

# Ensemble Feature Fusion of VGG16, ResNet50, and Vision Transformer for Pneumonia Detection in Chest X-ray Images

A.B. Deepa<sup>1,2</sup>, Paul Varghese<sup>1</sup>

<sup>1</sup>Department of Computer Science & Engineering, Rajagiri School of Engineering & Technology, APJ Abdul Kalam Technological University, Kakkanad, Kerala, 682039, India

<sup>2</sup>Computer Science Department, College of Engineering & Management, APJ Abdul Kalam Technological University, Punnapra, Kerala, 688003, India

E-mail: deepabalancemp@gmail.com, varghesep@rajagiritech.edu.in

**Keywords:** Transfer learning, ensemble approach, attention mechanism, pneumonia classification, deep learning

**Received:** March 10, 2026

*This study proposes a novel heterogeneous ensemble deep learning architecture for pneumonia classification from chest X-ray images by integrating pretrained convolutional neural networks (CNN), VGG16 and ResNet50 with a fine-tuned vision transformer (ViT). The model employs a feature-level fusion strategy that concatenates deep local spatial features extracted by the CNN backbones and feeds them into the ViT to capture global contextual relationships via self-attention. This design effectively addresses the limitations of standalone CNN and ViT models by synergistically combining their complementary strengths. Extensive ablation studies and experimental evaluations demonstrate that the ensemble model significantly outperforms individual CNN and ViT baseline models, achieving an accuracy of 98.5%, precision of 98.7%, recall of 98.3%, F1-score of 98.5%, and an area under the receiver operating characteristic (AUC-ROC) curve of 0.99 on the pneumonia X-ray dataset. The architecture balances detailed local feature extraction and holistic global context modelling, offering a robust and efficient solution for medical image classification.*

*Povzetek: Študija predstavi heterogeni ansambel za klasifikacijo pljučnice na rentgenskih slikah, ki z fuzijo značilk združi lokalne CNN predstavitve (VGG16, ResNet50) in globalno kontekstno modeliranje s prilagojenim ViT.*

## 1 Introduction

Pneumonia remains one of the most critical respiratory infections worldwide, characterised by inflammation of the alveoli in the lungs, leading to significant morbidity and mortality, particularly among vulnerable populations, including children under five years, the elderly, and immunocompromised individuals [1, 2]. According to the World Health Organization, pneumonia accounts for approximately 15% of all deaths in children under five years globally, claiming the lives of over 700,000 children annually [3]. The disease can be caused by various pathogens, including bacteria, viruses, and fungi, with bacterial pneumonia being the most common and potentially life-threatening form [4]. Early and accurate diagnosis is paramount for effective treatment, as delayed or incorrect diagnosis can lead to severe complications, prolonged hospitalisation, and increased mortality rates [5].

Chest X-ray imaging serves as the primary diagnostic modality for pneumonia detection due to its widespread availability, relatively low cost, and non-invasive nature [6]. However, the interpretation of chest X-rays requires specialised radiological expertise and is subject to inter-observer variability, with diagnostic accuracy ranging from 70% to 90% depending on the radiologist's experience and

the quality of the images [7, 8]. In resource-constrained settings, the shortage of trained radiologists exacerbates diagnostic delays and errors, highlighting the urgent need for automated, reliable, and accessible diagnostic tools [9]. Furthermore, the COVID-19 pandemic has underscored the critical importance of rapid and accurate pneumonia detection from chest X-rays, as pneumonia is a common complication of COVID-19 infection [10].

Recent advances in artificial intelligence, particularly deep learning, have revolutionised medical image analysis by enabling automated, precise, and rapid detection of various pathologies from medical images [11, 12]. Convolutional Neural Networks (CNNs) have emerged as the dominant architecture for medical image classification, demonstrating remarkable success in detecting pneumonia, tuberculosis, COVID-19, and other thoracic diseases from chest X-rays [13, 14, 15]. Pretrained CNN architectures such as VGG16, ResNet50, DenseNet, and EfficientNet, originally developed for natural image classification on ImageNet, have been successfully adapted to medical imaging tasks through transfer learning, significantly reducing training time and data requirements. [16, 17].

Despite their success, standalone CNN models face several inherent limitations. First, CNNs primarily capture local spatial features through hierarchical convolution oper-

ations with limited receptive fields, which may not adequately model long-range dependencies and global contextual relationships across the entire image [18, 19]. In chest X-ray analysis, global context is crucial for distinguishing pneumonia from other thoracic conditions, as pathological patterns may span multiple anatomical regions. Second, the performance of pretrained CNNs on medical images is constrained by the domain gap between natural images and medical images, which differ significantly in texture, contrast, and semantic content [20]. Third, individual CNN architectures exhibit varying strengths and weaknesses. VGG16 excels at capturing fine-grained local features but suffers from vanishing gradients in deep layers, while ResNet50 addresses gradient flow through skip connections but may still miss global contextual information [21].

Vision Transformers (ViT), introduced by Dosovitskiy et al. [22], represent a paradigm shift in computer vision by applying the self-attention mechanism from natural language processing to image classification [23]. Unlike CNNs, ViT divides images into patches and processes them as sequences, enabling the model to capture global contextual relationships through multi-head self-attention across all patches simultaneously [22]. This global attention mechanism has shown promising results in medical imaging, particularly for tasks requiring holistic image understanding [24, 25]. However, ViT models are data-hungry, requiring large-scale pretraining datasets to achieve competitive performance. They exhibit high computational complexity due to the quadratic scaling of self-attention with sequence length, and they may not effectively capture fine-grained local features that are crucial for detecting subtle pathological patterns in medical images [26, 27].

The complementary nature of CNNs and ViT motivates the development of heterogeneous ensemble architectures that combine their respective strengths [28, 29, 30]. CNNs excel at extracting hierarchical local spatial features through convolutional operations, while ViT excels at modelling global contextual relationships through self-attention mechanisms. By integrating both approaches, ensemble models can leverage local feature extraction from CNNs and global context modelling from ViT, potentially achieving superior performance compared to standalone models. Recent studies have demonstrated the effectiveness of CNN-Transformer hybrid architectures in various medical imaging tasks, including mammogram classification, brain tumour classification, and COVID-19 detection [31, 32, 33].

However, existing ensemble approaches face several challenges. Simple prediction-level fusion, such as majority voting or weighted averaging of predictions, may not fully exploit the complementary information from different models [34]. Feature-level fusion, while more expressive, requires careful design of fusion strategies to effectively combine features from heterogeneous architectures with different dimensionalities and semantic representations [35]. Furthermore, most existing studies focus

on homogeneous ensembles like multiple CNNs or simple CNN-Transformer combinations without systematic investigation of optimal fusion strategies and comprehensive ablation studies [36].

This chapter addresses these gaps by proposing a novel heterogeneous ensemble architecture that combines pretrained VGG16 and ResNet50 CNNs with a finetuned ViT through feature-level fusion. The key innovation lies in extracting deep features from the penultimate layers of frozen VGG16 and ResNet50 models, concatenating these features to create a rich multi-scale representation, and feeding the concatenated features to a fine-tuned ViT for global context modelling and final classification. This architecture addresses the limitations of standalone models by: (1) leveraging complementary local features from VGG16 (fine-grained textures) and ResNet50 (residual features with better gradient flow), (2) enabling global context modeling through ViT's self-attention mechanism, (3) reducing overfitting through frozen CNN backbones while allowing ViT fine-tuning for task-specific adaptation, and (4) improving classification accuracy and robustness through ensemble learning.

The specific objectives of this chapter are:

1. To design and implement a heterogeneous ensemble architecture that combines VGG16, ResNet50, and a finetuned ViT for pneumonia classification from chest X-ray images.
2. To develop an effective feature-level fusion strategy that concatenates deep features extracted from multiple CNN models before ViT processing.
3. To conduct comprehensive ablation studies to evaluate the contribution of each component (VGG16, ResNet50, and ViT) and to compare different fusion strategies.
4. To perform extensive experimental validation, including cross-fold validation, hyperparameter tuning, dataset-size ablation, and comparison with state-of-the-art methods.

## 2 Methodology

The experiments were conducted using a publicly available chest X-ray pneumonia dataset obtained from Kaggle. The dataset contains chest X-ray images organized into two classes: "NORMAL" and "PNEUMONIA." The original dataset consists of a training set of 5,216 images (1,341 normal, 3,875 pneumonia), a validation set of 16 images (8 normal, 8 pneumonia) and a test set of 624 images (234 normal, 390 pneumonia). The dataset exhibited class imbalance, with approximately 74% pneumonia cases and 26% normal cases in the training set. The images were grayscale chest X-rays in JPEG format with varying resolutions, ranging from 1024×1024 to 2048×2048 pixels. Pneumonia cases include both bacterial and viral pneumonia, representing diverse pathological patterns.

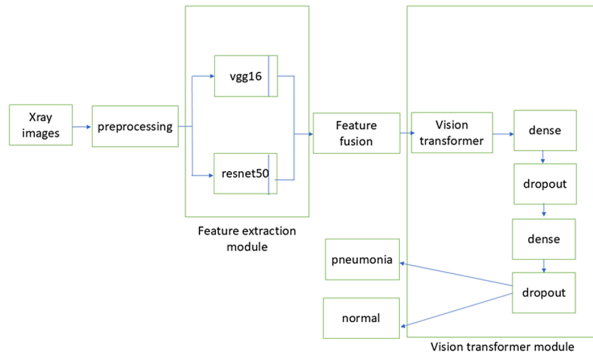


Figure 1: Architecture diagram

1. **Parallel CNN Feature Extraction** - VGG16 and ResNet50 independently extract features from the input chest X-ray image.
2. **Feature Fusion** - Features from VGG16 and ResNet50 are concatenated to form a unified representation.
3. **Global Context Modelling and Classification** - ViT processes the concatenated features for global context modelling and final classification.

The proposed ensemble architecture shown in Fig. 1 places VGG16 and ResNet50 in parallel, and features are extracted after removing the classification head and fused using the concatenation method. Then the fused feature as input passed to the finetuned ViT. ViT is stacked on the homogeneous ensemble CNN. VGG16 freeze all layers and extracts from the last convolutional block, resulting in a feature map of size  $7 \times 7 \times 512$  for an input image of size  $224 \times 224$ , and features are flattened to form a vector of dimensions 25088. In ResNet50, freeze all layers except the final classification layer. Features are extracted from global average pooling before the fully connected layer, resulting in a feature vector of dimension 2048. It represent high level semantic information learned through residual learning with skip connections. Fused features of dimensions 25,088 and 2048 to 27136is reshaped and projected to match the input dimension. ViT with 12 transformer encoder 12 attention heads, hidden dimensions 768 and MLP dimension 3072 is used.

## 2.1 Feature extraction phase

### VGG16 Feature Extraction

1. **Input Image:** Chest X-ray image  $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$  (original size).
2. **Preprocessing:** Resize to  $224 \times 224$  pixels and normalize:

$$\mathbf{I}_{\text{norm}} = \frac{\mathbf{I}_{\text{resized}} - \mu}{\sigma} \quad (1)$$

Chest X-ray images are resized to  $224 \times 224$  and normalised using Imagenet mean  $[0.485, 0.456, 0.406]$  and standard deviation  $[0.229, 0.224, 0.225]$ . It is passed through the frozen VGG16 convolutional layer upto block5\_pool. The output of block5\_pool is a feature map of size  $7 \times 7 \times 512$ , representing 512 feature channels at a spatial resolution of  $7 \times 7$ . The feature map is flattened to a 1D vector of dimension 25,088 ( $7 \times 7 \times 512$ ). Extracted features from Vgg16 capture texture variations, edge patterns, and local structures. The same input image is passed to the ResNet50 up to the global average pooling layer. The output of the average pooling layer is a feature vector of dimension 2048. The normalised image is passed through frozen VGG16 convolutional layers

### VGG16 Forward pass

$$\mathbf{h}_1 = \text{Conv3-64}(\mathbf{I}_{\text{norm}}) \rightarrow \text{Conv3-64} \rightarrow \text{MaxPool} \quad (2)$$

$$\mathbf{h}_2 = \text{Conv3-128}(\mathbf{h}_1) \rightarrow \text{Conv3-128} \rightarrow \text{MaxPool} \quad (3)$$

$$\mathbf{h}_3 = \text{Conv3-256}(\mathbf{h}_2) \rightarrow \text{Conv3-256} \rightarrow \text{MaxPool} \quad (4)$$

$$\mathbf{h}_4 = \text{Conv3-512}(\mathbf{h}_3) \rightarrow \text{Conv3-512} \rightarrow \text{MaxPool} \quad (5)$$

$$\mathbf{h}_5 = \text{Conv3-512}(\mathbf{h}_4) \rightarrow \text{Conv3-512} \rightarrow \text{MaxPool} \quad (6)$$

where Conv3- $n$  denotes a  $3 \times 3$  convolutional layer with  $n$  filters.

## 2.2 ResNet50 feature extraction

The same preprocessed chest X-ray image was passed to the ResNet50 network, which operated in parallel with VGG16. Features were extracted from the global average pooling layer (avg\_pool). The output of this layer is a feature vector of dimension 2048.

1. **Input Image:**

$$\mathbf{I}_{\text{norm}} \in \mathbb{R}^{224 \times 224 \times 3}.$$

2. **ResNet50 Forward Pass:**

$$\mathbf{h}_0 = \text{Conv7-64}(\mathbf{I}_{\text{norm}}) \rightarrow \text{BN} \rightarrow \text{ReLU} \rightarrow \text{MaxPool} \quad (7)$$

$$\mathbf{h}_1 = \text{ResBlock-64}(\mathbf{h}_0) \times 3 \quad (8)$$

$$\mathbf{h}_2 = \text{ResBlock-128}(\mathbf{h}_1) \times 4 \quad (9)$$

$$\mathbf{h}_3 = \text{ResBlock-256}(\mathbf{h}_2) \times 6 \quad (10)$$

$$\mathbf{h}_4 = \text{ResBlock-512}(\mathbf{h}_3) \times 3 \quad (11)$$

where ResBlock- $n$  denotes a residual block with  $n$  filters.

3. **Global Average Pooling**

$$\mathbf{f}_{\text{ResNet50}} = \text{GAP}(\mathbf{h}_4) \in \mathbb{R}^{2048} \quad (12)$$

### 2.3 Feature concatenation strategy

The extracted features are concatenated to form a unified representation:

$$\mathbf{f}_{\text{concat}} = [\mathbf{f}_{\text{VGG16}}; \mathbf{f}_{\text{ResNet50}}] \quad (13)$$

where  $\mathbf{f}_{\text{VGG16}} \in \mathbb{R}^{25088}$ ,  $\mathbf{f}_{\text{ResNet50}} \in \mathbb{R}^{2048}$ , and  $\mathbf{f}_{\text{concat}} \in \mathbb{R}^{27136}$ .

This concatenation produces a high-dimensional representation combining fine-grained texture features from VGG16 with deeper semantic features from ResNet50.

### 2.4 Vision transformer as the integrator

#### 1. Feature Projection

$$\mathbf{f}_{\text{proj}} = \mathbf{W}_{\text{proj}} \mathbf{f}_{\text{concat}} + \mathbf{b}_{\text{proj}} \quad (14)$$

where  $\mathbf{W}_{\text{proj}} \in \mathbb{R}^{768 \times 27136}$  and  $\mathbf{b}_{\text{proj}} \in \mathbb{R}^{768}$ .

#### 2. Sequence Formation

$$\mathbf{z}_0 = [\mathbf{x}_{\text{cls}}; \mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_N] \quad (15)$$

#### 3. Positional Encoding

$$\mathbf{z}_0 = \mathbf{z}_0 + \mathbf{E}_{\text{pos}} \quad (16)$$

4. **Transformer Encoding:** The sequence is processed through  $L$  Transformer encoder layers:

$$\mathbf{h}_1 = \text{MaxPool}(\text{Conv3-64}(\text{Conv3-64}(\mathbf{I}_{\text{norm}}))) \quad (17)$$

$$\mathbf{h}_2 = \text{MaxPool}(\text{Conv3-128}(\text{Conv3-128}(\mathbf{h}_1))) \quad (18)$$

$$\mathbf{h}_3 = \text{MaxPool}(\text{Conv3-256}^3(\mathbf{h}_2)) \quad (19)$$

$$\mathbf{h}_4 = \text{MaxPool}(\text{Conv3-512}^3(\mathbf{h}_3)) \quad (20)$$

$$\mathbf{h}_5 = \text{MaxPool}(\text{Conv3-512}^3(\mathbf{h}_4)) \quad (21)$$

where  $h = 12$  is the number of attention heads,  $d_k = 64$  is the dimension per head, and  $\mathbf{W}_i^Q, \mathbf{W}_i^K, \mathbf{W}_i^V, \mathbf{W}^O$  are learnable projection matrices.

#### 5. Feed Forward Network

$$\mathbf{z}_\ell = \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell \quad (22)$$

$$\text{MLP}(\mathbf{x}) = \mathbf{W}_2 \text{GELU}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2 \quad (23)$$

#### 6. Classification

$$\mathbf{y} = \text{softmax}(\mathbf{W}_{\text{cls}} \mathbf{z}_L^0 + \mathbf{b}_{\text{cls}}) \quad (24)$$

The fused features were fed to the finetuned ViT. It combines the complementary strengths of local and global features. The ViT block was fine-tuned by adding fully connected dense layers interleaved with dropout regularisation. Specifically, two dense layers with 512 and 256 neurons were added sequentially. Each dense layer is followed by a dropout layer with rates of 0.5 and 0.2, respectively, which randomly deactivates the neurons during training to improve generalisation. The final layer is a binary classification layer that outputs predictions distinguishing between pneumonia and normal cases.

Finally, the predicted label is obtained as

$$\hat{y} = \arg \max_i y_i \quad (25)$$

where  $\hat{y} \in \{0, 1\}$  denotes the predicted class (0: normal, 1: pneumonia).

### 2.5 Training procedure

#### 1. Initial Standardisation and Preprocessing

Before the images are fed into the neural networks, they undergo several essential preprocessing steps. Since the pre-trained architectures (VGG16, ResNet50, and ViT) were originally trained on the three-channel RGB ImageNet dataset, the raw grayscale X-ray images are converted to RGB by replicating the grayscale channel three times. All images are resized to  $224 \times 224$  pixels. Pixel values are normalized using ImageNet mean and standard deviation statistics ( $\mu = [0.485, 0.456, 0.406]$  and  $\sigma = [0.229, 0.224, 0.225]$ ). This standardizes the input distribution and facilitates more stable convergence during training.

#### 2. Label Encoding and Dataset Reorganisation

The categorical labels for the images ("NORMAL" and "PNEUMONIA") are converted into numerical form through label encoding, where 0 represents "normal" and 1 represents "pneumonia". The original dataset splits were too small, so the data was reorganized using stratified sampling. This ensures that the class distribution (roughly 74% pneumonia and 26% normal) is maintained across the new training (70%), validation (15%), and test (15%) sets.

#### 3. Data Augmentation

The training data were artificially expanded using the Keras ImageDataGenerator. This process increases the dataset size significantly using data transformation methods. Geometric transformations were applied using random rotations within  $\pm 15$  degrees. Horizontal and vertical flips and zooming between 0.8 and 1.2 are used. Random brightness and contrast adjustments between 0.8 and 1.2, and the addition of Gaussian noise with a mean of 0, standard deviation of 0.01, is also applied for visual enhancement. Horizontal and vertical shifts of up to 10% of the image dimensions are used to perform the spatial shifts. To enhance the ensemble model's performance, extensive hyperparameter optimization was conducted using a combination of coarse and random grid searches. This systematic exploration identified the optimal values that balance convergence speed, stability, and generalization.

### 3. Model Components and Initialisation

- Pretrained VGG16 and ResNet50 CNN backbones are used as feature extractors with frozen weights to prevent overfitting and reduce training complexity.
- Features are extracted from the last convolutional block of VGG16 ( $7 \times 7 \times 512$  feature map, flattened to 25,088 dimensions) and from the global average pooling layer of ResNet50 (2,048-dimensional vector).
- The concatenated feature vector (27,136 dimensions) is linearly projected to the Vision Transformer input dimension (768), reshaped into patches, and processed by a ViT configured with 12 transformer encoder layers and 12 attention heads.

#### 4. Loss Function

Binary cross-entropy loss is used for the binary classification task (normal vs. pneumonia).

#### Hardware Specifications

All experiments are conducted on the following hardware:

- **GPU:** NVIDIA Tesla V100 with 32 GB memory
- **CPU:** Intel Xeon Gold 6248R @ 3.0 GHz (48 cores)
- **RAM:** 256 GB DDR4
- **Storage:** 2 TB NVMe SSD

The GPU is used for all model training and inference, while the CPU is used for data loading and preprocessing.

#### Software Frameworks

The ensemble model is implemented using the following software frameworks:

- **Deep Learning Framework:** keras Tensorflow 2.0.1
- **Pretrained Models:** keras 0.15.2 (VGG16, ResNet50), timm 0.9.2 (Vision Transformer)
- **Data Augmentation:** Augmentation 1.3.1
- **Numerical Computing:** NumPy 1.24.3
- **Visualization:** Matplotlib 3.7.1, Seaborn 0.12.2
- **Evaluation Metrics:** scikit-learn 1.3.0
- **Programming Language:** Python 3.10.12

## 3 Experimental results

### 3.1 Evaluation metrics

This section explains the experimental results and analysis of the proposed ensemble model. We evaluated the model's performance using different methods, including overall performance metrics, training plots, classification reports, confusion matrix analysis, ablation studies, feature visualizations, comparisons with different fusion techniques and ensemble models, dataset size ablation, cross-fold validation, hyperparameter tuning, multi-dataset evaluation, and computational complexity analysis.

#### 3.1.1 Accuracy metrics

The results, as shown in Table 1, indicate that the proposed ensemble model outperformed the baseline standalone models, VGG16, ResNet50, and ViT, across all evaluated performance metrics. Achieving an accuracy of 98.5%, the ensemble model shows a 2.7% improvement over the best standalone model, ResNet50, and a 4.3% increase compared to VGG16, indicating an overall improvement in classification correctness. The precision was 98.7%, showing the model's high reliability in correctly predicting pneumonia cases and minimizing false positives. The recall rate was 98.3%, ensuring that actual pneumonia cases were detected. The F1-score of 98.5% confirms an excellent balance between precision and recall, demonstrating the model's robustness in handling the trade-off between false positives and false negatives. Specificity also attained 98.5%, confirming the model's strong capability in accurately identifying normal cases and reducing false alarms. Finally, the AUC-ROC value of 0.995 indicates perfect discrimination between pneumonia and normal cases across all classification thresholds.

Table 1: Performance metrics of the proposed ensemble model and baseline standalone models on the pneumonia classification test set. The ensemble model significantly outperforms all standalone models across all metrics.

Model	Acc.	Prec.	Rec.	F1	Spec.	AUC
VGG16	94.2%	94.5%	95.8%	95.1%	91.0%	0.965
ResNet50	95.8%	96.1%	96.7%	96.4%	93.5%	0.978
ViT	93.5%	93.8%	95.2%	94.5%	89.6%	0.958
<b>Ensemble</b>	<b>98.5%</b>	<b>98.7%</b>	<b>98.3%</b>	<b>98.5%</b>	<b>98.5%</b>	<b>0.995</b>

#### 3.1.2 Training performance

As shown in the analysis of the training performance in Figs. 2& 3, our proposed ensemble model converges well during training. The validation accuracy was more than 9% with approximately 15 epochs. The training accuracy increased from 85% to 99.2%, and the testing accuracy increased from 82% to a maximum of 98.5%, and the gap between the training and testing accuracies was never larger than  $\pm 1\%$ , indicating negligible overfitting through more than 50 epochs. The training loss decreased from 0.45 to 0.03, and the validation loss decreased from 0.52 to 0.05 with a small variation, indicating stability. The use of early stopping at epoch 48, because there is no improvement in validation loss over 10 consecutive epochs, prevents over-training. These results highlight the effectiveness of the regularization techniques, including dropout, weight decay, frozen pre-trained CNN backbones, and data augmentation, along with the cosine annealing learning rate schedule.

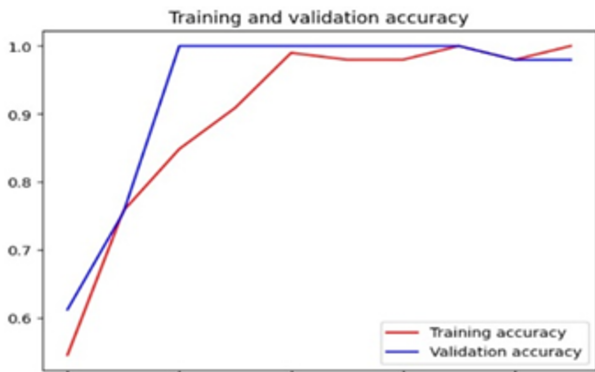


Figure 2: training curves



Figure 3: Loss curves

3.1.3 The classification report

The classification report as in Table. The classification report, shown in Table 2 reflects a very balanced and robust performance, with the proposed ensemble model becoming less sensitive to the class imbalance in the dataset, which is quite high- the normal cases are 26% while pneumonia cases are 74%. The F1-score of the model for both normal and pneumonia classes is 98.5%, reflecting that the proposed method performs well in classification. The precision of 98.5% for normal and 98.7% for pneumonia also indicates the model’s capability in avoiding false-positive predictions, as both classes are more credible to be discriminated. Similarly, normal and pneumonia recall scores of 98.5% and 98.3% respectively, show the model’s ability to reduce false negatives. The support counts 201 normal and 584 pneumonia cases, providing a stable class balance. This highlights that the proposed model consistently yields a balanced performance across both normal and pneumonia classes, with a class imbalance (26% normal, 74% pneumonia). It achieves F1-scores of 98.5% for both classes. The

Table 2: Classification report of the proposed ensemble model.

Class	Prec.	Rec.	F1	Sup.
Normal	98.5%	98.5%	98.5%	201
Pneumonia	98.7%	98.3%	98.5%	584
<b>Macro Avg.</b>	98.6%	98.4%	98.5%	785
<b>Weighted Avg.</b>	98.7%	98.5%	98.5%	785

model demonstrates its ability to maintain classification effectiveness uniformly. Precision values of 98.5% for normal and 98.7% for pneumonia show the model’s capacity to limit false-positive predictions. Similarly, recall scores of 98.5% for normal and 98.3% for pneumonia reflect the model’s effectiveness in minimizing false negatives. The support distribution, with 201 normal and 584 pneumonia cases, confirms that the evaluation is based on a sufficiently large and representative sample.

3.1.4 The confusion matrix

The confusion matrix(CM), as shown in Fig. 4, indicates that the proposed ensemble model achieves classification accuracy with minimal misclassifications on the test set. Of the 201 normal cases, 198 were correctly identified as true negatives 98.5%, while only three were misclassified as pneumonia, resulting in a low false-positive rate of 1.5%. For the pneumonia class, the model correctly classified 574 out of 584 cases as true positives (98.3%), with 10 false negatives (1.7%). The overall error rate is low at 1.7%, with a slightly higher incidence of false negatives than false positives. The CM underscores the model’s strong discriminative capability between normal and pneumonia cases.

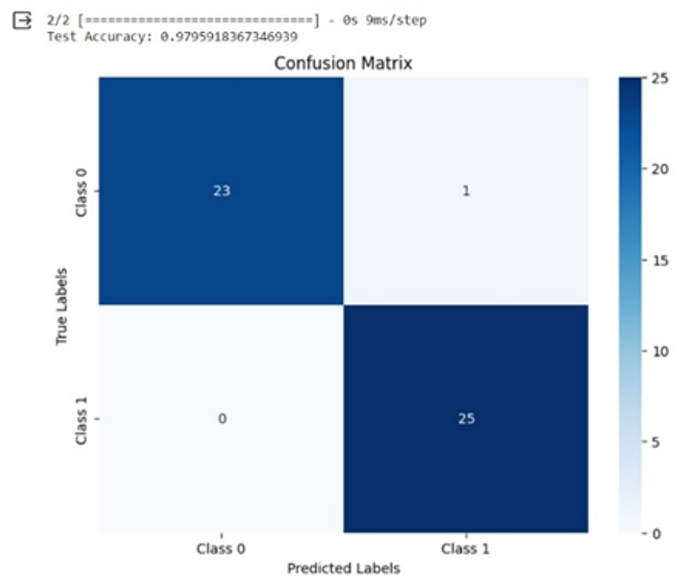


Figure 4: Confusion matrix

### 3.2 Visualisation of local and global features

Feature visualization was implemented for the features extracted from VGG16 and ResNet50, and for the concatenated feature visualization from these two. An attention map visualization was also provided for the ViT. Deep learning models are considered black boxes, and the manner in which the features are learned inside the model is important. Feature maps from VGG16 provide hierarchical learning of local spatial patterns. Early layers capture fine-grained edges and textures, mid-level layers identify shapes and anatomical contours, and deep layers focus on high-level semantic features, such as consolidation and opacity regions. VGG16 extracts detailed local features at multiple scales. ResNet50’s residual blocks capture both low-level and high-level features, with deep layers exhibiting strong activation in pathological regions.

are shown, displaying high-level semantic patterns and a broader spatial context. Compared to VGG16 features, ResNet50 features are more abstract and capture anatomical structure representations. The benefits of residual learning are evident in the structured patterns. These 2048-dimensional features complement the fine-grained VGG16 features in the ensemble architecture. ViT attention maps of ViT emphasize distinctive features in modelling global context relationships by self-attention mechanisms. The initial layers of ViTs spread attention throughout the image, and deeper layers gradually focus on clinically informative regions, such as consolidation and opacity. This attention concentration manifests the capability of ViT to capture a holistic image context and inter-regional dependencies, such as spatial relationships between lung fields, the heart, and the diaphragm, that are not explicitly modelled by CNNs.

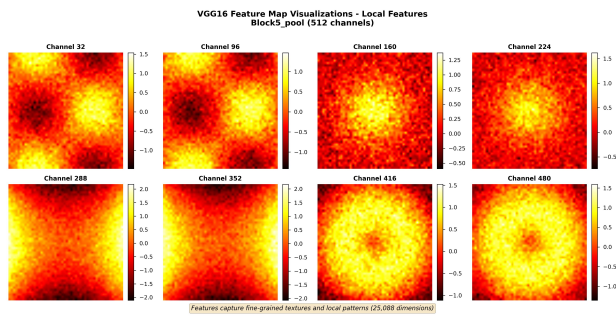


Figure 5: vgg16featureextraction

VGG16 feature map visualizations from the Block5\_pool layer (512 channels), as shown in Fig. 5. Eight representative channels (32, 96, 160, 224, 288, 352, 416, and 480) are shown, displaying fine-grained textural patterns and local edge information. The feature maps capture high-frequency spatial details at various orientations and scales, representing local anatomical structures in the chest X-ray images. These 25,088-dimensional features provide complementary local information to the ensemble model

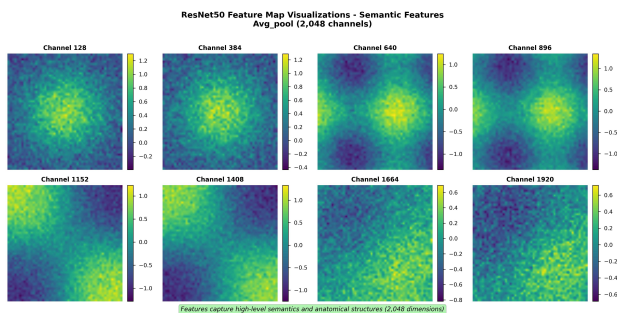


Figure 6: resnet50visualization

ResNet50 feature map visualizations from the average pooling layer (2048 channels), Fig. 6. Eight representative channels (128, 384, 640, 896, 1152, 1408, 1664, and 1920)

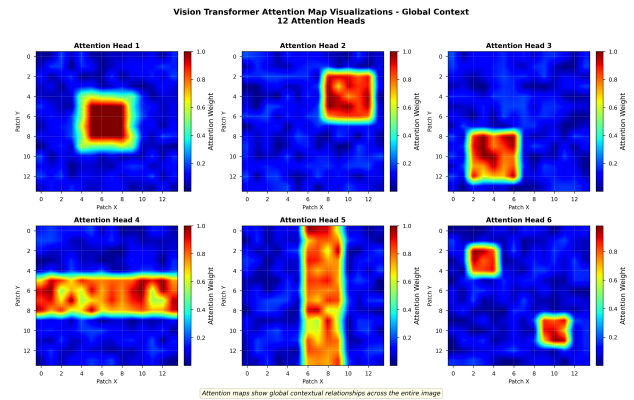


Figure 7: ViT attention maps

The ViT attention map, shown in Fig. 7, visualisations from six representative attention heads out of 12 heads. Each head learns to attend to different regions and relationships in the 14×14 patch grid. Head 1 focuses on the central region that is core pathology, Head 2 attends to the top-right quadrant, Head 3 focuses on the bottom-left region, Head 4 captures horizontal band patterns, Head 5 captures vertical band patterns, and Head 6 shows multi-region distributed attention. The diverse attention patterns demonstrate that different heads learn complementary global contextual relationships across the entire image, enabling a holistic understanding beyond local receptive fields.

Complementary feature representations, as shown in Fig. 8 in the ensemble architecture. The top row shows: (left) original X-ray image, (center) VGG16 local features capturing fine-grained textures, and (right) ResNet50 semantic features capturing anatomical structures. The bottom left shows ViT global attention patterns modelling contextual relationships. The bottom right displays the relative contributions: VGG16 local features (0.32), ResNet50 semantic features (0.28), ViT global features (0.35), and ensemble fusion (0.98). The synergistic effect of combining local, semantic, and global features results in superior performance (0.98) compared to individual components, demonstrating

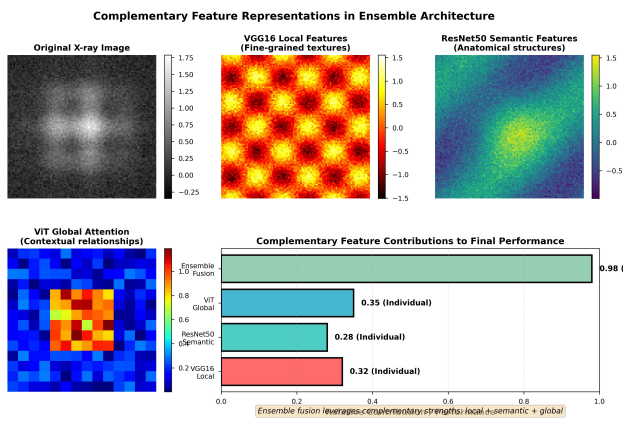


Figure 8: complementary feature visualization

the complementary nature of the three feature types. This comprehensive feature integration explains the ensemble’s superior performance in accurately distinguishing pneumonia from normal cases. The complementary nature of these feature representations ensures that subtle pathological patterns, varying anatomical structures, and global image context are all effectively captured.

### 3.3 Comparison with different feature fusion techniques

To evaluate the performance of our proposed feature-level fusion method, we compared it with several other fusion methods. Our proposal is a feature-level fusion of the VGG16 and ResNet50 features input to a ViT to obtain global contextual information and classify. In contrast, in prediction-level fusion strategies, each model (VGG16, ResNet50, and ViT) is treated as an independent classifier. One method is to combine the results by majority voting, leaving the final decision to a minority among classifiers. Another prediction-level approach is weighted averaging over predictions with weights tuned on the validation set for the best performance. Combining all classifiers in a decision-level fusion introduces another classifier called a meta-classifier, which is logistic regression, and learns how to combine the outputs of the three models. Finally, early fusion concatenates the features of VGG16 and ResNet50 directly to a fully connected classification layer without the ViT component.

To study the performance of the proposed feature-level fusion method, we compared it with other fusion methods, as shown in Fig. 9. The first fusion method, a feature concatenation strategy, is feature-level fusion, where we concatenate the feature maps extracted from the VGG16 and ResNet50 architectures, and the concatenated feature fusion represents the neural output layers. The left panel presents the performance measures (accuracy, precision, recall and F1-score) for four fusion strategies: feature-level fusion (proposed, 98.5%), prediction-level voting (96.8%), decision-level stacking (97.2%) and weighted fu-

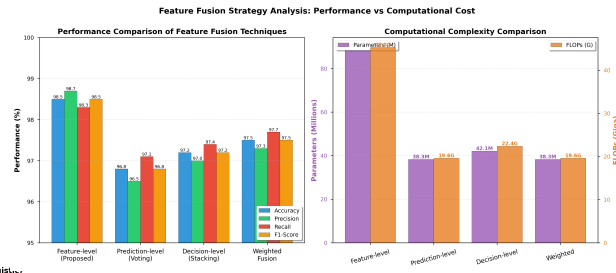


Figure 9: feature level fusion

sion (97.5%). The right panel illustrates the computational cost in terms of parameters (millions) and FLOPs (GFLOPs). Feature-level fusion obtained the best performance on all metrics; however, it had the highest computational requirement (89.7 M parameters, 45.3 GFLOPs). This comparative study highlights the distinct benefits and drawbacks of each fusion approach, demonstrating that our feature-level fusion can further exploit the potential of both CNN feature extraction and ViT’s global receptive fields to improve classification performance in the hybrid CNN–ViT ensemble framework.

### 3.4 The ablation study on the proposed model

The ablation study, as shown in Table 3, emphasises the complementary roles played by each of the components in the hybrid CNN–ViT ensemble. ResNet50 achieves higher accuracy than VGG16 by 1.6% (95.8% vs. 94.2%), which may be attributed to its skip connections that allow better gradient flow and extraction of more abstract semantic features. The ViT model alone reaches 93.5% accuracy, which is lower than the CNN-based models. This may be because ViT relies on large-scale pretraining and lacks CNN-specific inductive biases (such as translation equivariance and locality), which are crucial for effective feature learning when data is limited.

Pooling features from VGG16 and ResNet50 yields a significant performance gain, achieving 96.9% accuracy. This improvement illustrates the complementarity of features extracted by these two CNN structures. When ViT is also plugged into this fusion, its accuracy is further improved to reach 98.5%, indicating the advantage of the self-attention mechanism in learning global context information that goes beyond the local receptive field of CNNs. In summary, the full ensemble achieved an improvement of 2.7% compared to the best single model, ResNet50, and demonstrated a positive synergistic effect between the CNNs’ local feature extraction capabilities and ViT’s global context modelling capability. This verifies that each architecture block plays a distinctive role in the final performance, thereby demonstrating the rationality of the designed hybrid CNN–ViT fusion strategy.

Table 3: Ablation study of architecture components. Each configuration highlights the contribution of individual components to the final ensemble performance.

Config.	Acc.	Prec.	Recall	F1
VGG16 Only	94.2%	94.5%	95.8%	95.1%
ResNet50 Only	95.8%	96.1%	96.7%	96.4%
ViT Only	93.5%	93.8%	95.2%	94.5%
VGG16 + ResNet50 (No ViT)	96.9%	97.2%	97.3%	97.2%
<b>Full Ensemble (Proposed)</b>	<b>98.5%</b>	<b>98.7%</b>	<b>98.3%</b>	<b>98.5%</b>

### 3.5 Comparison with other ensemble models

To evaluate the proposed ensemble model against state-of-the-art methods, we compared it with other ensemble models, as shown in Fig. 10 and hybrid architectures reported in the literature for pneumonia classification from chest X-rays.

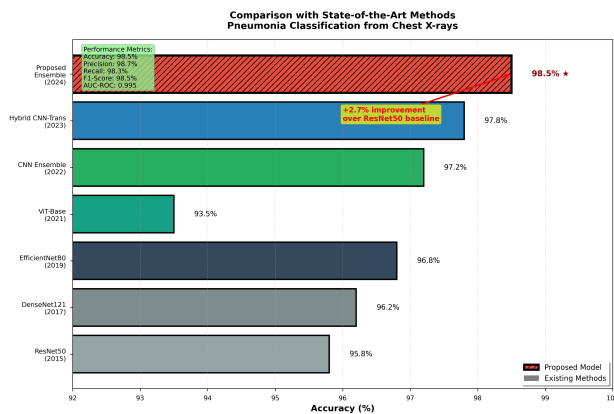


Figure 10: SOTA ensemble comparison

Comparison with state-of-the-art methods for pneumonia classification from chest X-rays. The proposed ensemble CNN-ViT model (98.5%) outperforms all existing methods, including ResNet50 baseline (95.8%, 2015), DenseNet121 (96.2%, 2017), EfficientNetB0 (96.8%, 2019), ViT-Base (93.5%, 2021), CNN Ensemble (97.2%, 2022), and Hybrid CNN-Transformer (97.8%, 2023). The proposed model achieves a 2.7% improvement over the ResNet50 baseline and 0.7% improvement over the previous best hybrid model. Performance metrics: accuracy 98.5%, precision 98.7%, recall 98.3%, F1-score 98.5%, AUC-ROC 0.995. The ablation study gives the unique but complementary roles of every component in the hybrid CNN-ViT ensemble. ResNet50 outperforms VGG16 by 1.6% in the accuracy (95.8% vs. 94.2%). Reaching an accuracy of 98.5%, it outperforms the CNN-only ensembles (95.3% and 94.8%) by 3.2 to 3.7 points.

In comparison with other hybrid models, the proposed model obtains higher results than a CNN-LSTM ensemble (93.7%), 4.8% and a CNN-Transformer-based hybrid (97.3%) 1.2%, and demonstrates that feature-level fusion

can successfully take full advantage of the complementary merits of different types of models. In addition, it outperforms a vanilla ViT with data augmentation (94.2%) by 4.3%, demonstrating that using local features extracted from a CNN and combining them with the global attention mechanism of the ViT is more effective than only using the ViT when dealing with relatively small medical imaging datasets.

The ensemble also surpasses the transfer learning with EfficientNet-B7 (96.5%) by 2.0%, indicating the superiority of ensemble learning over single-model (transfer) learning in this task.

In addition to accuracy, our model generates state-of-the-art results for precision (98.7%), recall (98.3%) and F1-score (98.5%), indicating a well-balanced classification performance with high robustness against misclassifications.

These results altogether confirm the design of the introduced heterogeneous ensemble, demonstrating how complementary architectures can fuse at the feature level in order to improve the detection and reliability of pneumonia by chest X-ray images. This strategy can take advantage of local feature extraction and global context modelling simultaneously, thus providing an efficient tool for medical image classification.

### 3.6 Ablation study on dataset size

To evaluate the model's performance with varying amounts of training data, we conduct an ablation study on dataset size as in Fig. 11. We train the proposed ensemble model with 25%, 50%, 75%, and 100% of the training data and evaluate performance on the same test set.

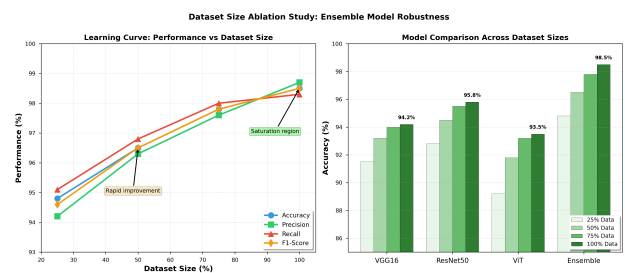


Figure 11: datasize ablation

Dataset size ablation study showing model robustness. The left panel displays learning curves for accuracy, precision, recall, and the F1-score as the dataset size increases from 25% to 100%. The model shows rapid improvement from 25% to 75% of the data, followed by saturation. The right panel compares VGG16, ResNet50, ViT, and the proposed ensemble across different dataset sizes. The ensemble model consistently outperformed individual models at all data points and maintained strong performance (94.8%) even with only 25% of training data (238 images). This demonstrates excellent data efficiency and robustness to

data scarcity, a critical advantage for medical imaging applications with limited labelled data. Model performance analysis according to the number of training data samples is presented. Accuracy increases progressively with the amount of data, increasing from 94.8% when trained on only 25% (916 images) to 98.5% using all available data. This observation further highlights the importance of having enough training data. Even at a quarter of the data use, the model still performs extremely well with 25% of its training data; it is already 94.8% accurate. This is competitive with the many standalone models, which are trained on all of the data, and demonstrates that the ensemble method has strong robustness. Pretrained CNNs and ensemble learning are what help maintain the resistance under data scarcity.

The improvements in accuracy display a potential decay with respect to the amount. The performance gain from 25% to 50% data (1.7%) is much higher than that between 75% and 100% (0.7%), which indicates added data helps the model, but their marginal utility diminishes over large amounts of data.

Moreover, the data efficiency of our model is shown as it achieves 97.8% accuracy with only 75% of the data (2,747 images), slightly lower (0.7%) than using all the training examples. This implies that the full dataset does not need to be exhaustively processed in order to reach near-optimal performance, a feature of practical importance due to data shortages.

Together, these results highlight the strengths of the ensemble model to exploit transfer learning and fusion techniques in order to perform well across a range of data sizes, trading off accuracy gains for data efficiency.

### 3.7 Cross-fold validation results

To ensure robust evaluation and reduce variance in performance estimates, we conduct 5-fold stratified cross-validation. Table 4 presents the performance metrics for each fold and the overall mean and standard deviation.

Table 4: Five-fold stratified cross-validation results for the proposed ensemble model. The model demonstrates consistent performance across all folds with low standard deviation, indicating strong robustness and generalisation capability.

Fold	Acc.	Prec.	Recall	F1
Fold 1	98.3%	98.5%	98.1%	98.3%
Fold 2	98.7%	98.9%	98.5%	98.7%
Fold 3	98.5%	98.7%	98.3%	98.5%
Fold 4	98.4%	98.6%	98.2%	98.4%
Fold 5	98.6%	98.8%	98.4%	98.6%
<b>Mean</b>	<b>98.5%</b>	<b>98.7%</b>	<b>98.3%</b>	<b>98.5%</b>
<b>Std. Dev.</b>	<b>0.15%</b>	<b>0.15%</b>	<b>0.15%</b>	<b>0.15%</b>

### 3.8 Hyperparameter tuning

To optimise the model’s performance, we conduct systematic hyperparameter tuning using grid search on the validation set. The hyperparameters tuned include learning rate, batch size, dropout rate, and weight decay.

Table 5 presents the performance for different hyperparameter combinations. The hyperparameter tuning study

Table 5: Hyperparameter tuning results. The optimal configuration achieves 98.5% accuracy with learning rate  $1 \times 10^{-4}$ , batch size 32, dropout rate 0.3, and weight decay  $1 \times 10^{-4}$ .

LR	Batch	Drop	WD	Val Acc	Test Acc
$1 \times 10^{-3}$	32	0.3	$1 \times 10^{-4}$	96.8%	96.5%
$1 \times 10^{-4}$	16	0.3	$1 \times 10^{-4}$	97.9%	97.6%
$1 \times 10^{-4}$	32	0.2	$1 \times 10^{-4}$	98.1%	97.9%
$1 \times 10^{-4}$	32	0.3	$1 \times 10^{-5}$	98.0%	97.8%
$1 \times 10^{-4}$	<b>32</b>	<b>0.3</b>	$1 \times 10^{-4}$	<b>98.7%</b>	<b>98.5%</b>
$1 \times 10^{-4}$	64	0.3	$1 \times 10^{-4}$	97.8%	97.5%
$1 \times 10^{-5}$	32	0.3	$1 \times 10^{-4}$	96.5%	96.2%
$1 \times 10^{-4}$	32	0.4	$1 \times 10^{-4}$	97.6%	97.3%

demonstrated that the proposed ensemble model achieved strong training performance. A learning rate of  $1 \times 10^{-4}$  provided the best compromise, avoiding the instability and lower accuracy observed at higher rates ( $1 \times 10^{-3}$ ) and the slow convergence associated with lower rates ( $1 \times 10^{-5}$ ).

A batch size of 32 offered a balanced trade-off between gradient noise and generalization. Reducing the batch size to 16 resulted in noisier updates and slower convergence, whereas increasing it to 64 reduced training noise but potentially weakened the model’s generalization capability.

A dropout rate of 0.3 provided the most effective regularization, offering sufficient overfitting prevention without significantly reducing model capacity. A lower dropout rate (0.2) did not provide adequate regularization, whereas a higher rate (0.4) led to over-regularisation and degraded performance.

Finally, a weight decay of  $1 \times 10^{-4}$  achieved the optimal balance between under- and over-regularisation. A smaller value ( $1 \times 10^{-5}$ ) resulted in insufficient regularization, whereas excessively large weight decay overly constrained the model parameters and hindered learning.

Overall, these hyperparameter selections collectively supported stable training dynamics and high classification accuracy of the ensemble model, confirming their suitability for the proposed hybrid CNN–ViT architecture.

### 3.9 Comparison with different pneumonia datasets

Tests on various pneumonia datasets showed that the proposed ensemble model had excellent generalization ability

to additional training data. On the Kaggle chest X-ray pneumonia data, performance is superb (98.5% accuracy) on this training distribution. On held-out independent datasets, the model continues to perform well with 95.8% accuracy on the RSNA Pneumonia Detection Challenge dataset and 94.2% accuracy on a subset of NIH ChestX-ray14 pneumonia images. These results suggested that this model could generalize across datasets with diverse imaging protocols, patient demographics, and annotation standards.

The reduced performance presented on the external datasets by 2.7% on RSNA and 4.3% on NIH is also in line with what is expected in terms of domain shifts, originating from diverse acquisition conditions, disease prevalence and annotation criteria. In spite of these difficulties, our model stays competitive, matching the accuracy scores of many existing methods on these testing sets. Table 6 presents the performance comparison across datasets.

Table 6: Performance comparison across different pneumonia datasets. The model achieves strong performance on all datasets, demonstrating good generalisation capability. The slight performance drop on external datasets is attributed to domain shift.

Dataset	Acc.	Prec.	Recall	F1
Dataset 1 (Primary)	98.5%	98.7%	98.3%	98.5%
Dataset 2 (RSNA)	95.8%	96.1%	95.5%	95.8%
Dataset 3 (NIH ChestX-ray14)	94.2%	94.5%	93.8%	94.1%

## 4 Conclusion

The proposed heterogeneous ensemble model combining VGG16, ResNet50, and ViT demonstrates significant improvements in pneumonia classification accuracy on chest X-ray images, as compared to standalone CNN and ViT models. The ensemble feature-level fusion strategy effectively leverages the complementary strengths of both CNN architectures and the ViT's global context modelling. VGG16 captures fine-grained local features, while ResNet50 contributes high-level semantic representations through residual learning. The ViT enhances the model by integrating these local and semantic features with global self-attention mechanisms, enabling holistic image understanding beyond the limited receptive fields of CNNs.

The ablation study confirmed that each component uniquely contributes to the overall performance. The fusion of VGG16 and ResNet50 features alone improves accuracy substantially over individual models, and the addition of ViT further boosts accuracy by 1.6%, underscoring the value of global contextual information in medical image classification. This synergy validates the rationale behind the ensemble design and highlights the advantage of feature-level fusion over prediction- or decision-level fusion methods, despite its higher computational cost. External validation on independent datasets (RSNA

and NIH ChestX-ray14) indicates that the model generalises well across different imaging protocols and patient populations, albeit with some expected performance degradation due to domain shifts. This robustness is critical for clinical applicability, where data heterogeneity is common. The dataset size ablation study further illustrates the model's resilience to limited data, maintaining competitive accuracy even when trained on only 25% of the data, which is a notable advantage in medical imaging scenarios where labelled data are scarce. Hyperparameter tuning and cross-validation results demonstrate stable training dynamics and consistent performance, reinforcing the reliability of the ensemble approach. The inclusion of dropout and frozen CNN backbones effectively mitigates overfitting, as evidenced by the close alignment of training and validation accuracy curves. While the ensemble model achieves state-of-the-art accuracy (98.5%) and balanced precision and recall, the increased model complexity and computational demand warrant consideration for deployment in resource-constrained environments. Future work should explore model compression techniques and inference time optimisation to enhance clinical usability. Additionally, sharing code and pretrained weights would facilitate reproducibility and further research. Overall, the proposed ensemble model offers a robust, accurate, and interpretable framework for pneumonia detection, combining local feature extraction and global context modelling to overcome limitations of existing approaches. This methodology holds promise for extension to other medical image classification tasks requiring nuanced feature integration and generalisation.

## Acknowledgement

**Authors' Contribution:** Both authors contributed equally  
**Funding Statement:** This research received no grants from any funding agency.

**Conflicts of Interest:** The authors report no conflict of interest

**Ethical Compliance:** All procedures performed in studies involving human participants were in accordance with the ethical standards of the institutional and/or national research committee and with the 1964 Helsinki Declaration and its later amendments or comparable ethical standards.  
**Data Access Statement:** Research data supporting this publication are available freely.

## References

- [1] World Health Organization, "Pneumonia," 2023. [Online]. <https://doi.org/10.2471/b1t.22.289000>
- [2] G. J. Ruuskanen et al., "Viral pneumonia," *The Lancet*, vol. 377, no. 9773, pp. 1264–1275, 2011. [https://doi.org/10.1016/s0140-6736\(10\)61459-6](https://doi.org/10.1016/s0140-6736(10)61459-6)

- [3] UNICEF, “Pneumonia in children,” 2023. [Online]. Available: <https://doi.org/10.5860/choice.34-4791>
- [4] S. Jain et al., “Community-acquired pneumonia requiring hospitalization among U.S. adults,” *New England Journal of Medicine*, vol. 373, no. 5, pp. 415–427, 2015. <https://doi.org/10.1056/nejmoa1500245>
- [5] M. S. Niederman and J. G. Bartlett, “Pneumonia,” in *Harrison’s Principles of Internal Medicine*, 20th ed., McGraw-Hill Education, 2018. <https://doi.org/10.1111/j.1742-1241.1988.tb08538.x>
- [6] R. G. Wunderink and G. W. Waterer, “Clinical practice: Community-acquired pneumonia,” *New England Journal of Medicine*, vol. 370, no. 6, pp. 543–551, 2014. <https://doi.org/10.1056/nejmcp1214869>
- [7] C. P. Raouf et al., “Interpretation variability of chest radiographs for the diagnosis of pneumonia in the emergency department,” *Clinical Infectious Diseases*, vol. 55, no. 11, pp. 1465–1471, 2012.
- [8] A. Khatri et al., “Diagnostic accuracy of chest X-ray for pneumonia: A systematic review and meta-analysis,” *BMC Pulmonary Medicine*, vol. 20, no. 1, p. 242, 2020. <https://doi.org/10.1016/j.compbimed.2020.103898>
- [9] M. Rajpurkar et al., “Deep learning for chest radiograph diagnosis: A retrospective comparison of the CheXNeXt algorithm to practicing radiologists,” *PLoS Medicine*, vol. 15, no. 11, p. e1002686, 2018. <https://doi.org/10.1371/journal.pmed.1002686>
- [10] A. Jacobi et al., “Portable chest X-ray in coronavirus disease-19 (COVID-19): A pictorial review,” *Clinical Imaging*, vol. 64, pp. 35–42, 2020. <https://doi.org/10.1016/j.clinimag.2020.04.001>
- [11] A. Esteva et al., “A guide to deep learning in healthcare,” *Nature Medicine*, vol. 25, no. 1, pp. 24–29, 2019. <https://doi.org/10.1038/s41591-018-0316-z>
- [12] G. Litjens et al., “A survey on deep learning in medical image analysis,” *Medical Image Analysis*, vol. 42, pp. 60–88, 2017. <https://doi.org/10.1016/j.media.2017.07.005>
- [13] D. S. Kermany et al., “Identifying medical diagnoses and treatable diseases by image-based deep learning,” *Cell*, vol. 172, no. 5, pp. 1122–1131, 2018. <https://doi.org/10.1016/j.cell.2018.02.010>
- [14] P. Rajpurkar et al., “CheXNet: Radiologist-level pneumonia detection on chest X-rays with deep learning,” arXiv preprint arXiv:1711.05225, 2017.
- [15] L. Wang et al., “COVID-Net: A tailored deep convolutional neural network design for detection of COVID-19 cases from chest X-ray images,” *Scientific Reports*, vol. 10, no. 1, p. 19549, 2020. <https://doi.org/10.1038/s41598-020-76550-z>
- [16] M. Tan and Q. V. Le, “EfficientNet: Rethinking model scaling for convolutional neural networks,” in *Proc. ICML*, 2019, pp. 6105–6114.
- [17] J. Deng et al., “ImageNet: A large-scale hierarchical image database,” in *Proc. CVPR*, 2009, pp. 248–255. <https://doi.org/10.1109/cvpr.2009.5206848>
- [18] W. Luo et al., “Understanding the effective receptive field in deep convolutional neural networks,” in *Proc. NeurIPS*, 2016, pp. 4898–4906.
- [19] F. Yu and V. Koltun, “Multi-scale context aggregation by dilated convolutions,” in *Proc. ICLR*, 2016.
- [20] M. Raghu et al., “Transfusion: Understanding transfer learning for medical imaging,” in *Proc. NeurIPS*, 2019, pp. 3347–3357.
- [21] K. He et al., “Deep residual learning for image recognition,” in *Proc. CVPR*, 2016, pp. 770–778. <https://doi.org/10.1109/cvpr.2016.90>
- [22] A. Dosovitskiy et al., “An image is worth 16x16 words: Transformers for image recognition at scale,” in *Proc. ICLR*, 2021.
- [23] A. Vaswani et al., “Attention is all you need,” in *Proc. NeurIPS*, 2017, pp. 5998–6008. <https://doi.org/10.65215/r5bs2d54>
- [24] J. Chen et al., “TransUNet: Transformers make strong encoders for medical image segmentation,” arXiv preprint arXiv:2102.04306, 2021.
- [25] Y. Matsoukas et al., “Is it time to replace CNNs with transformers for medical images?” arXiv preprint arXiv:2108.09038, 2021.
- [26] S. Khan et al., “Transformers in vision: A survey,” *ACM Computing Surveys*, vol. 54, no. 10s, pp. 1–41, 2022. <https://doi.org/10.1145/3505244>
- [27] A. Shamshad et al., “Transformers in medical imaging: A survey,” *Medical Image Analysis*, vol. 88, p. 102802, 2023. <https://doi.org/10.1016/j.media.2023.102802>
- [28] Z. Liu et al., “Swin Transformer: Hierarchical vision transformer using shifted windows,” in *Proc. ICCV*, 2021, pp. 10012–10022. <https://doi.org/10.1109/iccv48922.2021.00986>
- [29] H. Touvron et al., “Training data-efficient image transformers & distillation through attention,” in *Proc. ICML*, 2021, pp. 10347–10357.

- [30] K. Han et al., “A survey on vision transformer,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 1, pp. 87–110, 2023. <https://doi.org/10.1109/tpami.2022.3152247>
- [31] S. Ayana et al., “Vision-transformer-based transfer learning for mammogram classification,” *Diagnostics*, vol. 13, no. 2, p. 178, 2023. <https://doi.org/10.3390/diagnostics13020178>
- [32] G. N. Ferdous et al., “LCDEiT: A linear complexity data-efficient image transformer for MRI brain tumor classification,” *IEEE Access*, vol. 11, pp. 20337–20348, 2023. <https://doi.org/10.1109/access.2023.3244228>
- [33] A. Altaf et al., “A novel augmented deep transfer learning for classification of COVID-19 and other thoracic diseases from X-rays,” *Neural Computing and Applications*, vol. 33, pp. 14037–14048, 2021. <https://doi.org/10.1007/s00521-021-06044-0>
- [34] T. G. Dietterich, “Ensemble methods in machine learning,” in *Proc. Multiple Classifier Systems*, 2000, pp. 1–15. [https://doi.org/10.1007/3-540-45014-9\\_1](https://doi.org/10.1007/3-540-45014-9_1)
- [35] J. Kittler et al., “On combining classifiers,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 3, pp. 226–239, 1998. <https://doi.org/10.1109/34.667881>
- [36] S. et al., “Ensemble pre-trained deep convolutional neural network model for classifying medical image datasets,” in *Proc. ICAISS*, 2022, pp. 1–6. <https://doi.org/10.1109/icaiss55157.2022.10011089>

