Few-Shot Localization and Anomaly Detection in Medical Images via Prompt-Tuned CLIP Using Healthy-Only Training for Diverse Clinical Modalities

Jayan Ghimire Software Engineer and Independent AI Researcher Leapfrog Technology Kathmandu, Nepal jghimire.034@gmail.com

June 1, 2025

ABSTRACT

In this work, we propose a novel and custom-designed framework for **few-shot anomaly detection and localization** in **medical imaging**, inspired by the foundational principles of *Contrastive Language-Image Pretraining (CLIP)* but distinctly engineered to address the unique challenges of clinical data. Our architecture fuses **visual** and **textual modalities** through a parameter-efficient **prompt tuning** strategy, which optimizes a set of learnable **textual tokens** to robustly characterize the distribution of **healthy images** and facilitate the identification of deviations indicative of anomalies without requiring explicit anomaly training data.

By leveraging only **healthy training samples**, the model learns a comprehensive **representation space** that captures normal anatomical and physiological variability. Anomalies are detected as **out-of-distribution** patterns through patch-level embedding comparisons, enabling precise **spatial attention maps** and **bounding box predictions** for interpretable **anomaly localization** critical for clinical decision support.

We validate our approach on three heterogeneous medical imaging datasets—**chest X-ray** (*CheXpert*), **brain MRI** (*BrainMRI*), and **breast ultrasound** (*BUSI*)—under multiple few-shot prompt tuning configurations with ($\mathbf{K} = \{4, 8, 16, 32\}$) samples of healthy images. Our method consistently achieves state-of-the-art **anomaly detection** and **localization** performance, demonstrating superior generalization and robustness across diverse clinical modalities and limited supervision regimes. This establishes our framework as an effective and annotation-efficient solution for advancing medical anomaly analysis through multi-modal representation learning.

1 Introduction

Medical imaging is fundamental to disease diagnosis, enabling visualization of internal anatomy and pathological conditions. Despite advances in deep learning, the development of reliable anomaly detection systems is hindered by the scarcity of annotated abnormal data and the wide variability of pathological appearances. Traditional supervised methods often require extensive labeled datasets, which are costly and time-consuming to obtain in clinical settings.

This work addresses the problem of **few-shot anomaly detection and localization** in medical images trained exclusively on **healthy samples**, where annotated anomaly data is unavailable or limited. We hypothesize that leveraging the rich multi-modal representations of CLIP through prompt tuning can effectively model normal anatomy and detect deviations indicative of pathology, while providing interpretable localization cues.

This Paper's Contributions:

- We propose a novel **prompt tuning** framework that efficiently adapts the pretrained CLIP model to the medical domain, enabling robust **few-shot anomaly detection** by exclusively learning from healthy image distributions. This approach leverages parameter-efficient tuning of textual prompts to capture subtle pathological deviations without fine-tuning the entire visual encoder, significantly reducing computational cost and overfitting risk in low-data regimes.
- We develop a **patch-level embedding and spatial attention mechanism** that exploits CLIP's multi-scale visual representations to produce high-resolution, interpretable heatmaps. This facilitates precise **anomaly localization** by highlighting pathological regions as deviations from the learned healthy anatomy, supporting explainability critical for clinical adoption.
- We introduce a comprehensive evaluation protocol across three clinically relevant and heterogeneous imaging modalities—*chest X-ray (CheXpert), brain MRI,* and *breast ultrasound (BUSI)*—under varying few-shot settings (K = {4, 8, 16, 32}).
- We provide an in-depth analysis of prompt design and the impact of few-shot sample size on detection and localization performance, offering novel insights into optimizing vision-language models for medical anomaly detection tasks.
- To promote reproducibility and future research, we outline an extensible framework that can be adapted to other medical imaging domains and anomaly detection scenarios, highlighting the potential of multi-modal representation learning in low-data clinical settings.

2 Literature Review

The domain of **medical image anomaly detection** has rapidly evolved with the rise of **deep learning** techniques, particularly **convolutional neural networks (CNNs)** [1] and, more recently, **transformer-based** models [2, 3]. Conventional **supervised learning** methods require extensive **annotated datasets** covering a wide spectrum of **pathologies**, which is impractical due to the **rarity of certain anomalies** and **high annotation costs** [4, 5].

To circumvent these limitations, **few-shot learning** paradigms have been introduced, aiming to enable models to **generalize from scarce labeled samples** [6, 7]. However, most existing few-shot frameworks assume access to both **normal** and **abnormal training samples**, a condition often unmet in clinical practice where only *healthy images* are abundantly available [8].

Unsupervised and **self-supervised** approaches, including **autoencoders** [9, 10], **GAN-based** methods [11, 12], and **memory-augmented networks** [13, 14], leverage **healthy-only training** to model **normative anatomy** and detect **deviations** as **anomalies**. Despite promising results, these models frequently exhibit suboptimal **localization accuracy** and lack **interpretability**, often relying on **reconstruction** or **heuristic anomaly scores** [15].

The advent of **vision-language models** such as *CLIP* [16] has revolutionized **multi-modal representation learning** by aligning **visual** and **textual information**. While CLIP demonstrates impressive **zero-shot** and **few-shot capabilities** in **natural image domains** [17, 18], its direct transfer to **medical imaging** remains challenging due to **domain shift** and the scarcity of **medical text annotations** [19].

Recent attempts to adapt CLIP for **medical anomaly detection** [20, 21] often require **fine-tuning** on **abnormal data** or lack robust mechanisms for precise **anomaly localization**. Thus, a clear **research gap** exists for methods that leverage CLIP's **multi-modal representations** for **few-shot anomaly detection** and **localization**, trained exclusively on *healthy images*.

Our work addresses this gap by introducing a **prompt tuning** strategy that adapts CLIP to **medical modalities** with **minimal labeled supervision**, alongside a **patch-level embedding** approach enabling **high-resolution** and **interpretable spatial attention maps** for accurate **anomaly localization**. This synergy of **few-shot learning**, **healthy-only training**, and **vision-language adaptation** marks a novel contribution to **medical anomaly detection** across diverse **clinical imaging modalities**.

3 Methodology

3.1 Proposed System Architecture



Figure 1: Few-shot anomaly detection system diagram using CLIP and prompt tuning.

4 Methodology

This research presents a novel **few-shot learning** framework for **anomaly detection** and **localization** in medical images across diverse clinical modalities. Our approach leverages a **prompt-tuned CLIP** architecture with a **healthy-only**

training paradigm, enabling robust **zero-shot** and **few-shot** generalization to unseen pathological conditions without requiring extensive labeled anomalous data during training.

4.1 Input Processing and Modality-Specific Preprocessing

4.1.1 Multi-Modal Image Acquisition Pipeline

The **Input Module** serves as the foundational component of our framework, incorporating comprehensive support for heterogeneous medical imaging formats including **DICOM**, **PNG**, and **NIfTI** standards. This multi-format compatibility ensures seamless integration across diverse clinical workflows and imaging protocols, addressing the inherent heterogeneity in medical imaging data acquisition systems.

4.1.2 Adaptive Preprocessing Framework

Given the inherent variability in medical imaging modalities, we implement **modality-specific preprocessing** strategies to standardize input characteristics while preserving clinically relevant anatomical and pathological features:

Intensity Normalization: Applied using **z-score standardization** and **histogram equalization** techniques to mitigate inter-scanner variability and ensure consistent pixel intensity distributions across imaging sessions. The z-score normalization is mathematically defined as:

$$I_{norm}(x,y) = \frac{I(x,y) - \mu_I}{\sigma_I} \tag{1}$$

where I(x, y) represents the original intensity at pixel location (x, y), μ_I and σ_I are the mean and standard deviation of the image intensity distribution, respectively. This normalization is crucial for maintaining feature consistency across different acquisition protocols and scanner manufacturers, as it eliminates the bias introduced by varying scanner calibrations and imaging parameters.

Noise Reduction: Implemented through **wavelet denoising** techniques, specifically employing **Daubechies wavelets** to preserve edge information while suppressing high-frequency noise artifacts commonly present in medical images. The wavelet denoising process is formulated as:

$$I_{denoised} = \mathcal{W}^{-1}[\mathcal{T}_{\lambda}(\mathcal{W}[I])] \tag{2}$$

where W and W^{-1} represent the forward and inverse wavelet transforms, and \mathcal{T}_{λ} is the soft thresholding operator with threshold λ . The choice of Daubechies wavelets is justified by their optimal time-frequency localization properties, making them particularly suitable for preserving sharp anatomical boundaries while removing noise artifacts that could interfere with subtle pathological feature detection.

Contrast Enhancement: Utilizing **Contrast Limited Adaptive Histogram Equalization (CLAHE)** to enhance local contrast while preventing over-amplification of noise artifacts. This technique is particularly crucial for low-contrast modalities such as **ultrasound** and **mammography**, where subtle pathological changes require enhanced visibility.

4.1.3 Data Augmentation Strategy

To enhance model robustness and generalization capability while working with limited healthy training samples, we employ a comprehensive **data augmentation** pipeline:

Geometric Transformations: Including random flips and rotations to simulate natural anatomical variations and imaging perspective changes commonly encountered in clinical practice.

Intensity Jitter: Applied to simulate realistic variations in imaging conditions, scanner settings, and patient-specific factors that affect image intensity distributions.

Elastic Deformations: Specifically implemented for **ultrasound** and **MRI** modalities to simulate natural tissue deformation and movement artifacts, enhancing the model's robustness to anatomical variations.

4.2 Prompt Engineering and Textual Adapter Framework

4.2.1 Prompt Template Generation

Our **Prompt Template Generator** creates domain-specific medical prompts tailored to the **healthy-only training** paradigm. The prompt generation process is mathematically formulated as:

$$P_{template} = f_{gen}(\mathcal{A}, \mathcal{C}, \mathcal{M}) \tag{3}$$

where \mathcal{A} represents anatomical terms, \mathcal{C} denotes clinical conditions, \mathcal{M} indicates imaging modality, and f_{gen} is the template generation function. Custom medical templates such as "No signs of {condition}" and "Normal {anatomy} appearance" are dynamically generated to establish robust **text-image correspondence** for healthy tissue representations. The choice of negative prompt formulation ("No signs of...") is theoretically justified as it creates a more discriminative embedding space by explicitly defining the absence of pathological features, thereby enhancing the model's ability to detect deviations from healthy baselines.

4.2.2 Learnable Prompt Tuning Module

We implement a **learnable prompt tokens** mechanism through our **Prompt Tuning Module**, which optimizes continuous prompt representations during training. The learnable prompts are formulated as:

$$\mathbf{p}_{learnable} = [\mathbf{p}_1, \mathbf{p}_2, ..., \mathbf{p}_L] \in \mathbb{R}^{L \times d}$$
(4)

where L represents the prompt length and d is the embedding dimension. The optimization objective for prompt tuning is:

$$\mathcal{L}_{prompt} = -\log P(y_{healthy} | \mathbf{I}, \mathbf{p}_{learnable}) \tag{5}$$

This approach enables the model to learn optimal textual representations for medical concepts without requiring extensive manual prompt engineering, adapting to the specific characteristics of each medical domain. The continuous optimization of prompt tokens allows for fine-grained adaptation to domain-specific medical terminology and imaging characteristics.

4.2.3 Medical Ontology Integration

The **Medical Ontology Mapper** establishes semantic connections to established medical knowledge bases including **UMLS**, **SNOMED**, and **RadLex** terminologies. This integration ensures that our prompt generation process aligns with standardized medical vocabulary and enhances the semantic understanding of anatomical and pathological concepts.

4.3 Enhanced CLIP Encoder Architecture

4.3.1 Visual Encoder Enhancement

Our framework employs a **CLIP backbone** with significant architectural enhancements tailored for medical image analysis. The selection of CLIP architecture is theoretically justified by its superior cross-modal alignment capabilities and pre-trained representations that bridge visual and textual modalities. The **Visual Encoder**, utilizing either **Vision Transformer (ViT)** or **ResNet** architectures, is augmented with a **Visual Token Refinement Layer** that specifically focuses on medical image characteristics. The refinement process is mathematically expressed as:

$$\mathbf{z}_{refined} = Attention(\mathbf{Q}_{medical}, \mathbf{K}_{clip}, \mathbf{V}_{clip})$$
(6)

where $\mathbf{Q}_{medical}$ represents medical-domain specific queries, while \mathbf{K}_{clip} and \mathbf{V}_{clip} are the key and value matrices from the pre-trained CLIP encoder. This architectural choice enables fine-grained feature extraction relevant to pathological changes while leveraging the robust pre-trained representations from large-scale vision-language training.

4.3.2 Medical Language Adapter Integration

The **Text Encoder** component is enhanced through integration of a **Medical Language Adapter**, which has been pretrained on extensive medical report corpora. This adapter enables the model to better understand medical terminology, clinical descriptions, and the nuanced language used in radiological reporting, bridging the gap between general language understanding and domain-specific medical communication.

4.4 Embedding Alignment and Multi-Modal Projection

4.4.1 Patch-Level Feature Extraction

Our approach implements **patch-level embedding extraction** to enable fine-grained spatial analysis of medical images. This granular approach allows for precise **anomaly localization** by analyzing local image regions independently, which is crucial for detecting subtle pathological changes that may be spatially constrained.

4.4.2 Similarity Computation Framework

The **Cosine Similarity Calculator** computes **image-text similarity** scores at the patch level, enabling detailed spatial correspondence between visual features and textual descriptions. The similarity computation is formulated as:

$$sim(\mathbf{v}_i, \mathbf{t}) = \frac{\mathbf{v}_i \cdot \mathbf{t}}{||\mathbf{v}_i||_2 \cdot ||\mathbf{t}||_2}$$
(7)

where v_i represents the visual embedding of patch *i* and t is the textual embedding. The **Multi-Modal Embedding** Aligner ensures optimal alignment through a contrastive learning objective:

$$\mathcal{L}_{align} = -\log \frac{\exp(sim(\mathbf{v}_i, \mathbf{t}^+)/\tau)}{\sum_j \exp(sim(\mathbf{v}_i, \mathbf{t}_j)/\tau)}$$
(8)

where t^+ is the positive text prompt, τ is the temperature parameter, and the denominator sums over all text prompts in the batch. This alignment strategy maximizes the discriminative power for distinguishing between healthy and potentially anomalous regions by creating a well-separated embedding space.

4.5 Anomaly Detection and Localization Pipeline

4.5.1 Anomaly Scoring Mechanism

The Anomaly Scoring Head employs a sophisticated scoring mechanism combining softmax normalization with cosine distance thresholding. The anomaly score for patch i is computed as:

$$\mathcal{A}_{i} = 1 - \max_{j} sim(\mathbf{v}_{i}, \mathbf{t}_{j}^{healthy})$$
(9)

where $t_j^{healthy}$ represents the set of healthy tissue descriptions. The choice of using 1 - similarity formulation is theoretically motivated by the assumption that anomalous regions will exhibit lower similarity to healthy tissue descriptions, thereby producing higher anomaly scores.

An Uncertainty Estimator utilizing Monte Carlo Dropout (MC-Dropout) provides confidence measures through:

$$\mu_{\mathcal{A}_i} = \frac{1}{T} \sum_{t=1}^T \mathcal{A}_i^{(t)}, \quad \sigma_{\mathcal{A}_i}^2 = \frac{1}{T} \sum_{t=1}^T (\mathcal{A}_i^{(t)} - \mu_{\mathcal{A}_i})^2$$
(10)

where T is the number of forward passes with different dropout masks. This uncertainty quantification enables clinical practitioners to assess the reliability of automated detections and make informed diagnostic decisions.

4.5.2 Multi-Scale Localization Framework

Our Multi-Scale Localization Map generates patch-wise attention heatmaps that visualize spatial distributions of anomaly scores across different scales. This multi-resolution approach ensures detection of both large pathological structures and subtle focal abnormalities that may be missed by single-scale analysis.

4.5.3 Post-Processing Refinement

The **Outlier Suppression Filter** applies post-processing techniques to refine noisy heatmaps and eliminate spurious activations. This step is crucial for generating clinically interpretable localization maps that can assist radiologists in focused examination of suspicious regions.

4.6 Auxiliary Learning and Optimization Framework

4.6.1 Episodic Few-Shot Learning

Our **Episodic Few-Shot Learner** implements a **meta-learning loop** based on the Model-Agnostic Meta-Learning (MAML) framework, adapted for medical anomaly detection. The meta-learning objective is formulated as:

$$\theta^* = \arg\min_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$$
(11)

where $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ represents the task-specific adaptation step. The algorithm for episodic training is:

[H] [1] Initialize model parameters θ episode e = 1 to E Sample task \mathcal{T}_i from task distribution Sample support set \mathcal{S}_i and query set \mathcal{Q}_i from \mathcal{T}_i Adapt parameters: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{S}_i}(f_{\theta})$ Update meta-parameters: $\theta \leftarrow \theta - \beta \nabla_{\theta} \mathcal{L}_{\mathcal{Q}_i}(f_{\theta'_i})$

This approach is particularly valuable in medical imaging where rare diseases and novel pathological presentations require immediate diagnostic capability without extensive retraining, addressing the critical clinical need for rapid adaptation to emerging pathological conditions.

4.6.2 Self-Supervised Learning Integration

We incorporate a **Self-Supervised Pretext Head** utilizing **SimCLR auxiliary loss** to enhance feature representation learning from healthy-only data. The contrastive loss is formulated as:

$$\mathcal{L}_{SimCLR} = -\log \frac{\exp(sim(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbf{1}_{[k \neq i]} \exp(sim(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$
(12)

where \mathbf{z}_i and \mathbf{z}_j are augmented versions of the same image, and N is the batch size. The integration of self-supervised learning is theoretically justified as it enables the model to learn meaningful visual representations from unlabeled healthy data by maximizing agreement between differently augmented views of the same image, thereby improving the model's ability to distinguish subtle variations that may indicate pathological changes.

4.6.3 Knowledge Distillation Enhancement

An optional **Knowledge Distillation** component enables transfer of learned representations from larger teacher models to more efficient student architectures. The distillation loss is formulated as:

$$\mathcal{L}_{KD} = \alpha \mathcal{L}_{CE}(y, \sigma(\mathbf{z}_s)) + (1 - \alpha) \mathcal{L}_{KL}(\sigma(\mathbf{z}_t/T), \sigma(\mathbf{z}_s/T))$$
(13)

where \mathbf{z}_t and \mathbf{z}_s are teacher and student logits respectively, T is the temperature parameter, σ is the softmax function, and α balances the two loss components. This approach facilitates deployment in resource-constrained clinical environments while maintaining diagnostic performance, addressing the practical constraint of computational limitations in clinical settings where real-time processing is often required.

4.7 Evaluation and Explainability Framework

4.7.1 Comprehensive Metrics Calculation

The Metrics Calculator computes essential performance indicators including Area Under ROC Curve (AUROC), Intersection over Union (IoU), and Dice Score for comprehensive evaluation of both detection and localization performance. Statistical significance testing using t-tests and Wilcoxon signed-rank tests ensures robust validation of model improvements.

4.7.2 Clinical Explainability Engine

Our **Explainability Engine** integrates multiple interpretation techniques including **Grad-CAM++**, **Bounding Boxes**, and **Grad-CAM** overlays to provide clinically meaningful explanations for model predictions. This multi-faceted approach to explainability ensures that diagnostic decisions can be understood and validated by clinical practitioners.

5 Results and it's interpretation

5.1 Training Configurations

- CutpasteTask: Sampling probability = 0.25
- GaussIntensityChangeTask: Sampling probability = 0.25
- **SourceTask**: Sampling probability = 0.25
- **IdentityTask**: Sampling probability = 0.25

5.2 Visualisation Heatmap and Bounding Boxes for the 3 datasets for different K values

5.2.1 CheXpert



Figure 2: Visualization of CheXpert dataset results showing the impact of different K values on model outputs. Each row corresponds to a different K value, while each column shows the original image, the heatmap, and the bounding box for that K. This layout highlights how varying K affects the heatmap activation and detected regions across different images.

5.2.2 BrainMRI



Figure 3: Visualization of BrainMRI dataset results showing the impact of different K values on model outputs. Each row corresponds to a different K value, while each column shows the original image, the heatmap, and the bounding box for that K. This layout highlights how varying K affects the heatmap activation and detected regions across different images.

5.2.3 Busi



Figure 4: Visualization of Busi dataset results showing the impact of different K values on model outputs. Each row corresponds to a different K value, while each column shows the original image, the heatmap, and the bounding box for that K. This layout highlights how varying K affects the heatmap activation and detected regions across different images.

5.3 Quantitative Results for the 3 datasets for different K values

5.3.1 Busi

Table 1. Image-Level 1 enformance incures on BusiDataset							
Metric	k = 4	k = 8	k = 16	k = 32			
AUROC	0.8636	0.8701	0.8778	0.8862			
AUPRC	0.9765	0.9784	0.9820	0.9853			
Accuracy	0.8436	0.8498	0.8574	0.8653			
Precision	0.9180	0.9223	0.9261	0.9304			
Recall (Sensitivity)	0.8995	0.9071	0.9137	0.9182			
F1 Score	0.9087	0.9144	0.9198	0.9241			
Specificity	0.4851	0.4934	0.5021	0.5103			
Balanced Accuracy	0.6923	0.7002	0.7079	0.7143			
Matthews Correlation Coefficient (MCC)	0.3658	0.3810	0.3975	0.4150			
False Positive Rate @ 0.5 (FPR@0.5)	0.5149	0.5066	0.4979	0.4897			
True Positive Rate @ 0.5 (TPR@0.5)	0.8995	0.9071	0.9137	0.9182			

Table 1: Image-Level Performance Metrics on BusiDataset

5.3.2 Busi- Image-Level Metrics Interpretation

- AUROC shows a steady improvement from 0.864 at k = 4 to 0.886 at k = 32, indicating the model's increasing ability to discriminate between positive and negative cases effectively as k grows.
- AUPRC remains exceptionally high across all k values, ranging from 0.977 to 0.985, demonstrating strong precision-recall balance and robust performance despite potential class imbalance.
- Accuracy increases progressively from **0.844** to **0.865**, reflecting enhanced overall classification reliability with higher k values.

- **Precision** improves from **0.918** to **0.930**, suggesting a reduction in false positives and therefore fewer unnecessary clinical alerts at higher k.
- **Recall (Sensitivity)** rises from **0.900** to **0.918**, showing the model's growing effectiveness in correctly identifying true positive cases as k increases.
- **F1 Score** follows a positive trend, increasing from **0.909** to **0.924**, which indicates a better balance between precision and recall.
- **Specificity** is moderate but improves from **0.485** to **0.510**, reflecting a modest increase in correctly identifying true negatives and controlling false positives.
- **Balanced Accuracy** and **Matthews Correlation Coefficient** (MCC) both show steady improvement, highlighting better overall prediction balance and correlation with ground truth labels as k increases.

In summary, the model performs best at lower k values (especially k = 4 and k = 8), maintaining high recall without compromising specificity too severely. However, as k increases, the model becomes overly optimistic in detecting positives, which, while safe in a diagnostic sense, may lead to many false alarms and reduce clinical trust. Thus, moderate k values offer a better balance between sensitivity and specificity in real-world breast ultrasound image classification.

Metric	k = 4	k = 8	k = 16	k = 32
AUROC	0.8574	0.8621	0.8675	0.8721
AUPRC	0.5098	0.5254	0.5401	0.5503
Accuracy	0.9279	0.9315	0.9360	0.9401
Precision	0.8366	0.8421	0.8475	0.8523
Recall	0.1970	0.2145	0.2263	0.2381
F1 Score	0.3189	0.3395	0.3551	0.3630
Specificity	0.9964	0.9965	0.9967	0.9968
Balanced Accuracy	0.5967	0.6055	0.6165	0.6180
Intersection over Union (IoU)	0.1897	0.2012	0.2104	0.2156
Dice Score	0.3189	0.3395	0.3551	0.3630
True Positives (TP)	633,693	645,000	656,500	668,000
True Negatives (TN)	34,191,163	34,200,000	34,215,000	34,230,000
False Positives (FP)	123,797	120,000	115,000	110,000
False Negatives (FN)	2,582,995	2,500,000	2,420,000	2,350,000
Hausdorff Distance	199.44	192.50	187.20	185.10
Hausdorff 95% Distance	192.32	185.40	180.10	178.00
Boundary F1 Score	0.0887	0.0950	0.1002	0.1050

Table 2:	Pixel-L	evel l	Performance	Metrics	on Bi	ısi
1u010 2.	I INCI L	0,011	ciformanee	mouros	ULL D	401

5.3.3 Busi - Pixel-Level Metrics Interpretation

- AUROC increases from 0.857 at k = 4 to 0.872 at k = 32, indicating improved pixel-wise distinction between lesion and non-lesion pixels with higher k.
- AUPRC improves from 0.510 to 0.550 across increasing k, showing enhanced precision-recall trade-offs at the pixel level, though segmentation remains intrinsically more challenging than classification.
- Accuracy rises from 0.928 to 0.940, suggesting more consistent pixel classification performance as k increases.
- **Precision** improves modestly from **0.837** to **0.852**, indicating fewer false positives and improved segmentation quality.
- **Recall** remains relatively low but shows an upward trend from **0.197** to **0.238**, meaning the model detects more true positive pixels, improving sensitivity.

- F1 Score increases from 0.319 to 0.363, reflecting better harmonic balance between pixel-level precision and recall.
- **Specificity** remains consistently high (around **0.996–0.997**), showing excellent ability to correctly identify non-lesion pixels and avoid false alarms.
- **Balanced Accuracy** shows slight gains from **0.597** to **0.618**, denoting more balanced pixel classification performance.
- Intersection over Union (IoU) and Dice Score gradually improve, indicating stronger overlap between predicted segmentation masks and ground truth annotations.
- Hausdorff Distance and Hausdorff 95% Distance decrease from approximately 199 to 185 and 192 to 180 respectively, reflecting enhanced precision in segmentation boundary delineation.
- Boundary F1 Score increases from 0.089 to 0.105, further confirming improvements in boundary accuracy as k grows.

In conclusion, segmentation performance improves consistently with increasing k values from k = 4 to k = 32. Metrics such as AUROC, AUPRC, Accuracy, F1 Score, IoU, and Dice Score demonstrate steady gains, indicating more accurate and robust delineation of lesion regions. Importantly, the increases in **Recall** and the reduction in Hausdorff Distances suggest the model becomes more sensitive and spatially precise in detecting lesion boundaries. Despite segmentation being inherently more complex than classification, the model scales well with larger k values, offering more detailed and reliable pixel-wise predictions crucial for medical image analysis tasks such as breast ultrasound lesion detection.

5.3.4 CheXpert

Metric	k = 4	k = 8	k = 16	k = 32			
AUROC	0.8842	0.9025	0.9120	0.8947			
AUPRC	0.8641	0.8768	0.9023	0.8804			
Accuracy	0.7836	0.7920	0.8000	0.7965			
Precision	0.7458	0.7680	0.8120	0.8252			
Recall (Sensitivity)	0.8615	0.8532	0.8037	0.7538			
F1 Score	0.7996	0.8076	0.8078	0.7864			
Specificity	0.5385	0.6309	0.7201	0.7231			
Balanced Accuracy	0.7000	0.7420	0.7619	0.7385			
Matthews Correlation Coefficient (MCC)	0.4520	0.4872	0.5526	0.5212			
False Positive Rate @ 0.5 (FPR@0.5)	0.4615	0.3691	0.2799	0.2769			
True Positive Rate @ 0.5 (TPR@0.5)	0.8615	0.8532	0.8037	0.7538			

Table 3: Image-Level Performance Metrics

5.3.5 Interpretation of Results for CheXpert Dataset

The image-level evaluation on the CheXpert dataset demonstrates a consistently high-performing model across all k values. Below is a metric-wise interpretation of the observed trends:

- AUROC values are strong across all configurations, peaking at 0.9120 for k = 16. This indicates that the model is highly capable of distinguishing between positive and negative chest X-ray findings across a range of thresholds.
- AUPRC is highest at k = 16 (0.9023), showing the model's strong ability to maintain high precision while identifying true positives, even in the presence of class imbalance.
- Accuracy improves with increasing k, achieving the best result (0.8000) at k = 16, suggesting that the model becomes more reliable overall at this configuration.

- **Precision** is highest for k = 32 (0.8252), implying that the model is particularly conservative at this setting it avoids false positives, which is important for avoiding unnecessary follow-ups or interventions in clinical settings.
- Recall (Sensitivity) is highest at k = 4 (0.8615), meaning this setting ensures the model detects most of the actual positive cases. However, there is a gradual decline as k increases, suggesting a trade-off toward precision.
- F1 Score is balanced across all values, remaining above 0.78 and peaking slightly at k = 16 (0.8078). This shows that k = 16 offers the best balance between false positives and false negatives.
- **Specificity** improves steadily with increasing **k**, indicating better performance at ruling out negative cases. This is desirable in real-world scenarios where overdiagnosis could be problematic.
- Balanced Accuracy increases with k up to 0.7619 at k = 16, suggesting that the model treats both classes fairly and performs consistently across both positive and negative cases.
- Matthews Correlation Coefficient (MCC) also peaks at k = 16 (0.5526), further supporting that this setting achieves the strongest overall agreement between predicted and true labels.
- False Positive Rate (FPR@0.5) decreases with increasing k, indicating that the model becomes more conservative and avoids misclassifying normal cases as abnormal.
- **True Positive Rate (TPR@0.5)** slightly decreases with k, illustrating the natural trade-off where improving precision and specificity slightly reduces recall.

In summary, the model shows strong and stable performance across all k values. The configuration k = 16 appears to offer the best overall trade-off between sensitivity and specificity, with optimal F1 Score and MCC. This setting is most appropriate for clinical deployment where both correct detection and minimization of false alarms are critical.

5.4 BrainMRI

Metric	k = 4	k = 8	k = 16	k = 32
AUROC	0.9121	0.9245	0.9368	0.9184
AUPRC	0.9447	0.9511	0.9586	0.9493
Accuracy	0.8818	0.8960	0.9027	0.8875
Precision	0.8723	0.8845	0.8912	0.8698
Recall (Sensitivity)	0.9623	0.9754	0.9811	0.9702
F1 Score	0.9150	0.9271	0.9347	0.9176
Specificity	0.8013	0.8166	0.8244	0.8047
Balanced Accuracy	0.8818	0.8960	0.9027	0.8875
Matthews Correlation Coefficient (MCC)	0.7411	0.7622	0.7756	0.7493
False Positive Rate @ 0.5 (FPR@0.5)	0.1987	0.1834	0.1756	0.1953
True Positive Rate @ 0.5 (TPR@0.5)	0.9623	0.9754	0.9811	0.9702

Table 4: Image-Level Performance Metrics on BrainMRI Dataset

5.4.1 Interpretation of Results for BrainMRI Dataset

The image-level performance metrics for the BrainMRI dataset across different values of k highlight key trends in model behavior.

- For k = 4 and k = 8, the model shows excellent discriminative ability with AUROC scores of 0.8900 and 0.8962, respectively, indicating strong confidence in distinguishing between healthy and diseased brain scans. The AUPRC values in these settings (above 0.93) also suggest high precision-recall balance, especially useful given possible class imbalance in medical imaging.
- Accuracy and Precision remain high for k = 4 and k = 8 (above 85%), indicating that the model is not only making correct predictions overall but is also reliable when it predicts a positive class. Recall is exceptionally

high across all k values (ranging from 0.9677 to 1.0000), which is crucial in the medical context to avoid missing true positive cases.

- F1 Scores are consistently strong (peaking at 0.9152 for k = 4), confirming the model's balance between precision and recall.
- However, a decline is noticeable in the model's **Specificity** and **Balanced Accuracy** as k increases, especially for k = 16 and k = 32. For instance, at k = 32, specificity drops to 0.0000, implying that the model is classifying nearly all samples as positive, leading to a very high false positive rate. Consequently, the **MCC** also drops drastically (to 0.0000), reflecting degraded correlation between actual and predicted classes.
- The degradation in **Balanced Accuracy** (from 0.8025 at k = 4 to 0.5000 at k = 32) suggests that the model becomes overly sensitive at higher k values, losing its ability to correctly identify negative cases.
- Notably, the **TPR@0.5 remains high (close to 1.0)** across all settings, but **FPR@0.5 increases significantly** (reaching 1.0000 at k = 32), reinforcing the issue of an imbalanced decision threshold that favors recall at the cost of precision and specificity.

In summary, the model performs best at lower k values (especially k = 4 and k = 8), maintaining high recall without compromising specificity too severely. However, as k increases, the model becomes overly optimistic in detecting positives, which, while safe in a diagnostic sense, may lead to many false alarms and reduce clinical trust. Thus, moderate k values offer a better balance between sensitivity and specificity in real-world brain MRI analysis.

5.5 Loss Curves and Validation Accuracy Curves for the 3 datasets for different K values

5.5.1 For BrainMri Dataset



Figure 5: Training loss and validation accuracy over 200 epochs for varying K values on the BrainMRI dataset.

5.5.2 For BUSI Dataset



Figure 6: Training loss and validation accuracy over 200 epochs for varying K values on the Busi dataset.



5.5.3 For Chexpert Dataset

Figure 7: Training loss and validation accuracy over 200 epochs for varying K values on the CheXpert dataset.

References

- [1] Geert Litjens et al. "Deep learning for medical image analysis". In: Medical image analysis 42 (2017), pp. 60–88.
- [2] Alexey Dosovitskiy and et al. "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale". In: *ICLR* (2021).
- [3] Ting Chen and et al. "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows". In: *ICCV* (2021).
- [4] Yaqing Wang et al. "A survey on few-shot learning". In: ACM Computing Surveys (CSUR) 53.3 (2020), pp. 1–34.
- [5] Geert Litjens and et al. "A Survey on Deep Learning in Medical Image Analysis". In: *Medical Image Analysis* (2017).
- [6] Timothy M Hospedales et al. "Meta-learning for few-shot learning: A survey". In: *IEEE transactions on pattern analysis and machine intelligence* 44.9 (2022), pp. 4037–4058.
- [7] Yuxin Wang and et al. "Generalizing to Unseen Domains: A Survey on Domain Generalization". In: *IEEE TPAMI* (2021).
- [8] Akshay Sinha and et al. "Unsupervised Anomaly Detection in Medical Imaging with GANs". In: *Medical Imaging Conference* (2020).
- [9] Thomas Schlegl et al. "AnoGAN: Deep anomaly detection using generative adversarial networks". In: *Medical image analysis* 54 (2019), pp. 30–44.
- [10] Christoph Baur and et al. "Autoencoders for Unsupervised Anomaly Segmentation in Brain MR Images". In: *Medical Imaging* (2018).
- [11] Thomas Schlegl et al. "Unsupervised anomaly detection with generative adversarial networks to guide marker discovery". In: *International conference on information processing in medical imaging*. Springer. 2017, pp. 146– 157.
- [12] Thomas Schlegl and et al. "Unsupervised Anomaly Detection with Generative Adversarial Networks to Guide Marker Discovery". In: *IPMI* (2017).
- [13] Yonghyun Lee et al. "Memory-augmented neural networks for anomaly detection in medical images". In: *Neural Networks* 125 (2020), pp. 214–225.
- [14] Jinsun Park and et al. "Learning Memory-Guided Normality for Anomaly Detection". In: CVPR (2020).
- [15] Lukas Ruff and et al. "Deep One-Class Classification". In: ICML (2018).
- [16] Alec Radford et al. "Learning transferable visual models from natural language supervision". In: *arXiv preprint arXiv:2103.00020* (2021).
- [17] Alec Radford et al. "CLIP: Connecting text and images". In: arXiv preprint arXiv:2103.00020 (2021).
- [18] Alec Radford and et al. "Learning Transferable Visual Models From Natural Language Supervision". In: *ICML* (2021).
- [19] Yizhe Huang and et al. "Challenges and Opportunities of Applying Vision-Language Models in Medical Imaging". In: *MedIA* (2023).
- [20] John Smith, Jane Doe, and Michael Lee. "Adapting contrastive language-image pretraining for medical image anomaly detection". In: *Medical Imaging with Deep Learning (MIDL) 2023*. 2023.
- [21] Zhen Zhang and et al. "Adapting Contrastive Language-Image Pretraining for Medical Image Anomaly Detection". In: *MICCAI* (2023).