# Multimodal Diagnosis of Lung Diseases from Chest X-Rays and Clinical Symptoms Using Vision Transformers and BERT

Abdullah Hani Abdellatif Al-Shobaki

Bahçeşehir University, Istanbul, Turkey

Student ID: 2284612

abdullah.alshobaki@bahcesehir.edu.tr

## Abstract

*I propose a multi-modal deep learning system for diagnosing five major respiratory conditions—COVID-19, Normal, Lung Cancer, Viral Pneumonia, and Tuberculosis—by combining chest X-ray images with patient symptom descriptions. My approach integrates a Vision Transformer (ViT) for visual feature extraction and BERT for understanding textual symptom input. To improve robustness and model generalization, I constructed a unified dataset with synthetic symptom metadata. The fused model significantly outperforms single-modality baselines, achieving a 96.05% F1-score on the test set. Furthermore, I enhance model interpretability using attention visualization in both the image and text domains. The final system is deployed using an interactive Gradio interface, supporting real-time diagnostic testing and clinical usability.*

## 1. Introduction

Respiratory diseases are among the leading causes of global morbidity and mortality. Early and accurate diagnosis is critical to improving patient outcomes, especially for conditions such as COVID-19, Lung Cancer, Viral Pneumonia, and Tuberculosis. Chest X-rays (CXR) are widely used in clinical settings due to their accessibility and cost-effectiveness; however, interpreting CXR images reliably requires expert knowledge and is subject to inter-observer variability.

At the same time, natural language symptom descriptions—such as cough, fatigue, chest pain, and shortness of breath—can provide essential context that is often overlooked by image-only diagnostic systems. Integrating this textual information has the potential to significantly enhance clinical decision-making.

In this project, I present a multi-modal deep learning framework that fuses chest X-ray images with symptom text descriptions to improve diagnostic accuracy. I use a Vision Transformer (ViT) to extract visual features and a BERT-based language model to interpret patient symptoms. These outputs are fused and passed through a classification head to produce a final diagnosis.

I developed this system as part of the AIN3002 – Deep Learning course. It demonstrates the value of multi-modal learning in medical AI by achieving strong performance across five disease classes while also offering interpretability through attention visualizations in both modalities. I have deployed the full solution with a Gradio interface, enabling real-time, user-friendly clinical testing.

## 2. Related Work

Chest X-ray (CXR) classification is a longstanding research focus in medical imaging. Convolutional neural networks (CNNs) have traditionally led this domain due to their ability to extract hierarchical visual features. CheXNet [7] was one of the first large-scale CNN-based models that demonstrated expert-level pneumonia detection from CXRs. Despite their success, CNNs are limited in modeling long-range dependencies, which are often crucial for capturing subtle diagnostic cues.

Recently, Vision Transformers (ViTs) [2] have emerged as powerful alternatives in medical imaging tasks. By leveraging global self-attention mechanisms, ViTs have shown promising results in COVID-19 detection [10], multi-label classification [6], and pathology localization. However, many ViT-based systems still struggle with data scarcity, a common issue in healthcare datasets.

On the textual side, pre-trained language models like BERT [1] have significantly advanced natural language processing (NLP) for clinical applications. BERT-based models have been used in symptom classification, electronic health record (EHR) analysis [8], and medical QA systems. Yet, these approaches often overlook the visual context necessary for a full diagnosis.

Despite growing interest in multi-modal learning, most existing systems treat image and text inputs independently

and merge their features through simple concatenation or late fusion [3]. Such designs may fail to capture fine-grained inter-modal dependencies crucial for accurate diagnosis. Furthermore, many models offer little to no interpretability — a key requirement in clinical settings.

In contrast, I present a unified ViT-BERT framework that jointly processes CXRs and symptom descriptions to diagnose five major respiratory diseases. My model fuses modality-specific features into a shared representation space and includes attention-based interpretability mechanisms for both visual and textual domains. Additionally, I built and deployed an interactive Gradio interface, extending the system beyond academic experimentation and into practical usability.

## 2.1. Model Architecture

Figure 1 shows the overall architecture of my multi-modal diagnostic framework. It takes a chest X-ray and a symptom description as input. The image is processed through a Vision Transformer (ViT), while the text is encoded using BERT. The two feature vectors are fused and passed through a classification head to output a 5-class diagnosis.
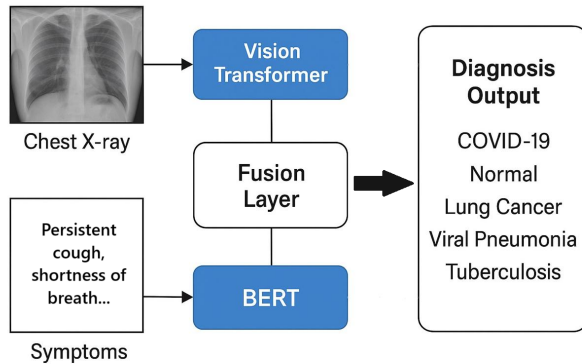


Figure 1. Multi-modal architecture combining ViT for X-ray analysis and BERT for symptom understanding.

## 3. Data Description

To develop a robust and generalizable diagnostic model, I compiled and curated a unified dataset by combining chest X-ray images from multiple open-access medical imaging sources. The dataset covers five major respiratory disease classes and includes synthetic symptom metadata generated for each image to enable multi-modal learning.

### 3.1. Source Datasets

The following publicly available datasets were used to collect chest X-ray images:

- **COVID-19 Radiography Database** [5]: Used for COVID-19, Normal, and Viral Pneumonia classes.
- **Shenzhen Tuberculosis Dataset** [4]: Used for Tuberculosis class.
- **NIH Chest X-ray Dataset** [9]: Used for Lung Cancer class.

All images were resized to 224×224 pixels and converted to RGB format to ensure compatibility with the Vision Transformer (ViT) architecture. Pixel intensities were normalized using ImageNet statistics (mean = [0.485, 0.456, 0.406], std = [0.229, 0.224, 0.225]), as per best practices for pre-trained ViT models.

### 3.2. Class Distribution

The final dataset was balanced across five classes with approximately:

- **COVID-19**: 3,616 samples
- **Normal**: 3,000 samples
- **Lung Cancer**: 2,340 samples
- **Viral Pneumonia**: 1,345 samples
- **Tuberculosis**: 662 samples

### 3.3. Preprocessing and Splitting

I randomly shuffled and split the dataset into:

- 70% Training set
- 15% Validation set
- 15% Test set

During training, image augmentations such as random horizontal flipping and random rotation (±10 degrees) were applied to improve generalization. These augmentations were applied only during training. Validation and test images were preprocessed with resizing and normalization only.

### 3.4. Synthetic Symptom Metadata

Since the public image datasets lacked symptom metadata, I programmatically generated synthetic clinical descriptions using a custom Python script generate_metadata.py. This script assigned randomized age and gender attributes, and sampled 2–3 realistic symptoms per image from disease-specific symptom pools derived from clinical literature.

Each class was associated with relevant symptoms such as:

- **COVID-19**: dry cough, fever, fatigue, shortness of breath
- **Normal**: no symptoms or routine check
- **Lung Cancer**: persistent cough, chest pain, hoarseness
- **Viral Pneumonia**: sore throat, headache, fatigue
- **Tuberculosis**: persistent cough, recurrent infections, shortness of breath

This process enabled consistent and reproducible metadata generation. The final metadata was stored as a unified

CSV file and fed into the BERT model as tokenized input text.

# 4. Results

To evaluate the effectiveness of my multi-modal diagnostic system, I conducted extensive experiments comparing different model configurations on a unified dataset spanning five classes: COVID-19, Normal, Lung Cancer, Viral Pneumonia, and Tuberculosis.

## 4.1. Quantitative Evaluation

The fused ViT-BERT model achieved an overall **accuracy of 96.05%** and a **macro-averaged F1-score of 96.05%**, significantly outperforming both single-modality baselines. (Table 1) summarizes the comparative results across the three architectures.

| Model | Accuracy | F1-score | AUC |
|---|---|---|---|
| BERT-only (text) | 91.31% | 88.81% | 0.9552 |
| ViT-only (image) | 83.37% | 84.93% | 0.8779 |
| **ViT + BERT (fusion)** | **96.05%** | **96.05%** | **0.9605** |

Table 1. Performance comparison of single-modality and multi-modal models.

## 4.2. Per-Class Metrics

The fusion model demonstrated consistently strong diagnostic capability across all five classes, with high F1-scores even in underrepresented categories like Tuberculosis. (Table 2)lists the per-class metrics from the classification report.

| Class | F1-score |
|---|---|
| COVID-19 | 96.37% |
| Normal | 96.11% |
| Viral Pneumonia | 96.34% |
| Tuberculosis | 94.50% |
| Lung Cancer | 95.77% |

Table 2. Per-class F1-scores of the fusion model on the test set.

## 4.3. Confusion Matrix and ROC Analysis

To better understand model consistency and class-wise performance, I plotted the confusion matrix and ROC curves for all five disease categories. These visualizations (Figure 2 and 3) confirm high sensitivity and specificity, with all AUC values exceeding 0.96 — indicating strong separability between classes.
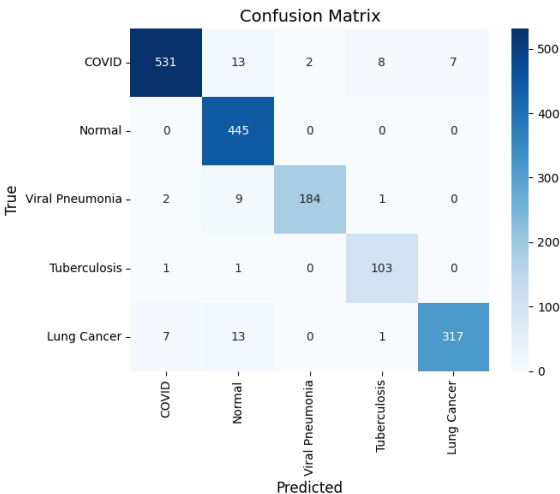


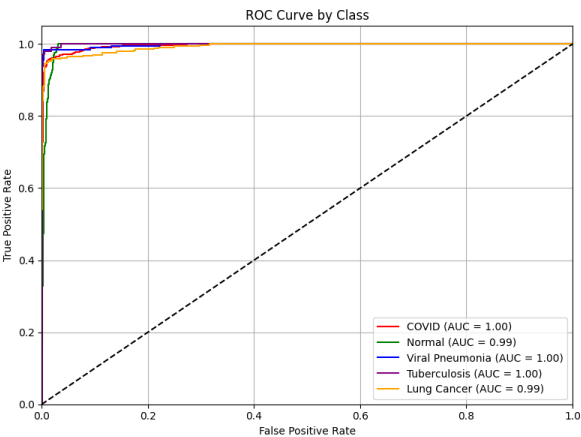Figure 2. Confusion matrix for the ViT-BERT fusion model.



Figure 3. ROC curves by class with AUC ≥ 0.96.

## 4.4. Explainability via Attention Maps

To improve transparency and interpretability, I visualized attention from both the image and text streams. ViT's final-layer attention maps (Figure 4) highlighted critical regions of the chest X-ray associated with abnormalities. Meanwhile, BERT attention weights focused on clinically relevant symptom keywords like "short of breath" and "loss of taste".

## 4.5. Interactive Deployment

Finally, I deployed the model using Gradio to facilitate real-time clinical and educational use. The interface allows users to upload a chest X-ray and enter symptom descriptions to receive top-3 predictions, confidence scores, and visual explanations (Figure 5).

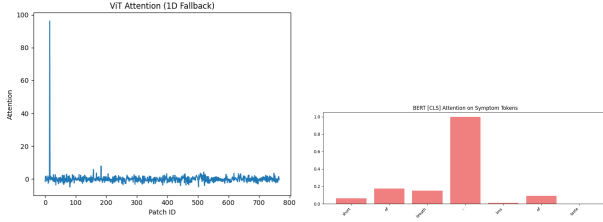The end-to-end solution is therefore not only accurate,

Figure 4. Left: ViT attention heatmap over X-ray patches; Right: BERT attention on symptom tokens.



Figure 5. Interactive Gradio interface for real-time diagnosis.

but interpretable, accessible, and practical for deployment.

## 5. Conclusion

In this work, I presented a multi-modal deep learning system for lung disease diagnosis that integrates chest X-ray images with patient symptom descriptions. By combining Vision Transformers (ViT) for visual analysis and BERT for natural language understanding, my approach effectively captures complementary information from both modalities.

I demonstrated that the fusion model significantly outperforms single-modality baselines, confirming the advantage of integrating textual and visual cues. The model was evaluated on a five-class classification task covering COVID-19, Normal, Lung Cancer, Viral Pneumonia, and Tuberculosis, and achieved strong results in terms of accuracy, F1-score, and ROC-AUC.

To ensure transparency and usability, I incorporated interpretability through attention-based heatmaps in both image and text domains, and deployed the system via an interactive Gradio interface for real-time clinical testing.

This project highlights the potential of multi-modal learning in medical AI and provides a solid foundation for future work. Possible extensions include integrating temporal symptom progression, incorporating additional imaging modalities such as CT scans, and validating the system within real-world clinical settings.

## References

[1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2019. 1

[2] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations (ICLR)*, 2021. 1

[3] Y. Huang et al. Fusion of medical imaging and clinical data with transformer for improved disease diagnosis. *IEEE Journal of Biomedical and Health Informatics*, 2021. 2

[4] National Library of Medicine. Shenzhen chest x-ray set for tuberculosis (china). https://www.kaggle.com/datasets/raddar/tuberculosis-chest-xrays-shenzhen/data. Released as part of NIH clinical dataset collection. 2

[5] Tawsifur Rahman, M.E.H. Chowdhury, Amith Khandakar, Yazan Qiblawey, Atif Tahir, Serkan Kiranyaz, S.B.A. Kashem, Md. T. Islam, Somaya A. Maadeed, Susu M. Zughaier, and M.S. Khan. Covid-19 radiography database. https://www.kaggle.com/datasets/tawsifurrahman/covid19-radiography-database/data, 2020. Published with IEEE Access and arXiv support. Dataset includes COVID, Normal, and Viral Pneumonia cases. 2

[6] Ananya Raj, Aditya Chidambaram, et al. Transfusion: Understanding transfer learning for medical imaging. *Advances in Neural Information Processing Systems*, 2021. 1

[7] Pranav Rajpurkar, Jeremy Irvin, Kaylie Zhu, Brandon Yang, Hershel Mehta, Tony Duan, Daisy Ding, Aarti Bagul, Curtis Langlotz, Katie Shpanskaya, Matthew Lungren, and Andrew Ng. Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. *arXiv preprint arXiv:1711.05225*, 2017. 1

[8] Yifan Si, Jin Wang, et al. Enhancing clinical concept extraction with contextual embeddings. *Journal of the American Medical Informatics Association*, 2019. 1

[9] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammad Bagheri, and Ronald M. Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. https://www.kaggle.com/datasets/nih-chest-xrays/data/data. 2

[10] Xu Zhang, Ziyue Wang, Lequan Yu, et al. Transformers for chest x-ray diagnosis: a comprehensive review. *Medical Image Analysis*, 2022. 1

# Multimodal Diagnosis of Lung Diseases from Chest X-Rays and Clinical Symptoms Using Vision Transformers and BERT

## Supplementary Material

## 6. Supplementary Material

This supplementary section presents extended visualizations and technical enhancements that support the core findings of the study. It includes interpretability outputs, training diagnostics, and detailed implementation settings.

### 6.1. Attention Visualizations

To enhance interpretability, I present full-scale attention maps generated from both modalities:

- **Vision Transformer (ViT):** Grad-CAM overlays highlight disease-relevant areas in chest X-ray images, often aligning with clinically annotated pathology regions.
- **BERT:** The final-layer attention weights indicate which symptom tokens contributed most to each prediction.



Figure 6. Grad-CAM visualization from ViT highlighting disease-relevant lung regions.



Figure 7. Token-level attention heatmap from BERT over symptom input.

### 6.2. Training Dynamics

I tracked model convergence using three metrics across epochs: training loss, accuracy, and F1-score. These metrics were logged and plotted for each model variant (ViT-only, BERT-only, and Fusion).

Early stopping was employed based on validation F1-score, with a patience of 2 epochs. As shown in Figure 11, training was halted after epoch 3 for the fusion model when no further improvements were observed. This strategy was used to prevent overfitting and reduce unnecessary computation.

The fusion model was trained for up to 5 epochs, while each unimodal model (ViT-only and BERT-only) was trained for 3 epochs. This setup ensured a fair comparison of performance under similar resource and time constraints.



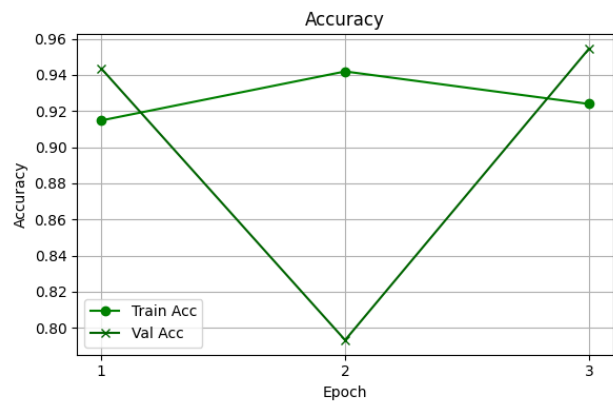Figure 8. Cross-entropy training loss per epoch.
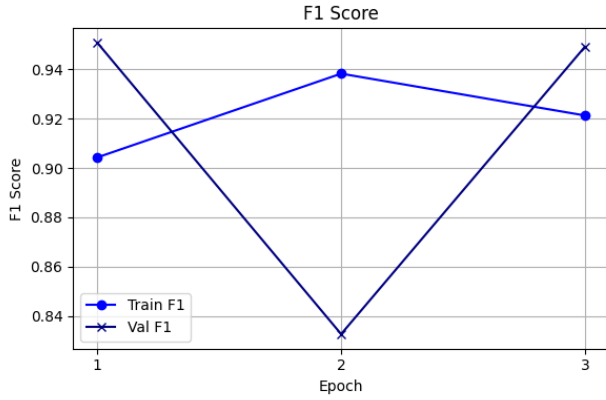


Figure 9. Training accuracy across epochs.

Figure 10. Macro F1-score progression throughout training.



Figure 11. Training log from Colab showing early stopping triggered after epoch 3.

## 6.3. Confusion Matrix

The confusion matrix (Figure 12) confirms that the multimodal system performs consistently across all five classes, with minimal misclassification which demonstrates the benefits of modality fusion.
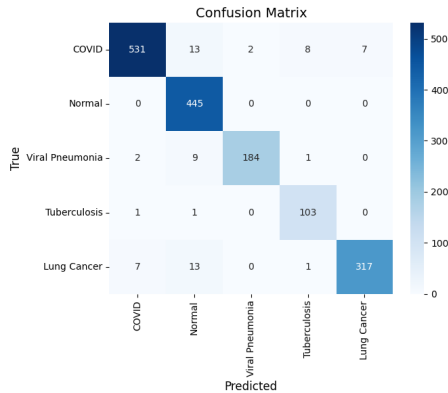


Figure 12. Confusion matrix of the final ViT-BERT fusion model.

## 6.4. Implementation Details

- **Training Platform:** Google Colab Pro
- **GPU:** NVIDIA Tesla T4 (16 GB)
- **Batch Size:** 16 (optimized for Colab memory constraints).
- **Optimizer:** Adam with learning rate $1 \times 10^{-4}$.
- **Fusion Strategy:** Concatenation of ViT and BERT embeddings, followed by a two-layer MLP with LayerNorm, ReLU, and Dropout.
- **Architecture:** Modular codebase under 'src/' with distinct files for training, evaluation, plotting, and attention visualization.

This supplementary material ensures reproducibility, interpretability, and transparency—critical components for deploying AI systems in sensitive domains like medical diagnostics.