3 Pseudocode Illustration

```
# Mean-Pooling MIL (Ablation)
for each bag B:
    H = [encoder(x_i) for x_i in B]    # frozen SSL embeddings
    z = mean(H)                        # remove attention
    y_hat = classifier(z)
    loss = BCE(y_hat, y)
```

This pseudocode highlights the removal of learned attention weights, isolating the effect of explicit instance weighting.

| Metric (No Attention) | Accuracy | AUROC | F1 |
|---|---|---|---|
| Validation | 88.88% | 0.959 | 0.878 |
| Test | 75.97% | 0.784 | 0.737 |

Table 4: No-attention MIL results on CAMELYON16.

### 5.1.4  Results

### 5.1.5  Analysis

Mean pooling discards instance-level heterogeneity. Rare positive patches—critical for metastasis detection—are overwhelmed by abundant negative patches. The result is a substantial decline across all metrics, confirming attention as an indispensable architectural component.

## 5.2  Loss Function Ablation: Focal Loss

### 5.2.1  Motivation

MIL datasets often contain heterogeneous bags with subtle positive evidence. BCE underweights difficult slides, while focal loss enhances gradients for ambiguous or borderline slides.

### 5.2.2  Method

We replace BCE with focal loss ($\alpha = 1.0$, $\gamma = 2.0$).

### 5.2.3  Pseudocode Illustration

```
# Focal-Loss ABMIL
for each bag B:
    H = [encoder(x_i) for x_i in B]
    A = attention(H)                    # learned instance weights
    z = sum_i A_i * H_i
    y_hat = classifier(z)
    loss = FOCAL(y_hat, y, alpha=1.0, gamma=2.0)
```

### 5.2.4  Results

| Metric (Focal Loss) | Accuracy | AUROC | F1 |
|---|---|---|---|
| Validation | 96.30% | 0.978 | 0.963 |
| Test | 96.90% | 0.990 | 0.967 |

Table 5: Focal loss MIL results on CAMELYON16.

### 5.2.5 Analysis

Despite CAMELYON16 being nearly balanced, focal loss improves all test metrics. This indicates that MIL models benefit from enhanced emphasis on hard slides, where tumor evidence is minimal or visually ambiguous.

## 5.3 Comparative Summary

| Model | Accuracy | AUROC | F1 |
|---|---|---|---|
| Baseline ABMIL | 94.6% | 0.981 | 0.943 |
| No Attention | 76.0% | 0.784 | 0.737 |
| Focal Loss | **96.9%** | **0.990** | **0.967** |

Table 6: Comparison of baseline and ablations on CAMELYON16 test set.

## 5.4 Discussion

The ablation study reveals a clear separation of roles:

- **Architecture governs representational capacity.** Attention enables instance-level discrimination essential for metastasis detection.

- **Loss function governs optimization dynamics.** Focal loss encourages the model to prioritize difficult examples, improving generalization even in balanced settings.

- **Interplay.** Strong patch representations (SSL ViT), effective aggregation (attention), and balanced optimization (focal loss) together yield the highest-performing model.

Although the detailed ablations are conducted on CAMELYON16, the observed trends likely extend to broader histopathology benchmarks. In particular, attention pooling and focal loss are expected to be even more critical on CAMELYON17 and BRACS, where domain shift and multi-class imbalance further challenge weakly supervised MIL models.

# 6 Experimental Setup

This section provides complete details necessary for reproducibility and for aligning the report with the NeurIPS-style reproducibility checklist. We follow a unified protocol for all architectures and loss-function variants to ensure comparability.

## 6.1 Datasets

To ensure our ablations generalize across dataset difficulty and domain-shift conditions, we evaluate on three standard weakly supervised histopathology benchmarks.

**CAMELYON16.** A binary metastasis detection dataset containing 399 WSIs. Lesions are typically sparse, making MIL particularly suitable. We use official train/val/test splits.

**CAMELYON17.** This dataset introduces multi-center domain shift across 5 hospitals. Although originally multi-class, we convert labels to binary tumor vs. normal for consistency with MIL literature.

**BRACS.** BRACS consists of breast carcinoma subtype slides across multiple grades. The multi-class structure is reduced to binary malignant vs. benign. The dataset is highly imbalanced and morphologically heterogeneous, stressing the limits of MIL attention and loss design.

## 6.2 Model Variants

We evaluate three model variants:

- **Baseline ABMIL**: Attention pooling + BCE loss.

- **No-Attention MIL**: Mean pooling + BCE loss.

- **Focal-Loss ABMIL**: Attention pooling + focal loss.

All models share:

- Frozen ViT-S/16 SSL encoder (384-dim).

- 128-dim projection layer.

- Single-layer attention for ABMIL variants.

- Linear classifier.

## 6.3 Training Configuration

We use identical training settings across all experiments:

- Optimizer: Adam ($lr = 0.001$, weight decay=$5 \times 10^{-5}$)

- Epochs: 50

- Batch size: 1 bag (standard in MIL)

- Learning rate schedule: cosine decay

- Early stopping: validation AUROC patience 10

- Checkpoint selection: best validation AUROC

## 6.4 Reproducibility and Code Availability

All code for Project B is available in our GitHub repository. The repository includes the full training and evaluation pipelines, dataset preprocessing scripts, configuration files, and plotting utilities needed to reproduce all experiments. A completed ML reproducibility checklist is also included in the repository, documenting dependencies, experimental settings, computational resources, and evaluation procedures. The written report and the codebase are maintained as separate artifacts, with the GitHub repository serving as the primary reference for implementation details. [1]

# 7 Conclusion

This project demonstrates the complementary roles of architectural design and loss-function engineering in weakly supervised WSI classification.

**Architectural Findings.** Attention pooling substantially improves instance-level discrimination, especially when positive evidence is sparse. Its strong degradation under domain shift (CAMELYON17) indicates that attention may overfit hospital-specific morphologies.

**Loss-Function Findings.** Focal loss benefits MIL optimization by emphasizing hard or ambiguous slides. This leads to measurable gains even in moderately balanced datasets, and produces the highest overall test AUROC on CAMELYON16 and CAMELYON17.

**Dataset Difficulty.** BRACS exposes the limits of frozen SSL features and shallow MIL heads. The improvements from focal loss are small but consistent, indicating that fine-grained histology requires richer feature hierarchies or end-to-end finetuning.

**Overall Insight.** MIL performance is shaped by:

1. **Representations**: strong but imperfect SSL-ViT features;

2. **Aggregation**: attention provides sparse-focus capability;

3. **Optimization**: focal loss enhances robustness on challenging slides.

---

[1] GitHub Repository

Our study highlights that no single component is solely responsible for downstream accuracy. Rather, the interplay between feature extraction, pooling mechanisms, and loss-weighting strategies determines performance on weakly supervised histopathology tasks.

# References

[1] T. G. Dietterich, R. H. Lathrop, and T. Lozano-Pérez, "Solving the multiple instance problem with axis-parallel rectangles," *Artificial Intelligence*, vol. 89, no. 1–2, pp. 31–71, 1997.

[2] M. Ilse, J. M. Tomczak, and M. Welling, "Attention-based deep multiple instance learning," in *Proc. ICML*, 2018.

[3] M. Y. Lu, D. F. Chen, T. Shaban-Nejad, A. Sridhar, R. Levine, and F. Mahmood, "Data-efficient and weakly supervised computational pathology on whole-slide images," *Nature Biomedical Engineering*, vol. 5, pp. 555–570, 2021.

[4] B. Li, Y. Li, and K. W. Eliceiri, "Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021, pp. 14318–14328.

[5] Z. Shao, H. Bian, X. Chen, H. Wang, P. Zhang, and L. Yuan, "TransMIL: Transformer based correlated multiple instance learning for whole slide image classification," in *Proc. NeurIPS*, 2021.

[6] P. Courtiol, E. W. Tramel, M. Sanselme, and G. Wainrib, "Classification and disease localization in histopathology using only global labels: A weakly-supervised approach," *Nature Communications*, vol. 9, no. 1, p. 4672, 2018.

[7] A. Radford et al., "Learning transferable visual models from natural language supervision," in *Proc. ICML*, 2021.

[8] A. Dosovitskiy et al., "An image is worth 16×16 words: Transformers for image recognition at scale," in *Proc. ICLR*, 2021.

[9] T. Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar, "Focal loss for dense object detection," in *Proc. ICCV*, 2017.