

Abstract

I propose a deep learning pipeline that fuses chest X-ray images and patient symptom data to diagnose five respiratory diseases. Utilizing Vision Transformers (ViT) for image analysis and BERT for textual data, the model achieves a 96.05% F1-score. The system includes a real-time Gradio interface and provides visual explanations through attention heatmaps.

Introduction

Respiratory diseases such as COVID-19, Lung Cancer, Viral Pneumonia, and Tuberculosis are leading causes of mortality worldwide. Early and accurate diagnosis is crucial. Combining imaging data with patient-reported symptoms can enhance diagnostic accuracy, but integrating these modalities poses challenges.

Motivation & Challenges

- Limitations of X-rays:** Sole reliance on imaging may miss contextual clinical information.
- Integration Difficulty:** Combining heterogeneous data types (images and text) is non-trivial.
- Need for Explainability:** Clinicians require transparent AI decisions for trust and adoption.

Key Contributions

- Developed a ViT + BERT fusion model for multimodal diagnosis.
- Synthesized a multimodal dataset with realistic symptom metadata.
- Achieved a 96.05% F1-score across five classes.
- Created a Gradio-based real-time user interface.
- Implemented attention heatmaps for model explainability.

Methodology

Model Architecture:

- Image Encoder:** ViT (Base 16x224)
- Text Encoder:** BERT (Base Uncased)
- Fusion:** Concatenation followed by a 2-layer MLP with LayerNorm, ReLU, and Dropout

Data Pipeline:

- Sources:** COVID-19 Radiography Database, NIH Chest X-rays, Shenzhen TB Dataset
- Preprocessing:** Resized images to 224x224, normalized using ImageNet statistics
- Metadata Synthesis:** Generated age, gender, and symptom data using `generate_metadata.py`

Training Setup:

- Loss Function:** CrossEntropyLoss
- Optimizer:** Adam with a learning rate of 1e-4
- Batch Size:** 16
- Early Stopping:** Implemented to prevent overfitting

Model Comparison

Model	Accuracy	F1-Score
ViT-only (image only)	83.37%	84.93%
BERT-only (text only)	91.31%	88.81%
Fusion (ViT + BERT)	96.0%	96.0%

Fusion significantly outperforms single-modality models, validating the power of multimodal learning.

Experimental Results

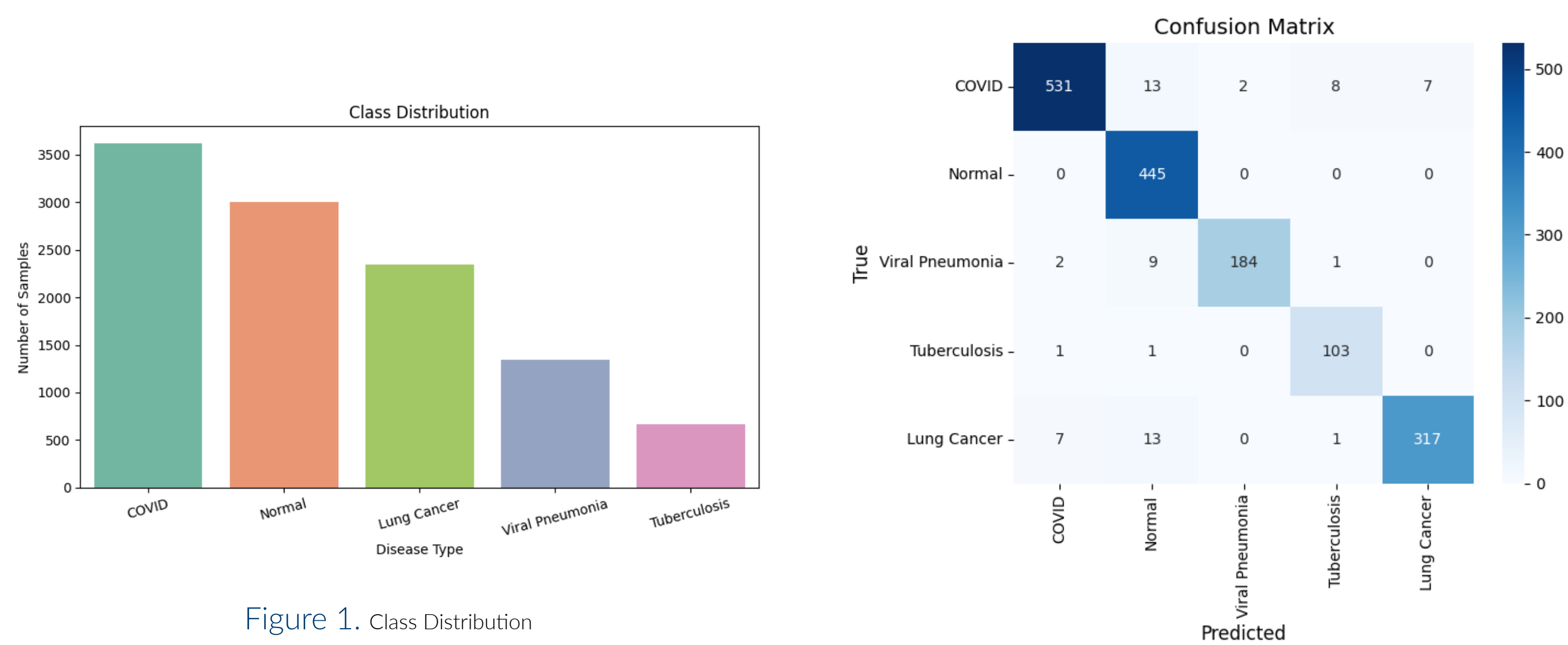


Figure 1. Class Distribution

Figure 2. Confusion Matrix

Visual Explainability

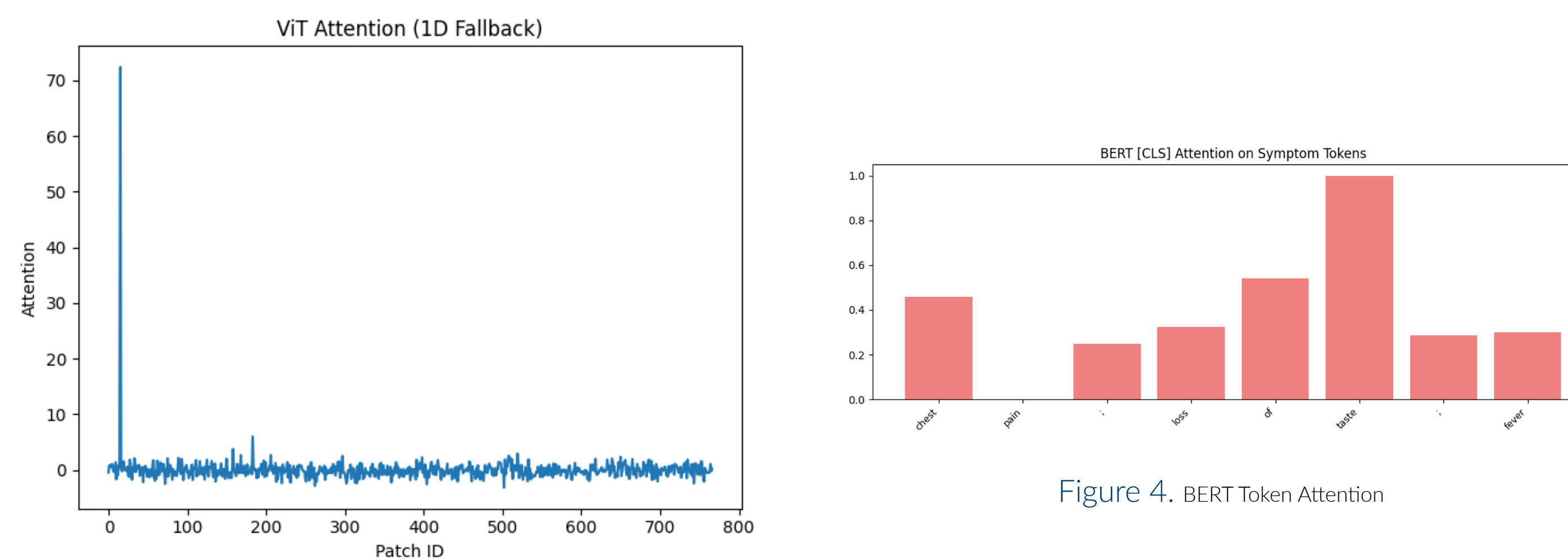


Figure 3. ViT Attention Heatmaps

Figure 4. BERT Token Attention

Live Demo Deployment

Gradio UI Features:

- Upload chest X-ray images and input patient symptoms
- Receive top-3 predictions with confidence scores
- Visualize attention heatmaps and download detailed diagnosis reports

Training Diagnostics

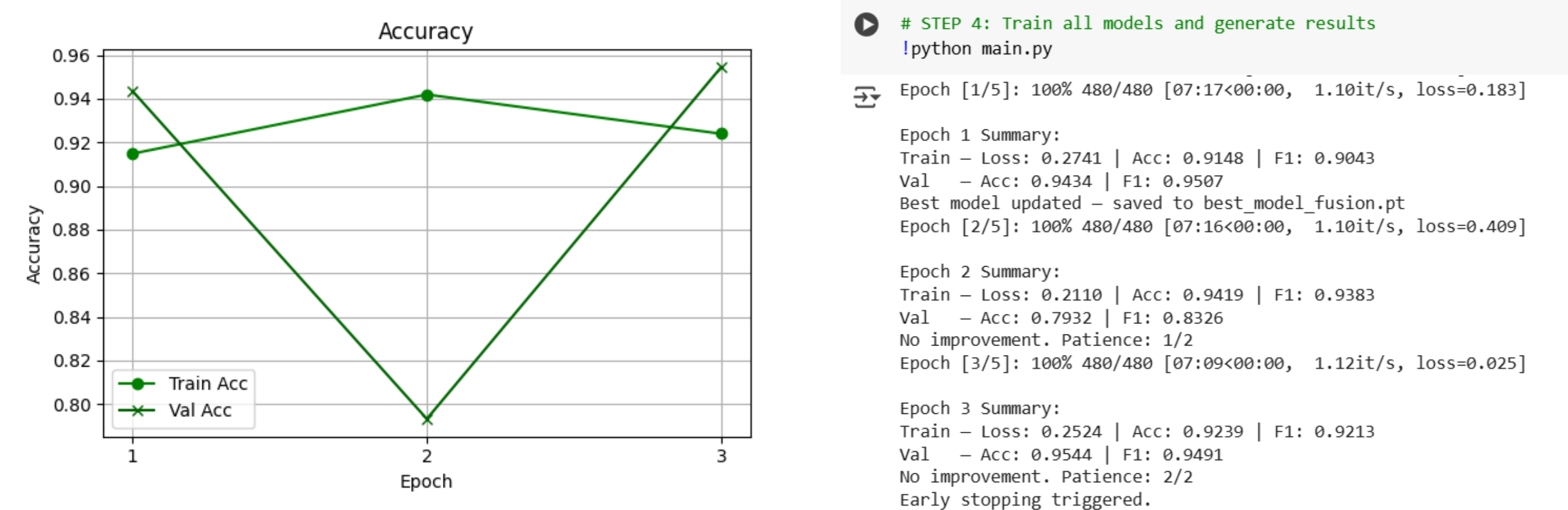


Figure 5. Training Metrics Over Epochs

Figure 6. Early Stopping Visualization

Reproducibility Notes:

- Training conducted on Google Colab Pro with Tesla T4 GPU
- Random seeds set for reproducibility

Conclusion

- Successfully integrated multimodal data for respiratory disease diagnosis
- Achieved high accuracy and F1-score across all classes
- Provided model interpretability through attention mechanisms
- Developed a user-friendly interface for real-time diagnosis

Future Work

- Incorporate CT images and electronic health records (EHRs)
- Validate the model in clinical settings
- Explore integration with hospital information systems

Why Multi-Modal AI Matters

- X-rays alone are limited by ambiguity in visual patterns.
- Symptoms add valuable context — especially in early-stage or ambiguous cases.
- ViT+BERT fusion bridges visual and semantic understanding, improving diagnostic decisions.

References

- Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers," NAAACL, 2019.
- Dosovitskiy et al., "An Image is Worth 16x16 Words," ICLR, 2021.
- Wang et al., "ChestX-ray8: Hospital-scale Dataset," CVPR, 2017.
- Rahman et al., "COVID-19 Radiography Database," IEEE Access, 2020.
- Huang et al., "Fusion of Imaging and Clinical Data with Transformer," IEEE JBHI, 2021.
- Raj et al., "TransFusion: Transfer Learning for Medical Imaging," NeurIPS, 2021.