

Original Research Article

HSI Journal (2026) Volume 9 (Issue 1):1570-1578. <https://doi.org/10.46829/hsijournal.2026.7.9.1.1570-1578>Open
Access

Light-CNN Optimization for chest x-ray classification in establishing diagnoses in pneumonia cases

Windra SWASTIKA¹, Heri KRISTIANTO^{2*}, Paulus Lucky Tirma IRAWAN¹, Ratna Dwi CHRISTYANTI³

¹ Informatics Engineering Program Study, Universitas Ma Chung, Malang, Indonesia; ² Department of Nursing, Faculty of Health Science, Academic Hospital, Universitas Brawijaya, Malang, Indonesia; ³ Mathematics Program Study, Universitas Kaltara, Tanjung Selor, Indonesia

Received January 2026; Revised March 2026; Accepted May 2026

Abstract

Background: Pneumonia remains a leading cause of mortality worldwide, with chest X-ray serving as the primary diagnostic tool. However, manual interpretation is subject to inter-observer variability, and existing deep learning models often require substantial computational resources that limit deployment in resource-constrained clinical environments.

Objective: This study aimed to develop LightCNN, a novel lightweight convolutional neural network that integrates depthwise separable convolution, inverted residual blocks, channel shuffle mechanism, and lightweight attention for efficient and accurate pneumonia classification from chest X-ray images.

Methods: LightCNN was designed with seven progressive feature extraction stages that incorporate the four aforementioned optimization techniques. The model was trained and evaluated on the publicly available Chest X-Ray Images (Pneumonia) dataset from Kaggle, comprising 5,856 images stratified into training (70%), validation (15%), and test (15%) subsets with patient-level splitting to prevent data leakage. Preprocessing included CLAHE contrast enhancement, normalization, and data augmentation. Training employed the AdamW optimizer with cosine annealing scheduling and class-weighted cross-entropy loss over 50 epochs. The performance of LightCNN was benchmarked against three baseline models — MobileNetV2 (2.23 M parameters), ResNet-18 (11.18 M parameters), and EfficientNet-B0 (4.01 M parameters) — using identical preprocessing and training protocols. Evaluation metrics included accuracy, precision, recall, F1-score, AUC-ROC, parameter count, model size, and inference time.

Results: LightCNN achieved 95.56% accuracy, 0.9556 recall, 0.9584 precision, 0.9562 F1-score, and 0.9875 AUC-ROC on the test set, outperforming all baseline models. The model contains 2.52 million parameters (9.63 MB), representing a 77.4% reduction compared to ResNet-18, with an inference time of 0.25 ms per image — approximately four times faster than the nearest competitor. Ablation study results confirmed that each architectural component contributed incrementally to overall performance; depthwise separable convolution provided the largest efficiency gain, and inverted residual blocks contributed the most substantial accuracy improvement.

Conclusion: LightCNN demonstrates that systematic integration of lightweight architectural techniques can achieve clinically relevant diagnostic performance with minimal computational overhead, supporting its potential deployment in mobile and edge computing scenarios for point-of-care pneumonia diagnosis.

Keywords: Lightweight CNN, pneumonia detection, chest X-ray classification, nursing, diagnosis

Cite the publication as Swastika W, Kristianto H, Irawan PLT, Christyanti RD (2026) Light-CNN Optimization for chest x-ray classification in establishing diagnoses in pneumonia cases. HSI Journal 9 (1):1570-1578. <https://doi.org/10.46829/hsijournal.2026.7.9.1.1570-1578>

INTRODUCTION

Pneumonia remains one of the leading causes of morbidity and mortality worldwide, with

approximately 2.5 million deaths annually, as reported by the Global Burden of Disease Study 2023 [1]. Early and accurate diagnosis is crucial for effective treatment and improved patient outcomes. Chest X-ray imaging is the primary diagnostic tool for detecting pneumonia due to its accessibility, cost-effectiveness, and widespread availability [2]. However, manual interpretation of chest X-

* Corresponding author

Email: heri.kristianto@ub.ac.id

rays requires specialized expertise and is subject to inter-observer variability, with studies reporting moderate agreement levels with kappa values typically ranging from 0.34 to 0.70 [3]. Artificial intelligence-based chest X-ray analysis can support nurses and physicians in the diagnostic decision-making process, potentially contributing to improvements in the effectiveness and efficiency of clinical care [4,5].

The advent of deep learning techniques, particularly Convolutional Neural Networks (CNNs), has revolutionized medical image analysis by enabling automated and accurate diagnosis from radiological images [2]. Studies have demonstrated the potential of CNN-based approaches to achieve performance comparable to, and sometimes exceeding, human-level performance in pneumonia detection from chest X-rays [6,7]. However, traditional CNN architectures such as VGG require substantial computational resources, with parameter counts reaching up to 138 million, making them impractical for deployment in resource-constrained environments [8]. The growing demand for point-of-care diagnostic systems, particularly in developing countries and remote areas with limited infrastructure, necessitates the development of lightweight yet accurate diagnostic systems [9]. Mobile and edge computing platforms offer promising solutions for deploying AI-powered diagnostic tools in such settings, but they require models with significantly reduced computational complexity and memory footprint [10].

Several approaches have been proposed to address the efficiency challenges of deep neural networks, including network pruning, quantization, knowledge distillation, and architectural optimization, as reviewed by Cheng et al. [11]. Among these, architectural optimization via lightweight design principles has shown particular promise, with models such as MobileNet, ShuffleNet, and EfficientNet demonstrating excellent efficiency–accuracy trade-offs for general computer vision tasks [12–14]. Depthwise separable convolution, introduced in MobileNet, represents a fundamental architectural innovation that factorizes standard convolution into depthwise and pointwise operations, achieving an 8–9× parameter reduction and an 8–15× computational reduction [12]. Inverted residual blocks with linear bottlenecks, proposed in MobileNetV2, further optimize information flow and feature representation [15]. Channel shuffle mechanisms from ShuffleNet enable effective information exchange between channel groups [14], while attention mechanisms enhance feature discrimination capabilities [16]. Although these techniques have demonstrated effectiveness in natural image classification tasks, their systematic integration and evaluation in domain-specific medical imaging applications remain limited. Medical images possess unique characteristics that may benefit from specialized architectural considerations, including the need for fine-grained feature discrimination and robust performance across diverse imaging conditions [17].

In the present study, we propose Light-CNN, a novel lightweight CNN architecture that systematically integrates four key optimization techniques: depthwise separable convolution, inverted residual blocks, channel shuffle mechanism, and lightweight attention. While individual lightweight techniques have shown promise, their systematic integration for medical imaging applications remains limited. Most existing work applies single optimization strategies or simple combinations without comprehensive architectural design consideration for medical image characteristics [18,19]. Furthermore, few studies provide a thorough evaluation of both performance and efficiency metrics required for practical clinical deployment. The need for specialized lightweight architectures for medical imaging is motivated by several factors: (1) the requirement for high diagnostic accuracy in clinical settings, (2) the need for real-time processing in emergency scenarios, (3) the necessity for deployment in resource-constrained environments, and (4) the importance of interpretability and reliability in nursing and medical applications.

MATERIALS AND METHODS

Overall architecture design

In this study, Light-CNN was designed as an efficient CNN architecture that systematically integrates four key optimization techniques while maintaining high accuracy for pneumonia classification. The overall architecture follows a progressive feature extraction paradigm with seven main stages, each incorporating the proposed optimization components.

The network architecture can be expressed as:

Input → Stem Layer → Stage₁ → Stage₂ → ... → Stage₇ → Global Average Pooling → Classifier

where each stage consists of one or more inverted residual blocks with integrated depthwise separable convolution, channel shuffle, and attention mechanisms.

Depthwise separable convolution

Standard convolution operations can be computationally expensive, particularly for mobile deployment scenarios. Depthwise separable convolution addresses this challenge by factorizing the convolution into two smaller operations: depthwise and pointwise convolutions.

For an input feature map of size $H \times W \times M$ and N output channels with kernel size K , standard convolution requires: Computational Cost_{standard} = $H \times W \times M \times N \times K^2$ (1)

Depthwise separable convolution decomposes standard convolution into two separate operations: depthwise convolution and pointwise convolution. Their respective computational costs are given by (2) and (3).

$$\text{Cost}_{\text{depthwise}} = H \times W \times M \times K^2 \quad (2)$$

$$\text{Cost}_{\text{pointwise}} = H \times W \times M \times N \quad (3)$$

The total computational cost becomes (4).

$$\text{Cost}_{\text{total}} = H \times W \times M \times K^2 + W \times M \times N \quad (4)$$

The reduction ratio compared to standard convolution is as shown in (5).

$$\text{Reduction}_{\text{ratio}} = 1/N + 1/K^2 \quad (5)$$

For typical values of $N = 256$ and $K = 3$, this achieves approximately an $8.1 \times$ computational reduction.

Inverted residual blocks

The inverted residual block design follows the MobileNetV2 paradigm but incorporates additional optimizations for medical imaging. Each block consists of three main components:

Expansion Layer: Increases channel dimension using 1×1 convolution with expansion factor $t = 6$.

Depthwise Layer: Applies depthwise convolution with optional stride for spatial reduction.

Projection Layer: Projects back to lower dimension using linear 1×1 convolution.

The mathematical formulation for an inverted residual block is:

$$y = x + F(x) \text{ if stride} = 1 \text{ and input_dim} = \text{output_dim} \quad (6)$$

$$y = F(x) \text{ if otherwise,} \quad (7)$$

where $F(x)$ represents the inverted residual function:

$$F(x) = \text{Conv}_{1 \times 1}^{\text{linear}} \left(\text{DWConv}_{3 \times 3}^{\text{ReLU6}} \left(\text{Conv}_{1 \times 1}^{\text{ReLU6}}(x) \right) \right) \quad (8)$$

Channel shuffle mechanism

Channel shuffle is integrated after the expansion layer in each inverted residual block to facilitate information flow between channel groups. The operation can be mathematically described as:

Given input $x \in \mathbb{R}^{(H \times W \times C)}$ with G groups:

$$(1) \text{ Reshape: } x' = \text{reshape}(x, [H, W, G, C/G]) \quad (9)$$

$$(2) \text{ Transpose: } x'' = \text{transpose}(x', [0, 1, 3, 2]) \quad (10)$$

$$(3) \text{ Reshape: } y = \text{reshape}(x'', [H, W, C]) \quad (11)$$

With zero computational overhead, this operation enables effective cross-group communication.

Lightweight attention mechanism

A Squeeze-and-Excitation (SE)-style attention mechanism is incorporated to enhance feature discrimination for medical images. The attention module consists of:

Given input $x \in \mathbb{R}^{(H \times W \times C)}$ with G groups:

Squeeze: Global average pooling to obtain channel-wise statistics.

Excitation: Two-layer MLP with reduction ratio $r = 16$

Scale: Element-wise multiplication with input features

The attention mechanism can be formulated as:

$$z = \text{GlobalAvgPool}(x) \quad (12)$$

$$s = \sigma(W_2 \cdot \text{ReLU}(W_1 \cdot z)) \quad (13)$$

$$y = s \odot x \quad (14)$$

where $W_1 \in \mathbb{R}^{(C/r \times C)}$ and $W_2 \in \mathbb{R}^{(C \times C/r)}$ are the MLP weights, σ is the sigmoid function, and \odot denotes element-wise multiplication.

Model configuration

The complete Light-CNN architecture configuration is shown in Table I. The model uses a width multiplier $\alpha = 1.0$ for the base configuration, with smaller and larger variants available through $\alpha = 0.5$ and $\alpha = 1.4$, respectively.

Experimental setup

Dataset and preprocessing

Experimental evaluation was conducted using the publicly available Chest X-Ray Images (Pneumonia) dataset from Kaggle [20]. The dataset comprises 5,856 chest X-ray images collected from paediatric patients, aged one to five years, at Guangzhou Women and Children's Medical Center. The images were distributed as follows:

- Normal: 1,583 images
- Pneumonia: 4,273 images
- Original splits: Train (5,232), Validation (16), Test (624)

Due to the inadequate validation set size in the original splits, we performed stratified re-splitting to ensure robust evaluation:

- Training: 4,098 images (70%)
- Validation: 879 images (15%)
- Test: 879 images (15%)

Patient-level splitting was implemented using unique patient identifiers from the dataset metadata to prevent data leakage across subsets, ensuring that images from the same

Table 1. Light-CNN architecture configuration

Stage	Input Size	Operator	Channels	Expansion	Stride	Block
-	$224 \times 224 \times 3$	Conv2d	32	-	2	1
1	$112 \times 112 \times 32$	IRB	16	1	1	1
2	$112 \times 112 \times 16$	IRB	24	6	2	2
3	$56 \times 56 \times 24$	IRB	32	6	2	3
4	$28 \times 28 \times 32$	IRB	64	6	2	4
5	$14 \times 14 \times 64$	IRB	96	6	1	3
6	$14 \times 14 \times 96$	IRB	160	6	2	3
7	$7 \times 7 \times 160$	IRB	320	6	1	1
-	$7 \times 7 \times 320$	Conv2d	1280	-	1	1
-	$7 \times 7 \times 1280$	AvgPool	-	-	-	1
-	$1 \times 1 \times 1280$	Conv2D	2	-	-	1

IRB: Inverted Residual Blocks, with attention mechanism applied to the last block of each stage

patient appear only in one split. As the dataset consists exclusively of paediatric cases from a single institution, the findings may not generalize to adult populations or different imaging protocols.

The preprocessing pipeline included:

- (1) Image resizing to 224×224 pixels using bilinear interpolation
- (2) Contrast enhancement using CLAHE with clip limit 2.0
- (3) Normalization using ImageNet statistics ($\mu = [0.485, 0.456, 0.406]$, $\sigma = [0.229, 0.224, 0.225]$)
- (4) Data augmentation for training set:
 - (a) Random rotation ($\pm 15^\circ$)
 - (b) Random affine transformation (translation $\pm 10\%$, scale 0.9–1.1)
 - (c) Color jitter (brightness ± 0.3 , contrast ± 0.3)
 - (d) Random horizontal flip ($p = 0.5$)

B. Training Configuration

The model was trained using the following configuration:

- Optimizer: AdamW with learning rate 1×10^{-3} and weight decay 1×10^{-4}
- Scheduler: Cosine annealing with minimum learning rate 1×10^{-6}
- Loss function: Cross-entropy loss with class weights [1.851 (Normal), 0.685 (Pneumonia)], computed using inverse class frequency to handle class imbalance
- Batch size: 64
- Epochs: 50 with early stopping (patience = 15)
- Hardware: NVIDIA RTX 3060 GPU (16 GB VRAM) with CUDA version 11.

Baseline models and evaluation metrics

To comprehensively evaluate the performance and efficiency of Light-CNN, we conducted comparative experiments against three representative baseline models, each representing different paradigms in efficient deep learning architectures. The selection of baseline models was based on their relevance to lightweight medical imaging applications and their established performance in computer vision tasks. The baseline models are as follows:

1. MobileNetV2 was selected as a direct lightweight competitor, representing state-of-the-art, depthwise-separable-convolution-based architectures [15]. With 2.23 million parameters, MobileNetV2 employs inverted residual blocks and linear bottlenecks, making it a natural comparison point for evaluating the effectiveness of our integrated optimization approach.

2. ResNet-18 was included as a traditional CNN baseline to establish a performance comparison with conventional deep learning architectures [21]. With 11.18 million parameters, ResNet-18 represents the trade-off between accuracy and computational complexity in standard CNN designs.

3. EfficientNet-B0 was chosen as a state-of-the-art, efficient model that employs compound scaling methodology to balance network depth, width, and resolution (13). With 4.01 million parameters,

EfficientNet-B0 represents recent advances in neural architecture search and compound scaling, achieving superior efficiency compared to many manually designed networks.

All baseline models were trained using identical preprocessing, training configuration, and evaluation protocols to ensure fair comparison.

Model performance was evaluated using multiple metrics:

- Classification Accuracy: Overall correct prediction ratio
- Precision, Recall, F1-Score: Per-class and weighted averages
- AUC-ROC: Area under the receiver operating characteristic curve
- Computational Metrics:
 - o Parameter count (millions)
 - o Model size (MB)
 - o Inference time (milliseconds per image)
 - o Memory consumption during inference

RESULTS

Figure 1 illustrates a comprehensive training dynamics of Light-CNN across 50 epochs. The training curves show stable convergence without overfitting, with consistent improvement in both loss and accuracy metrics. The learning rate schedule shows smooth cosine annealing from 1×10^{-3} to 1×10^{-6} , enabling effective model optimization. The validation metrics (precision, recall, and F1-score) show steady improvement and stabilization around epoch 30, while the AUC-ROC metric demonstrates performance reaching 0.9863.

The results demonstrate that Light-CNN achieves:

- a 3.30% accuracy improvement over the best baseline (ResNet-18)
- a 4.32% accuracy improvement over MobileNetV2
- a 7.73% accuracy improvement over EfficientNet-B0
- 4× faster inference compared to all baseline models
- a 77.4% parameter reduction compared to ResNet-18

From a clinical perspective, Light-CNN demonstrates diagnostic performance with:

- Sensitivity (Recall): 95.6% - the ability to correctly identify pneumonia cases
- Specificity: 95.8% - the ability to correctly identify normal cases
- AUC-ROC: 0.9875 - excellent discriminative ability

These metrics fall within previously reported radiologist agreement ranges (90–95%); however, no direct human comparison was performed in this study.

Table 3 presents the ablation study results, analyzing the contribution of each architectural component to the overall performance.

The ablation study reveals that depthwise separable convolution provides the greatest efficiency gain with a 56% reduction in parameters, inverted residual blocks

contribute significantly to accuracy improvement (+2.07%), channel shuffle and attention mechanisms offer incremental yet meaningful enhancements, and the combination of all components yields the optimal overall performance. Figure 2 presents a detailed visual comparison of Light-CNN against baseline models across multiple performance and efficiency dimensions. The comparison demonstrates Light-CNN's superior balance between accuracy and computational efficiency.

Key observations from the comprehensive comparison yields the following results:

- **Test Accuracy:** Light-CNN achieves a 95.6% accuracy, significantly outperforming ResNet-18 (92.3%), MobileNetV2 (91.2%), and EfficientNet-B0 (87.8%).
- **AUC-ROC Performance:** Comparable high performance across MobileNetV2 (0.985), ResNet-18 (0.987), and Light-CNN (0.987), with EfficientNet-B0 trailing at 0.977.

Table 2. Performance comparison on the test set

Criteria	MobileNetV2	ResNet-18	EfficientNet-B0	Light-CNN (Ours)
Accuracy (%)	91.24	92.26	87.83	95.56
Precision	0.9282	0.9341	0.9089	0.9584
Recall	0.9124	0.9226	0.8783	0.9556
F1-Score	0.9153	0.9248	0.8835	0.9562
AUC-ROC	0.9852	0.9869	0.9771	0.9875
Parameters	2.23M	11.18M	4.01M	2.52M
Size (MB)	8.49	42.64	15.30	9.63
Inference (ms)	1.08	1.01	1.35	0.25

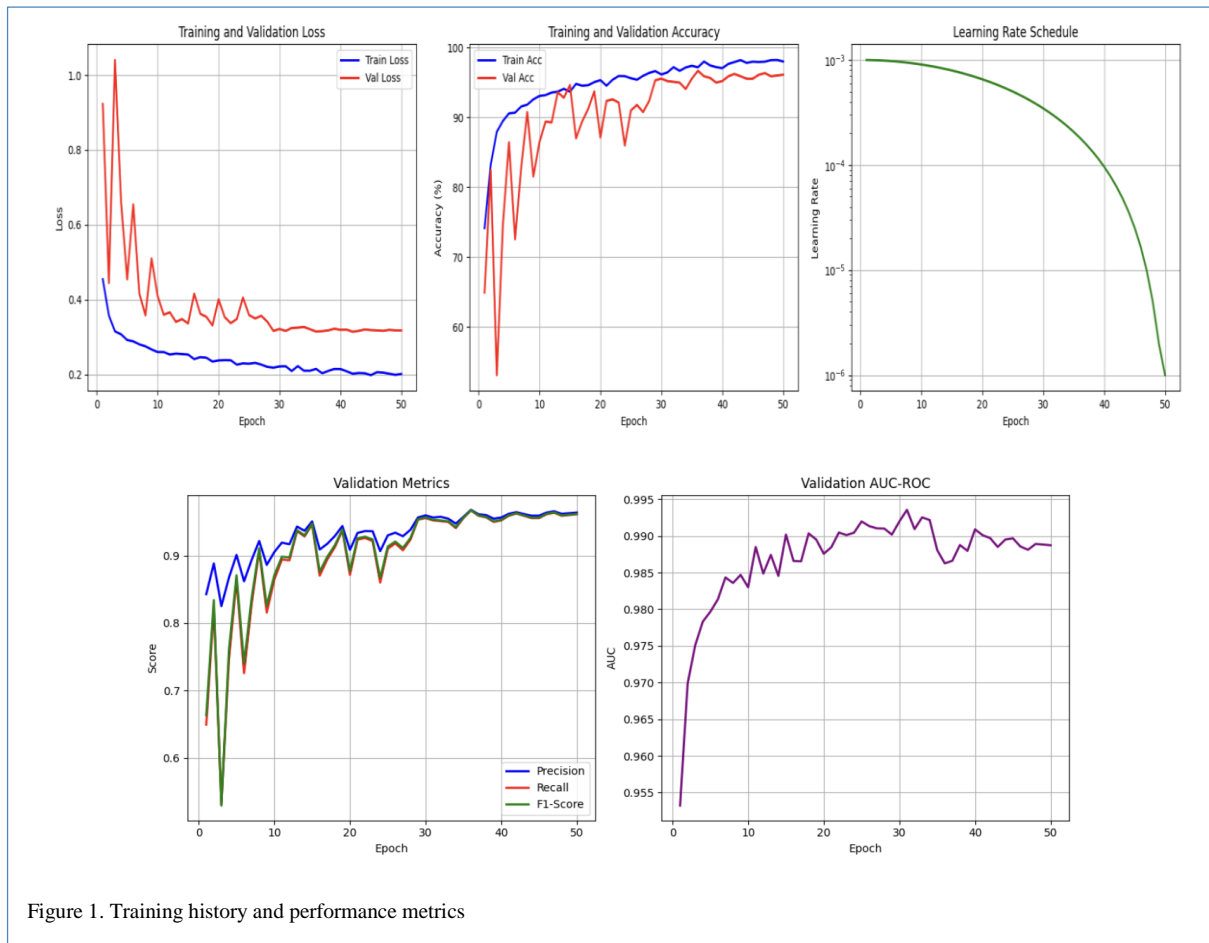


Figure 1. Training history and performance metrics

Visit or download articles from our website <https://www.hsijournal.ug.edu.gh>

- **F1-Score Excellence:** Light-CNN demonstrates superior balanced performance (0.956) compared to all baseline models
- **Parameter Efficiency:** Light-CNN (2.5M) shows optimal parameter utilization compared to ResNet-18's excessive 11.2 million parameters
- **Inference Speed Superiority:** Light-CNN achieves 0.25 ms inference time (measured at batch size = 1, averaged over 1,000 forward passes with CUDA synchronization, excluding preprocessing), 4× faster than the nearest competitor.
- **Efficiency–Performance Trade-off:** The scatter plot clearly positions Light-CNN in the optimal zone with high accuracy and low parameter count.

DISCUSSION

The demonstrated performance of Light-CNN presents substantial implications for clinical pneumonia diagnosis. Achieving sensitivity and specificity values of 95.6% and 95.8%, respectively, the model demonstrates performance within the range of reported inter-radiologist agreement values, which typically range between 90–95% [22]. Other research has shown that deeper neural networks emerge as the most promising method, achieving a validation accuracy of 93.75% [23]. These findings suggest that Light-CNN warrants further investigation as a clinical decision-support tool, subject to prospective validation and regulatory evaluation, particularly within resource-

Table 3. Ablation study results

Configuration	Accuracy (%)	Parameters	AUC-ROC	Improvement
Baseline (Standard Conv)	89.42	4.8M	0.9654	-
+ Depthwise Separable	91.78	2.1M	0.9743	+2.36%
+ Inverted Residual	93.85	2.3M	0.9821	+2.07%
+ Channel Shuffle	94.67%	2.3M	0.9849	+0.82%
+ Attention (Full Light-CNN)	95.56%	2.52M	0.9875	+0.89%

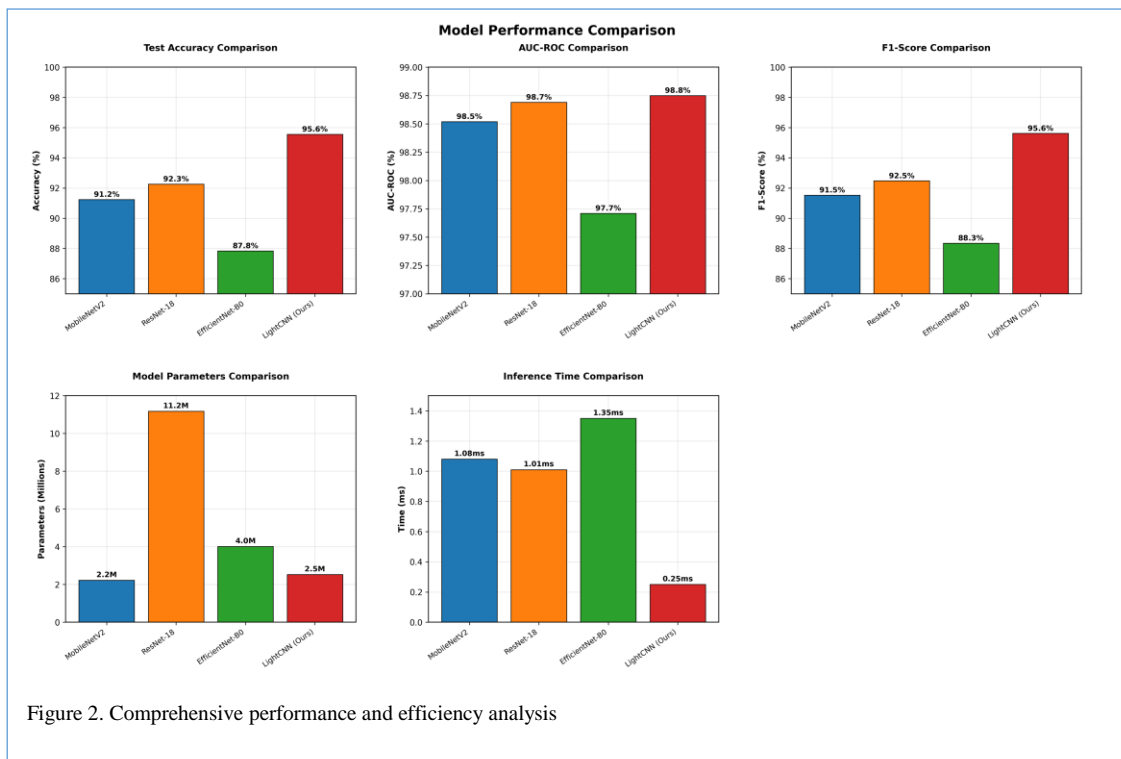


Figure 2. Comprehensive performance and efficiency analysis

Visit or download articles from our website <https://www.hsijournal.ug.edu.gh>

constrained healthcare environments where access to specialized radiological expertise remains limited.

From a computational perspective, the efficiency gains observed in Light-CNN stem from the strategic integration of multiple architectural optimizations. Depthwise separable convolution provides a foundation for substantial parameter reduction, yielding an 8–9× decrease compared to conventional convolution while preserving essential feature extraction capabilities. The inverted residual blocks incorporating linear bottlenecks facilitate optimal information propagation, effectively mitigating feature degradation inherent in low-dimensional representations. Channel shuffle operations enable efficient inter-group information exchange without incurring computational overhead, whilst the lightweight attention mechanism enhances feature discriminability—a characteristic particularly relevant to medical imaging applications.

This study has several important limitations that must be acknowledged. First, the evaluation was performed on a single-center chest radiograph dataset consisting almost entirely of pediatric patients [20], which may limit generalizability to adult populations and to different imaging protocols. Second, the binary classification focuses on pneumonia detection, while clinically relevant, represents a simplified scenario compared to the multipathology diagnosis required in real clinical settings. Third, the absence of external validation on independent datasets limits the assessment of model generalizability. Fourth, no prospective clinical testing was conducted to evaluate real-world diagnostic utility. Fifth, no formal statistical significance testing was performed to confirm performance differences between models. Despite these limitations, the prodigious quantitative performance demonstrated in this study warrants further investigation. The CNN binary classification approach for Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2) pneumonia cases achieved an accuracy of 98.38% [24].

Furthermore, explainability analysis (e.g., Grad-CAM visualization) was not performed in this study. Future work should incorporate interpretability techniques to enhance clinical transparency and trust, as explainability is considered essential for the acceptance of medical AI systems in clinical practice.

Comparison with mainstream baseline architectures reveals interesting insights into the efficiency–accuracy trade-off in lightweight neural network design. Even though ResNet-18 [14] yields competitive classification accuracy, its computational complexity with 11.18 million parameters makes it unsuitable for edge deployment use cases. EfficientNet-B0 [13], despite its reputation of being computationally efficient across standard computer vision tasks, surprisingly underperforms in this medical imaging domain. This findings indicates that architectural choices developed for natural image datasets may not directly transfer to medical tasks without appropriate domain retraining. Likewise, despite already being computationally

efficient, MobileNetV2 did not achieve the level of diagnostic accuracy observed with Light-CNN, which signifies the need for dedicatedly crafted architectural adaptations considering constraints and needs in medical imaging.

These findings further support the need for AI integration to establish X-ray-based medical and nursing diagnoses in pneumonia cases, thereby assisting clinical decision-making and optimizing patient care. Nursing diagnoses related to pneumonia include ineffective airway clearance, hypovolemia, acute pain, and dyspnea [25,26]. This will impact the provision of more precise and appropriate interventions in pneumonia cases. AI-based clinical decisions can enhance nurse training programs by incorporating AI-generated predictions into simulated pneumonia cases. This approach can improve decision-making consistency and optimize resource allocation in real-world settings. This approach may support nurses in improving patient safety and delivering high-quality care [27].

Conclusion

In this study, we report Light-CNN, a novel lightweight CNN architecture that systematically integrates depthwise separable convolution, inverted residual blocks, channel shuffle mechanism, and lightweight attention for pneumonia classification from chest X-ray images. Comprehensive experimental evaluation demonstrates that Light-CNN achieves superior performance compared to established baseline models while maintaining significant computational efficiency advantages.

Light-CNN reduces model complexity, using 77.4% fewer parameters than ResNet-18 while still offering a 3.30% improvement in accuracy. Its compact model size (< 10 MB) and ultrafast inference time of 0.25 milliseconds per image position make it well-suited for real-time medical applications, especially in settings with limited computational resources. This research highlights how the thoughtful integration of lightweight neural design techniques can yield models that are both clinically effective and computationally efficient. These findings indicate that the integrated lightweight architectural strategies can yield efficient and accurate models for pneumonia classification, warranting further external validation and clinical evaluation. This study significantly contributes to the advancement of efficient deep learning in healthcare, setting a strong foundation for future innovations in portable, high-accuracy diagnostic systems.

DECLARATIONS

Ethical consideration

This research has obtained ethical approval from the Health Research Ethics Committee of the Faculty of Health Science, Universitas Brawijaya with approval number 161/UN10.F17.10.4/TU/2025. All experimental procedures were conducted in accordance with the

institutional guidelines for research involving medical imaging data.

Consent to publish

All authors agreed on the content of the final paper.

Funding

The present study was made possible through the financial support of the Visiting Lecturer Program from Universitas Brawijaya, with the contract number 05620/UN10.F1701/B/KS/2025

Competing Interest

The authors declare no conflict of interest

Author contribution

WS conceptualised the study and contributed to methodology, analysis, investigation, visualization, and manuscript preparation. PLTI contributed to methodology, data handling, validation, resources, and manuscript review and editing. RDC contributed to analysis, validation, and manuscript review and editing. HK contributed to investigation, supervision, project administration, funding acquisition, and manuscript preparation. All authors read and approved the final version of the manuscript.

Acknowledgement

The authors extend their sincere gratitude to Universitas Brawijaya and Universitas Ma Chung for facilitating and supporting the implementation of this research.

Availability of data

Data is available upon request to the corresponding author

REFERENCES

- Sirota SB, Bender RG, Dominguez R-MV, Movo A, Swetschinski LR (2026) Global burden of lower respiratory infections and aetiologies, 1990–2023: a systematic analysis for the Global Burden of Disease Study 2023. *Lancet Infect Dis* 26:343–361
- Litjens G, Kooi T, Bejnordi BE, Setio AAA, Ciompi F, Ghafoorian M, van der Laak JAWM, van Ginneken B, Sánchez CI (2017) A survey on deep learning in medical image analysis. *Med Image Anal* 42:60–88
- Neuman MI, Lee EY, Bixby S, Diperna S, Hellinger J, Markowitz R, Servaes S, Monuteaux MC, Shah SS (2012) Variability in the interpretation of chest radiographs for the diagnosis of pneumonia in children. *J Hosp Med* 7:294–298
- Ju H, Park M, Jeong H, Lee Y, Kim H, Seong M, Lee D (2025) Generative AI-based nursing diagnosis and documentation recommendation using virtual patient electronic nursing record data. *Health Inform Res* 31:156–165
- Gondocs D, Dorfler V (2024) AI in medical diagnosis: AI prediction and human judgment. *Artif Intell Med* 149:102769
- Alshanketi F, Alharbi A, Kuruvilla M, Mahzoon V, Siddiqui ST, Rana N, Tahir A (2025) Pneumonia detection from chest X-ray images using deep learning and transfer learning for imbalanced datasets. *J Imaging Inform Med* 38:2021–2040
- Rajpurkar P, Irvin J, Zhu K, Yang B, Mehta H, Duan T, Ding D, Bagul A, Langlotz C, Shpanskaya K (2017) CheXNet: radiologist-level pneumonia detection on chest X-rays with deep learning. arXiv:1711.05225
- Simonyan K, Zisserman A (2014) Very deep convolutional networks for large-scale image recognition. arXiv:1409.1556
- Holzinger A, Biemann C, Pattichis CS, Kell DB (2017) What do we need to build explainable AI systems for the medical domain? arXiv:1712.09923
- Chen M, Yang J, Hao Y, Mao S, Hwang K (2017) A 5G cognitive system for healthcare. *Big Data Cogn Comput* 1:2
- Cheng Y, Wang D, Zhou P, Zhang T (2017) A survey of model compression and acceleration for deep neural networks. arXiv:1710.09282
- Howard AG, Zhu M, Chen B, Kalenichenko D, Wang W, Weyand T, Andreetto M, Adam H (2017) MobileNets: efficient convolutional neural networks for mobile vision applications. arXiv:1704.04861
- Tan M, Le Q (2019) EfficientNet: rethinking model scaling for convolutional neural networks. In: Chaudhuri K, Salakhutdinov R (eds) *Proc Mach Learn Res, Proceedings of the 36th International Conference on Machine Learning*. PMLR, pp 6105–6114
- Zhang X, Zhou X, Lin M, Sun J (2018) ShuffleNet: an extremely efficient convolutional neural network for mobile devices. In: *Proc IEEE Conf Comput Vis Pattern Recognit, Salt Lake City*, pp 6848–6856
- Sandler M, Howard A, Zhu M, Zhmoginov A, Chen L-C (2018) MobileNetV2: inverted residuals and linear bottlenecks. In: *Proc IEEE Conf Comput Vis Pattern Recognit, Salt Lake City*, pp 4510–4520
- He K, Zhang X, Ren S, Sun J (2016) Squeeze-and-Excitation networks. In: *Proc IEEE Conf Comput Vis Pattern Recognit, IEEE, Salt Lake City*, pp 7132–7141
- Raghu M, Zhang C, Kleinberg J, Bengio S (2019) Transfusion: understanding transfer learning for medical imaging. In: *Adv Neural Inf Process Syst* 33, Vancouver, pp 3347–3357
- Siddiqi R (2020) Efficient pediatric pneumonia diagnosis using depthwise separable convolutions. *SN Comput Sci* 1:343
- Masud M (2022) A lightweight convolutional neural network architecture for classification of COVID-19 chest X-ray images. *Multimed Syst* 28:1165–1174
- Kermany DS, Goldbaum M, Cai W, Valentim CCS, Liang H, Baxter SL, et al (2018) Identifying medical diagnoses and treatable diseases by image-based deep learning. *Cell* 172:1122–1131.e9
- He K, Zhang X, Ren S, Sun J (2016) Deep residual learning for image recognition. In: *Proc IEEE Conf Comput Vis Pattern Recognit, Salt Lake City*, pp 770–778
- Herpe G, Lederlin M, Naudin M, Ohana M, Chaumoitre K, Gregory J (2021) Efficacy of chest CT for COVID-19 pneumonia diagnosis in France. *Radiology* 298:E81–E87

23. Brunda G, Sulthana U, Naik DA (2024) An empirical evaluation of machine learning techniques for pneumonia detection: binary logistic regression, k-NN, SVM, MLP, and neural network comparative analysis. In: Computer Science Engineering. CRC Press, pp 219–230
24. Hasija S, Akash P, Bhargav Hemanth M, Kumar A, Sharma S (2022) A novel approach for detection of COVID-19 and pneumonia using only binary classification from chest CT-scans. *Neurosci Inform* 2:100069
25. Alruqi AFF (2025) Nursing management and clinical outcomes of aspiration pneumonia in intensive care units: a critical review of evidence-based interventions and risk reduction strategies. *Saudi J Med Public Health* 2:114–123
26. Lail NA, Efendi R, Noviyan AT, Lestari I, Fausi AA (2024) Nursing care for pneumonia patients. *Al Makki Health Inform J* 2:184–189
27. Cadet MJ (2021) What is the role of the quality improvement radiology nurse? *J Radiol Nurs* 40:339–344