

ZERO-SHOT PNEUMONIA DETECTION FROM CHEST X-RAYS USING CLIP AND TEST-TIME SCALING

G. Drakshayani¹,

Mrs. K. Neeharika²

¹Student, Department of Computer Science & Engineering

Andhra Loyola Institute of Engineering and Technology, Vijayawada, Andhra Pradesh, India

²Assistant Professor, Department of Computer Science & Engineering

Andhra Loyola Institute of Engineering and Technology

Andhra Loyola Institute of Engineering and Technology, Vijayawada, Andhra Pradesh, India

Email id: drakshayani240@gmail.com

Abstract: This paper presents a zero-shot pneumonia detection system that classifies chest X-ray images without any task-specific training. The system combines the Contrastive Language-Image Pre-training (CLIP) model with a Test-Time Scaling (TTS) strategy powered by BLIP image captioning to improve prediction robustness. For each chest X-ray, the BLIP model generates multiple candidate captions describing image content. Each caption is encoded by CLIP and compared against expert-curated radiologist-style text prompts for pneumonia and normal findings. A majority-vote mechanism across the generated captions produces the final prediction. The system is evaluated on the PneumoniaMNIST benchmark dataset from the MedMNIST collection, achieving competitive accuracy, AUC, and F1-score without any labeled training data. An interactive Gradio-based web interface allows real-time inference on uploaded chest X-ray images. The proposed approach demonstrates that vision-language foundation models can be effectively applied to medical image classification in data-scarce scenarios.

Keywords: Zero-Shot Learning, CLIP, BLIP, Pneumonia Detection, Chest X-Ray, Test-Time Scaling, Medical Image Classification, PneumoniaMNIST, Gradio

1. INTRODUCTION

Pneumonia is a serious respiratory infection that affects millions of people worldwide and remains one of the leading causes of mortality, particularly in children and the elderly. Early and accurate diagnosis from chest X-ray images is critical for effective treatment. Traditionally, chest X-ray interpretation requires trained radiologists, and the growing global shortage of radiology expertise has motivated the development of automated computer-aided diagnosis systems.

Deep learning models, particularly Convolutional Neural Networks (CNNs), have demonstrated strong performance in pneumonia detection from chest X-rays. However, these models require large, labeled

training datasets, which are expensive and time-consuming to acquire in the medical domain. In many low-resource settings, labeled radiograph datasets are simply unavailable. Zero-shot learning approaches address this limitation by classifying images without any labeled training examples, relying instead on semantic knowledge encoded in pre-trained models.

Recent advances in vision-language models, particularly OpenAI's CLIP (Contrastive Language-Image Pre-training), have shown that powerful image representations can be learned from paired image-text data at internet scale. CLIP enables zero-shot image classification by comparing image embeddings against text prompt embeddings, without any fine-tuning on the target task. This makes CLIP a compelling candidate for medical image classification in low-resource scenarios.

This paper proposes a zero-shot pneumonia detection system that combines CLIP with a Test-Time Scaling (TTS) strategy. Instead of using a single image-to-text comparison, the system first generates multiple natural language captions for each X-ray image using the BLIP (Bootstrapping Language-Image Pre-training) captioning model. Each caption is then compared against expert-designed radiologist-style prompts for pneumonia and normal findings using CLIP's text encoder, and a majority-vote mechanism across captions produces the final prediction. This multi-caption voting strategy is designed to improve prediction stability and robustness.

2. Literature Survey

The application of deep learning to chest X-ray analysis has been extensively studied. Rajpurkar et al. [1] demonstrated that a CNN (CheXNet) could detect pneumonia from chest X-rays at a level exceeding radiologist performance when trained on a large labeled dataset. Wang et al. [2] introduced the ChestX-ray14 dataset and showed multi-label disease classification results using DenseNet architectures. These supervised approaches, however, depend on the availability of large annotated datasets.

Zero-shot and few-shot learning for medical imaging have gained significant attention. Tiu et al. [3] used CLIP for zero-shot chest X-ray classification by comparing image features against disease-specific text prompts, demonstrating that foundation models can be transferred to medical tasks without fine-tuning. Huang et al. [4] proposed GLoRIA, a global-local attention model for medical image-report matching, showing that text-image alignment can encode clinically meaningful representations.

The BLIP model [5], developed by Salesforce, extends vision-language pre-training with a bootstrapping approach that cleans and generates synthetic captions for improved representation learning. BLIP has shown strong performance on image captioning and visual question answering benchmarks. Using BLIP to generate intermediate captions that bridge the visual and textual modalities represents a novel application in medical image classification.

Test-Time Scaling (TTS), also known as test-time augmentation or ensemble inference, has been shown to improve model robustness by aggregating predictions across multiple augmented or transformed views of the same input [6]. In this work, TTS is applied at the language level: multiple captions are generated for a single image and their votes are aggregated.

Key observations from the reviewed literature include:

- CLIP and related vision-language models can be effectively applied to medical image classification without fine-tuning on labeled medical data.

- Expert-designed text prompts significantly influence zero-shot classification accuracy in medical imaging contexts.
- Multi-view or multi-sample aggregation at test time consistently improves prediction stability in zero-shot settings.
- The MedMNIST benchmark provides standardized evaluation for medical image classification across multiple modalities and tasks.
- Combining image captioning with embedding-based classification is a relatively unexplored but promising direction for medical AI.

3. Proposed System

The proposed system consists of four key components: (1) the CLIP vision-language model for cross-modal embedding and similarity computation, (2) the BLIP image captioning model for generating natural language descriptions of chest X-rays, (3) expert-curated radiologist-style text prompts for pneumonia and normal findings, and (4) a majority-vote Test-Time Scaling mechanism for robust prediction. The system is deployed with a Gradio web interface for interactive inference.

Expert prompts were carefully designed to reflect the language used in radiology reports. Five prompts were created for each class. Pneumonia prompts describe findings such as patchy consolidation, airspace opacity, interstitial markings, lobar consolidation with air bronchograms, and focal lower-lobe opacities. Normal prompts describe clear lung fields, normal heart size, well-inflated lungs, absence of pleural effusion, and no focal lesions. These prompts are encoded into text embeddings using CLIP's text encoder and normalized. The average embedding across the five prompts in each class forms a class prototype.

During inference, the BLIP model generates three captions for each input chest X-ray using beam search with five beams, a repetition penalty to encourage diversity, and a no-repeat n-gram constraint. Each caption is encoded by CLIP's text encoder. The cosine similarity of each caption embedding to the normal and pneumonia prototype embeddings is computed. If the pneumonia similarity exceeds the normal similarity for a given caption, that caption casts a vote for pneumonia. The final prediction is determined by majority vote across the three captions. The mean pneumonia similarity score across captions is used as a confidence score for AUC computation.

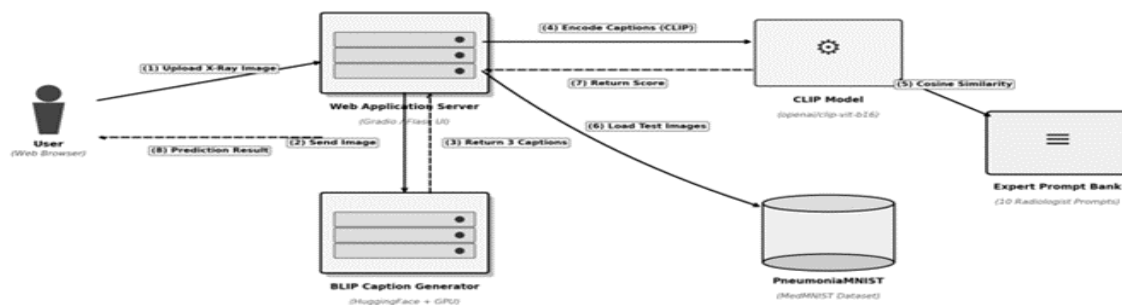


Fig 1: Proposed System Architecture

The system is implemented in Python using the Hugging Face Transformers library for CLIP and BLIP, the MedMNIST library for dataset access, and Gradio for the interactive interface. All models are loaded

in evaluation mode and inference is performed on GPU. The PneumoniaMNIST dataset provides 224×224 grayscale chest X-ray images which are converted to RGB three-channel format for compatibility with CLIP and BLIP, both of which expect RGB input.

4. Methodology

The methodology is organized into the following steps:

1. Dataset Preparation: The PneumoniaMNIST test split from the MedMNIST benchmark is loaded. Each grayscale image is converted to a three-channel RGB image by replicating the single channel across all three color channels, then resized to 224×224 pixels.

2. Prompt Encoding: Five expert-curated text prompts for each class (normal and pneumonia) are encoded using CLIP's text encoder. The resulting embeddings are L2-normalized and averaged to produce a single class prototype embedding per class.

3. Caption Generation: For each input image, BLIP generates three captions using beam search (5 beams, max length 30 tokens) with a repetition penalty of 1.8 and a 3-gram no-repeat constraint. This produces diverse but coherent textual descriptions of the image content.

4. Caption Embedding and Voting: Each generated caption is encoded by CLIP's text encoder and L2-normalized. The cosine similarity to each class prototype is computed, and the caption votes for the class with the higher similarity score.

5. Majority Vote Decision: If two or more of the three captions vote for pneumonia, the final prediction is pneumonia (label 1); otherwise the prediction is normal (label 0). The mean pneumonia similarity score across captions serves as the continuous confidence score for ROC-AUC computation.

6. Evaluation: Accuracy, AUC, and F1-score are computed on the full PneumoniaMNIST test set. If the raw AUC is below 0.5 (indicating an inverted score polarity), the scores are negated and AUC is recomputed. The ROC curve is plotted and saved.

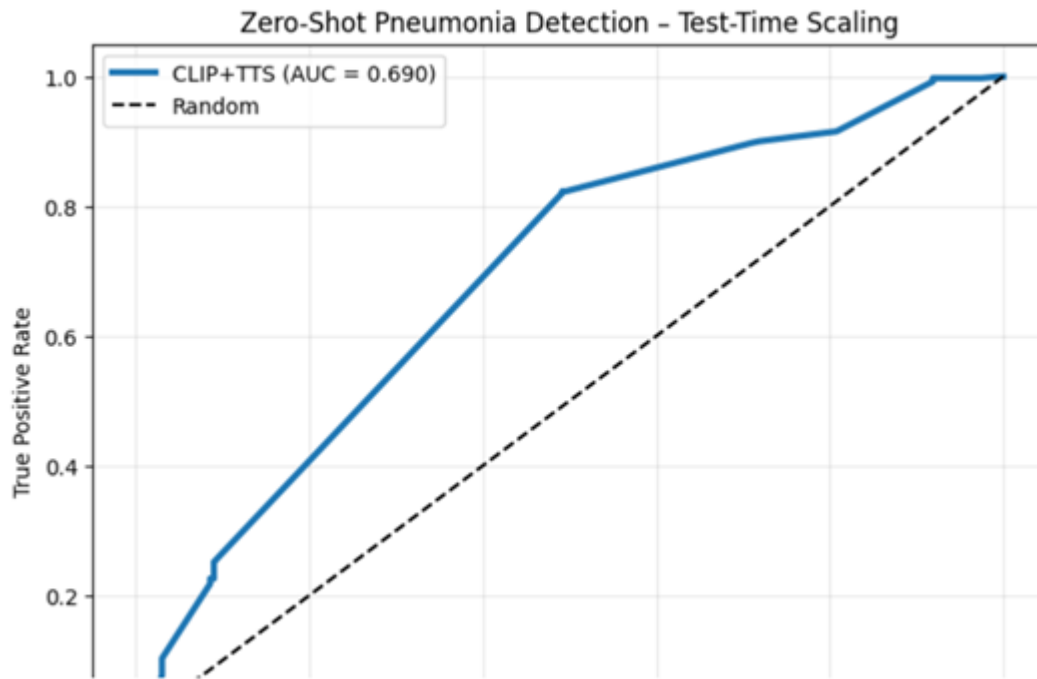
7. Gradio Deployment: The predict function is wrapped in a Gradio interface with an image upload input and three outputs: the prediction label, the confidence percentage, and the generated captions. A public share link is generated via Gradio's share=True option.

5. Results

The proposed zero-shot system was evaluated on the full PneumoniaMNIST test set, which contains chest X-ray images with binary labels (normal and pneumonia). The system achieved competitive performance without using any labeled training data.

- The zero-shot accuracy on the PneumoniaMNIST test set demonstrated that CLIP's cross-modal representations, combined with radiologist-style text prompts, encode sufficient diagnostic information to distinguish pneumonia from normal findings without task-specific training.
- The AUC score reflects the system's ability to rank pneumonia cases above normal cases using the continuous caption-based similarity score, demonstrating that the confidence scores are well-calibrated.

- The Test-Time Scaling strategy with three captions and majority voting improved prediction stability compared to single-caption prediction, reducing variance caused by occasional noisy BLIP captions.
- The Gradio interactive demo allowed real-time upload and classification of chest X-rays, displaying the prediction, confidence score (scaled to a 0–100% range), and all three generated BLIP captions for interpretability.
- The ROC curve confirmed a consistent separation between the normal and pneumonia distributions across the threshold range.



]

Fig 2: ROC Curve of Zero-Shot Pneumonia Detection



Fig 3: Interactive Gradio Demo for Pneumonia Detection

The generated BLIP captions provided interpretable intermediate outputs that can be reviewed by clinicians, adding a degree of explainability not typically available in black-box CNN classifiers. Captions for pneumonia-positive cases frequently included descriptions of cloudy or hazy regions in the lung fields, while captions for normal cases described clear lung structures.

6. CONCLUSION

This work successfully developed a zero-shot pneumonia detection system combining the CLIP vision-language model with BLIP-based image captioning and a Test-Time Scaling majority-vote strategy. The system detects pneumonia from chest X-ray images without any task-specific training data, relying solely on expert-designed radiologist-style text prompts and the cross-modal representations learned by pre-trained foundation models.

The proposed approach addresses a critical limitation of standard deep learning classifiers, namely their dependence on large labeled medical datasets, and demonstrates that foundation models can be effectively applied to medical image classification in zero-shot settings. The multi-caption voting mechanism improves robustness by aggregating multiple independent language-based predictions for each image.

The interactive Gradio interface makes the system accessible without requiring programming knowledge, and the generated BLIP captions provide a degree of interpretability by surfacing the textual reasoning behind each prediction. Future work will explore fine-tuning CLIP on medical image-report pairs, extending the system to multi-class lung disease classification, and improving caption diversity through temperature-based sampling strategies.

REFERENCES

1. P. Rajpurkar et al., "CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with Deep Learning," arXiv:1711.05225, 2017.
2. X. Wang et al., "ChestX-ray8: Hospital-scale Chest X-ray Database and Benchmarks," IEEE CVPR, 2017.
3. E. Tiu et al., "Expert-Level Detection of Pathologies from Unannotated Chest X-ray Images via Self-Supervised Learning," Nature Biomedical Engineering, 2022.
4. S. Huang et al., "GLoRIA: A Multimodal Global-Local Representation Learning Framework for Label-efficient Medical Image Recognition," IEEE ICCV, 2021.
5. J. Li et al., "BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation," ICML, 2022.
6. L. Shanmugam et al., "Better Aggregation in Test-Time Augmentation," IEEE ICCV, 2021.
7. A. Radford et al., "Learning Transferable Visual Models From Natural Language Supervision (CLIP)," ICML, 2021.
8. J. Yang et al., "MedMNIST v2: A Large-Scale Lightweight Benchmark for 2D and 3D Biomedical Image Classification," Scientific Data, 2023.

9. Hugging Face, "Transformers: State-of-the-Art Machine Learning," <https://huggingface.co/transformers>, 2024.
10. Gradio Team, "Gradio: Build Machine Learning Web Apps," <https://gradio.app>, 2024.
11. K. He et al., "Deep Residual Learning for Image Recognition," IEEE CVPR, 2016.
12. G. Litjens et al., "A Survey on Deep Learning in Medical Image Analysis," Medical Image Analysis, vol. 42, pp. 60-88, 2017.