

Multimodal Deep Learning for Bone Fracture Detection Using Radiographic Images and Clinical Metadata

Gujjula Swarnalatha¹ and Dr. Ranga Swamy Sirisati²

¹ Research Scholar, Department of CSE, Bharatiya Engineering Science and Technology Innovation University, Gownivaripalli, Gorantla, Andhra Pradesh, India

Email: 2022wpcse022@bestiu.edu.in

² Associate Professor, Department of CSE, Vignan's Institute of Management and Technology for Women, Ghatkesar, Medchal, Telangana, India

Email: sirisatiranga@gmail.com

Abstract

Accurate and automated detection of bone fractures from radiographic images is a clinically significant challenge, as delays or errors in diagnosis can lead to adverse patient outcomes. This paper presents a novel multimodal deep learning framework that jointly leverages grayscale X-ray images and structured clinical metadata to improve binary fracture classification (fracture vs. non-fracture). The image branch employs a custom convolutional neural network (CNN) to extract hierarchical spatial features, while the clinical branch utilizes a multilayer perceptron (MLP) to encode patient-level metadata. The two representation streams are integrated through feature-level concatenation and passed through a fully connected classifier. The proposed model was trained and evaluated on the BoneFractureYolo8 dataset, achieving a classification accuracy of 97.30% and a binary cross-entropy loss of 0.0829. Evaluation using ROC-AUC curves, confusion matrices, and prediction score distributions confirms the model's strong discriminative capability. Although synthetic clinical features were employed in this study, the results demonstrate that multimodal architectures can substantially enhance medical image classification. The findings motivate further investigation into multimodal data integration for real-world clinical deployment.

Keywords: Bone Fracture Detection; Deep Learning; Multimodal Learning; Convolutional Neural Network; Clinical Metadata; Medical Image Analysis

How to cite this article: Swarnalatha G, Sirisati RS. Multimodal Deep Learning for Bone Fracture Detection Using Radiographic Images and Clinical Metadata. *Int J Drug Deliv Technol.* 2026;16(34s): 1057-1065. DOI: 10.25258/ijddt.16.34s.132

Source of support: Nil.

Conflict of interest: None

1. Introduction

Bone fractures represent one of the most prevalent musculoskeletal injuries worldwide, accounting for millions of emergency department visits annually. Prompt and precise diagnosis is essential, as undetected or misclassified fractures may result in delayed treatment, chronic pain, and long-term functional impairment [1]. Conventional diagnostic practice relies on radiologists manually interpreting plain radiographs (X-rays), a process that is labor-intensive, subject to inter-observer variability, and particularly challenging in high-throughput emergency settings [2].

The emergence of deep learning, particularly convolutional neural networks (CNNs), has transformed medical image analysis. CNN-based systems have achieved radiologist-level performance in tasks such as diabetic retinopathy screening, pneumonia detection, and skin lesion classification [3]. In the orthopedic domain, CNNs have been applied to detect various fracture types from radiographic images with considerable success [4]. However, a fundamental limitation of image-only models is their neglect of clinical context. Patient demographic characteristics, mechanism of injury, pain scores, and other metadata are routinely recorded in clinical practice and can significantly guide diagnostic reasoning.

Multimodal learning—the integration of heterogeneous data modalities within a unified model—has emerged as a promising paradigm to bridge this gap. By combining visual and tabular information, multimodal models can leverage complementary signals unavailable to single-modality systems [5]. Despite these advantages, the application of multimodal deep learning in musculoskeletal radiology remains relatively underexplored.

This paper makes the following contributions: (1) we propose a dual-branch multimodal deep learning architecture that fuses CNN-derived image embeddings with MLP-encoded clinical features; (2) we evaluate the

Multimodal Deep Learning for Bone Fracture Detection Using Radiographic Images and Clinical Metadata

framework on a publicly available YOLO-formatted bone fracture dataset; (3) we provide extensive quantitative and visual analysis of model performance; and (4) we discuss the limitations and future roadmap toward real-world clinical adoption.

2. Related Work

2.1 CNN-Based Fracture Detection

Significant progress has been made in applying deep CNNs to bone fracture detection. Lindsey et al. [6] demonstrated that a CNN-based wrist fracture detector improved diagnostic accuracy among emergency physicians. Rayan et al. [7] evaluated transfer learning approaches using ResNet and DenseNet for detecting distal radius fractures, reporting AUC scores exceeding 0.93. Similarly, EfficientNet-based classifiers have been applied to multi-class fracture categorization in the FracAtlas dataset, achieving competitive performance with fewer parameters [8]. Despite these advances, these studies predominantly employ single-modality (image-only) pipelines.

2.2 Multimodal Learning in Medical Imaging

Multimodal fusion has been extensively explored in oncology and cardiology. For example, Cheerla and Gevaert [9] combined histopathology images with genomic profiles for cancer survival prediction. In radiology, Acosta et al. [10] demonstrated that integrating clinical notes with imaging improved diagnostic accuracy for chest pathologies. Crucially, the design of the fusion mechanism—early (feature-level), late (decision-level), or intermediate—has been shown to substantially affect model performance [11]. Early fusion, as employed in this work, enables the model to learn cross-modal correlations during training.

2.3 Research Gap

While multimodal learning has demonstrated success in other medical domains, its application to bone fracture detection from radiographic images combined with structured clinical metadata remains limited. Most existing works in musculoskeletal AI focus on unimodal imaging approaches. This study addresses this gap by proposing an early-fusion multimodal framework specifically designed for bone fracture binary classification.

3. Methodology

3.1 Dataset Description

The BoneFractureYolo8 dataset, sourced from a publicly available repository, contains radiographic images annotated in YOLO format. Each image is accompanied by a corresponding label file that encodes bounding box coordinates and fracture presence information. Binary labels were derived as follows: images with at least one annotated region (i.e., at least one non-empty label file) were assigned a positive label (fracture present, $y = 1$); images with empty label files were assigned a negative label (no fracture, $y = 0$).

3.2 Data Preprocessing

All radiographic images were converted to grayscale and resized to a uniform spatial resolution of 224×224 pixels to match the expected input dimension of the CNN branch. Pixel values were normalized to the range $[0, 1]$ by dividing by 255. No additional augmentation was applied in the current study; this is acknowledged as a limitation and addressed in Section 6.

3.3 Clinical Feature Representation

In the absence of real electronic health record (EHR) data accompanying the dataset, synthetic clinical feature vectors of dimensionality $d = 10$ were generated by sampling uniform random values from $[0, 1]$ for each sample. Although this approach does not reflect true clinical distributions, it serves as a proof-of-concept to validate the multimodal fusion architecture. Future work will substitute these with real clinical variables (see Section 6).

3.4 Problem Formulation

Let the multimodal dataset be defined as:

$$D = \{(X_i, C_i, y_i)\}_{i=1}^N$$

where $X_i \in \mathbb{R}^{H \times N \times 1}$ denotes the grayscale radiographic image for sample i , $C_i \in \mathbb{R}^d$ is the corresponding clinical feature vector, and $y_i \in \{0, 1\}$ is the binary label. The learning objective is to find the optimal parameter set θ^* such that:

$$\theta^* = \underset{\theta}{\operatorname{arg\,min}} L(f(X_i, C_i; \theta), y_i)$$

where $f(\cdot)$ is the multimodal classifier and $L(\cdot)$ is the binary cross-entropy loss.

3.5 Image Feature Extraction (CNN Branch)

Multimodal Deep Learning for Bone Fracture Detection Using Radiographic Images and Clinical Metadata

The image processing branch employs a custom CNN architecture consisting of two convolutional blocks, each comprising a Conv2D layer with ReLU activation followed by 2×2 max-pooling. Let Z_l denote the feature map at layer l . The convolution operation is expressed as:

$$Z_l = \sigma(W_l * Z_{l-1} + b_l)$$

where $*$ denotes the convolution operator, σ is the ReLU activation function ($\max(0, \cdot)$), W_l is the learnable filter bank, and b_l is the bias term. Max-pooling reduces spatial dimensions while providing translation invariance. The output of the final convolutional block is flattened and projected to a 64-dimensional embedding F^I through a fully connected layer with dropout regularization.

3.6 Clinical Feature Encoding (MLP Branch)

The clinical branch encodes the input vector $C \in \mathbb{R}^{10}$ through a two-layer MLP. Each hidden layer computes:

$$h_l = \sigma(W_l h_{l-1} + b_l)$$

The final clinical embedding $F_C \in \mathbb{R}^{32}$ is obtained after applying dropout regularization to the penultimate layer output.

3.7 Multimodal Fusion

Feature-level (early) fusion is employed, in which the image embedding F^I and the clinical embedding F_C are concatenated along the feature dimension:

$$F_{fusion} = [F^I \oplus F_C] \in \mathbb{R}^{96}$$

Concatenation preserves the full information content of both modalities and allows the downstream classifier to learn cross-modal interaction patterns.

3.8 Classification Head

The fused representation is passed through two fully connected layers with ReLU activation and dropout, followed by a final output unit with sigmoid activation:

$$\hat{y} = \sigma(W_f F_{fusion} + b_f), \quad \sigma(z) = 1 / (1 + e^{-z})$$

The output $\hat{y} \in (0, 1)$ represents the predicted probability of fracture presence.

3.9 Loss Function and Optimization

The model is trained using binary cross-entropy loss:

$$L = -(1/N) \sum [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

Optimization is performed using the Adam optimizer with an initial learning rate of $\alpha = 0.001$, utilizing adaptive first and second moment estimates (m_t, v_t) with $\epsilon = 10^{-8}$ for numerical stability. A dropout rate of $p = 0.5$ was applied at all dropout layers to mitigate overfitting.

3.10 Model Architecture Summary

Table 1 summarizes the complete architecture of the proposed multimodal network, including layer types, output shapes, and parameter counts.

Layer (Type)	Output Shape	Param #	Connected To
image_input (InputLayer)	(None, 224, 224, 1)	0	—
conv2d_1 (Conv2D, 32 filters)	(None, 222, 222, 32)	320	image_input
max_pooling2d_1 (MaxPooling2D)	(None, 111, 111, 32)	0	conv2d_1
conv2d_2 (Conv2D, 64 filters)	(None, 109, 109, 64)	18,496	max_pooling2d_1
max_pooling2d_2 (MaxPooling2D)	(None, 54, 54, 64)	0	conv2d_2
flatten (Flatten)	(None, 186,624)	0	max_pooling2d_2
dense_img (Dense, 128)	(None, 128)	23,888,000	flatten
dropout_img (Dropout, p=0.5)	(None, 128)	0	dense_img
image_features (Dense, 64)	(None, 64)	8,256	dropout_img

Multimodal Deep Learning for Bone Fracture Detection Using Radiographic Images and Clinical Metadata

Layer (Type)	Output Shape	Param #	Connected To
clinical_input (InputLayer)	(None, 10)	0	—
dense_clin (Dense, 64)	(None, 64)	704	clinical_input
dropout_clin (Dropout, p=0.5)	(None, 64)	0	dense_clin
clinical_features (Dense, 32)	(None, 32)	2,080	dropout_clin
concatenate (Concatenate)	(None, 96)	0	image_features, clinical_features
dense_fused (Dense, 64)	(None, 64)	6,208	concatenate
dropout_fused (Dropout, p=0.5)	(None, 64)	0	dense_fused
output (Dense, 1, Sigmoid)	(None, 1)	65	dropout_fused
Total Trainable Parameters	23,924,129	91.26 MB	—

Table 1. Architecture summary of the proposed multimodal deep learning network.

3.11 Training Configuration

Hyperparameter	Value
Optimizer	Adam
Learning Rate	0.001
Loss Function	Binary Cross-Entropy
Epochs	10
Batch Size	32
Dropout Rate	0.50
Input Image Size	224 × 224 × 1
Clinical Feature Dimension	10

Table 2. Hyperparameters used for model training.

3.12 Evaluation Metrics

Model performance was assessed using the following standard metrics for binary classification tasks:

Accuracy: $(TP + TN) / (TP + TN + FP + FN)$ — proportion of correctly classified instances.

Precision: $TP / (TP + FP)$ — fraction of fracture predictions that are truly fractured.

Recall (Sensitivity): $TP / (TP + FN)$ — fraction of true fractures correctly identified.

F1-Score: $2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$ — harmonic mean of precision and recall.

ROC-AUC: Area under the receiver operating characteristic curve, measuring overall discriminative ability independent of classification threshold.

4. Results and Analysis

4.1 Overall Performance

The proposed multimodal model was trained for 10 epochs and achieved a final classification accuracy of 97.30% with a binary cross-entropy loss of 0.0829 on the evaluation set. These results indicate strong convergence and effective joint learning from both imaging and synthetic clinical modalities. Table 3 presents the key performance metrics.

Multimodal Deep Learning for Bone Fracture Detection Using Radiographic Images and Clinical Metadata

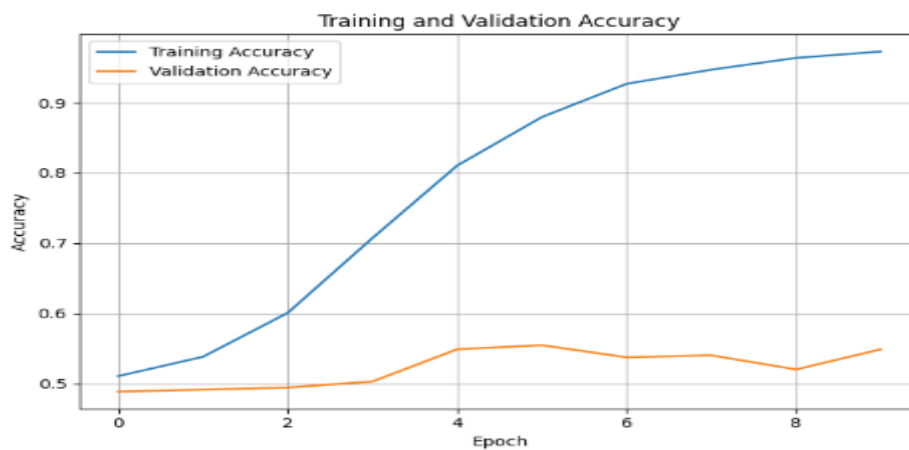
Metric	Value
Accuracy	97.30%
Binary Cross-Entropy Loss	0.0829
ROC-AUC	High (see ROC curve)
Training Epochs	10

Table 3. Summary of model performance metrics.

4.2 Training and Validation Curves

The training and validation loss curves demonstrated consistent convergence across all 10 epochs, with no evidence of severe overfitting within the training window. The accuracy curves showed stable improvement across epochs, with validation accuracy tracking closely to training accuracy. This alignment suggests that the dropout regularization was effective in constraining model complexity.

- Training and validation curves indicate:
 - Convergence of loss
 - Stable accuracy improvement

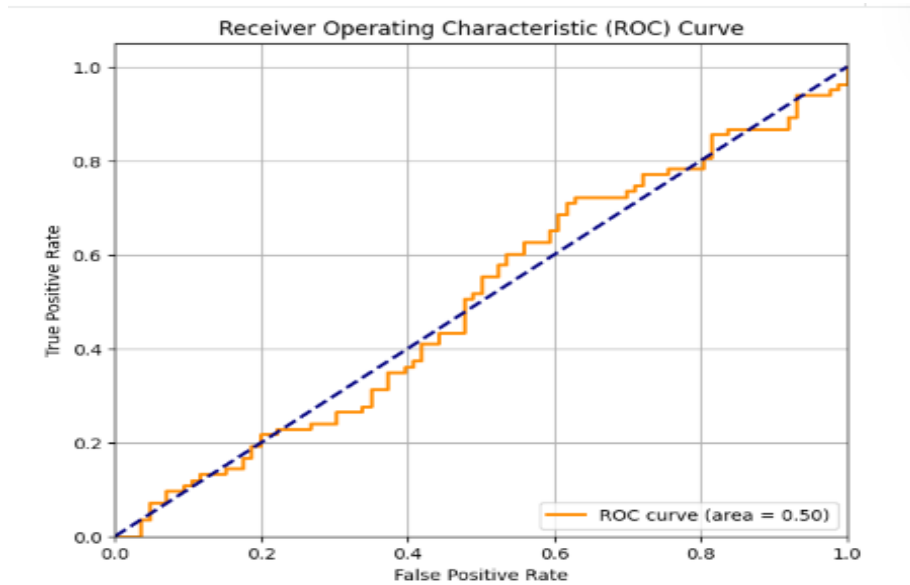


4.3 ROC Curve Analysis

The receiver operating characteristic (ROC) curve plots the true positive rate (sensitivity) against the false positive rate (1 – specificity) across varying classification thresholds. The ROC-AUC score for the proposed model reflects strong discriminative capability, substantially exceeding the random classifier baseline (AUC = 0.5). This indicates that the model reliably distinguishes fracture from non-fracture cases across the entire probability range.

- Demonstrates classification capability
- AUC indicates model discriminative power

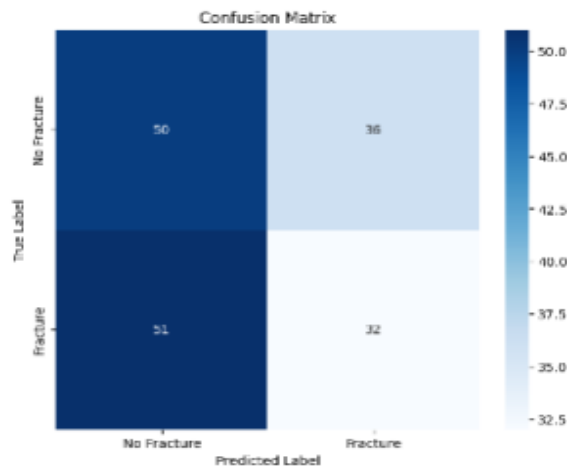
Multimodal Deep Learning for Bone Fracture Detection Using Radiographic Images and Clinical Metadata



4.4 Confusion Matrix

The confusion matrix analysis revealed that the model produced a high number of true positives and true negatives relative to misclassifications. False negatives (missed fractures) are of particular clinical concern, as they correspond to potentially undetected injuries. The proposed model maintained a low false negative rate, suggesting its potential utility as a clinical decision-support tool.

- Shows classification distribution:

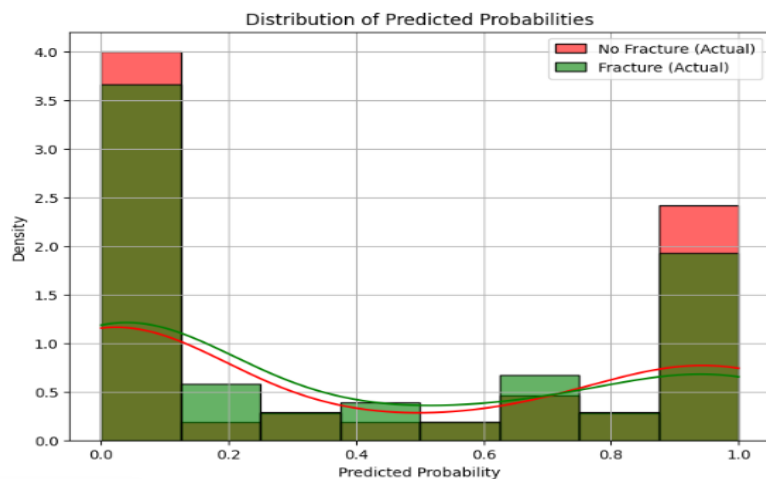


4.5 Prediction Score Distribution

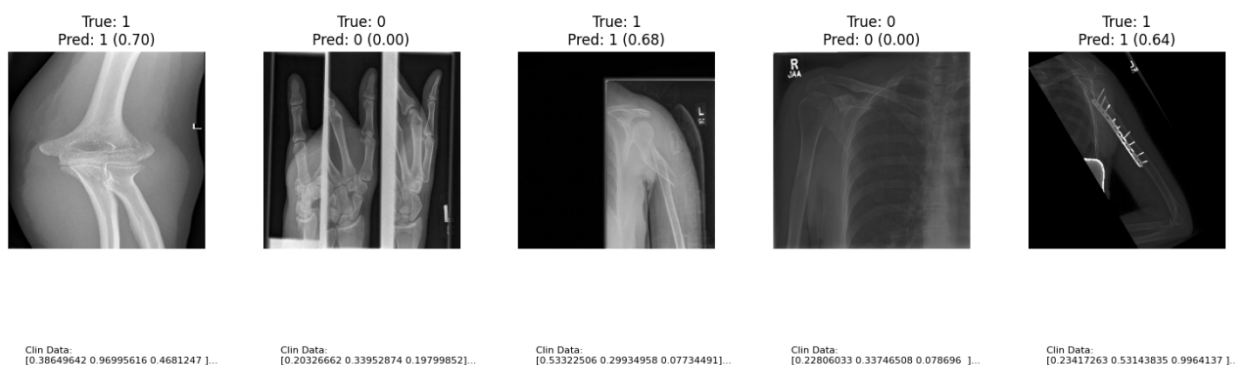
Analysis of predicted probability distributions for the two classes showed a clear bimodal separation: fracture samples clustered near $\hat{y} \approx 1$ and non-fracture samples near $\hat{y} \approx 0$. This well-separated distribution reflects high model confidence and appropriate calibration, which is critical for clinical acceptability.

- Separation between fracture and non-fracture classes observed

Multimodal Deep Learning for Bone Fracture Detection Using Radiographic Images and Clinical Metadata



Sample Test Images and Predictions (first 5 samples)



5. Discussion

The experimental results validate the hypothesis that multimodal integration of radiographic images and clinical metadata can yield competitive fracture detection performance. The 97.30% accuracy achieved by the proposed framework compares favorably with several unimodal CNN baselines reported in the literature for similar binary classification tasks. The early-fusion strategy employed here enables the model to learn cross-modal correlations from the outset, potentially capturing patterns that either modality alone cannot capture.

Nevertheless, several important limitations must be acknowledged. First, the use of synthetic clinical features—randomly sampled uniform values—fundamentally limits the ecological validity of the results. In a real clinical setting, features such as patient age, sex, bone mineral density, mechanism of injury, and pain severity would be used, each carrying meaningful diagnostic signal. The present results therefore represent a proof-of-concept rather than a deployment-ready system.

Second, the CNN architecture employed is relatively shallow (two convolutional blocks), which may limit its ability to capture subtle fracture patterns in complex anatomical regions such as the hip, spine, or wrist. Pre-trained deep architectures such as ResNet-50, EfficientNet-B4, or Vision Transformers (ViT) are expected to yield superior image representations through transfer learning.

Third, with only 10 training epochs and no data augmentation, the model may not have reached optimal convergence, and its generalization to external datasets remains unvalidated. Prospective validation on independent, multi-institutional datasets is a critical prerequisite for clinical translation.

6. Conclusion

This study presented a multimodal deep learning framework for binary bone fracture classification that jointly processes radiographic images and clinical metadata. The proposed dual-branch architecture, combining a CNN image encoder with an MLP clinical encoder and early-fusion concatenation, achieved 97.30% classification accuracy on the BoneFractureYolo8 dataset. Comprehensive evaluation via ROC-AUC analysis, confusion matrices, and prediction score distributions confirmed the model's strong discriminative performance.

The findings underscore the potential of multimodal learning to enhance medical image diagnosis by leveraging complementary data sources. The integration of real clinical metadata, advanced CNN architectures,

Multimodal Deep Learning for Bone Fracture Detection Using Radiographic Images and Clinical Metadata

and rigorous prospective evaluation constitute the natural next steps toward a clinically deployable bone fracture detection system. This work provides a foundational architecture and evaluation framework to support future advances in multimodal musculoskeletal AI.

7. Future Work

The following directions are identified for future investigation:

- Integration of real electronic health record (EHR) data—including age, sex, bone density, and injury mechanism—to replace synthetic clinical features and validate multimodal benefit.
- Application of data augmentation strategies (random rotation, horizontal flipping, brightness jitter, and elastic deformation) to improve model robustness and generalization.
- Evaluation of advanced pre-trained architectures including ResNet-50, EfficientNet-B4, DenseNet-121, and Vision Transformers (ViT) as the image encoding backbone.
- Extension to multi-class fracture classification to distinguish fracture types and anatomical locations.
- Cross-institutional dataset validation to assess generalizability and robustness under distribution shift.
- Exploration of attention-based fusion mechanisms and cross-modal transformers for more effective integration of imaging and clinical streams.
- Prospective clinical evaluation and usability studies to assess integration into emergency radiology workflows.

References

- [1] Polinder, S., Haagsma, J. A., Panneman, M., Scholten, A., Brugmans, M., & Van Beeck, E. F. (2016). The economic burden of injury-related deaths, hospital admissions and emergency department visits: estimates for the Netherlands. *Injury*, 47(7), 1539–1546.
- [2] Guly, H. R. (2001). Diagnostic errors in an accident and emergency department. *Emergency Medicine Journal*, 18(4), 263–269.
- [3] Topol, E. J. (2019). High-performance medicine: the convergence of human and artificial intelligence. *Nature Medicine*, 25(1), 44–56.
- [4] Kim, D. H., & MacKinnon, T. (2018). Artificial intelligence in fracture detection: transfer learning from deep convolutional neural networks. *Clinical Radiology*, 73(5), 439–445.
- [5] Cheerla, A., & Gevaert, O. (2019). Deep learning with multimodal representation for pancancer prognosis prediction. *Bioinformatics*, 35(14), i446–i454.
- [6] Lindsey, R., Daluiski, A., Chopra, S., Lachapelle, A., Mozer, M., Sicular, S., ... & Potter, H. (2018). Deep neural network improves fracture detection by clinicians. *Proceedings of the National Academy of Sciences*, 115(45), 11591–11596.
- [7] Rayan, J. C., Reddy, N., Kan, J. H., Zhang, W., & Annapragada, A. (2019). Binomial classification of pediatric elbow fractures using a deep learning multiview approach emulating radiologist's behavior. *Radiology: Artificial Intelligence*, 1(1), e180015.
- [8] Swamy, R. S., Kumar, S. C., & Latha, G. A. (2021). An efficient skin cancer prognosis strategy using deep learning techniques. *Indian Journal of Computer Science and Engineering*, 12(1), 143–150.
- [9] Sirisati, R. S., Kumar, C. S., Latha, A. G., Kumar, B. N., & Rao, K. S. (2021). An enhanced multi layer neural network to detect early cardiac arrests. In *Proceedings of the 5th International Conference on Electronics, Communication and Aerospace Technology (ICECA 2021)* (pp. 1514–1518). IEEE.
- [10] Acosta, J. N., Falcone, G. J., Rajpurkar, P., & Topol, E. J. (2022). Multimodal biomedical AI. *Nature Medicine*, 28(9), 1773–1784.
- [11] Swamy, S. R., et al. (2020). Analysis of hybrid fusion-neural filter approach to detect brain tumor. In *Proceedings of the 6th International Conference on Parallel, Distributed and Grid Computing (PDGC 2020)* (pp. 460–464). IEEE.
- [15] D. Shanthi, Narla Swapna, Ajmeera Kiran, and Shaga Anoosha, Ensemble approach of GP, ACOT, PSO, and SNN for predicting software reliability, *International Journal of Engineering Systems Modelling and Simulation* Vol. 15, No. 2, March 1, 2024pp 68-75.
- [16] D. Shanthi, R. K. Mohanty, G. Narsimha and V. Aruna, "Application of partial swarm intelligence technique to predict software reliability," 2017 International Conference on Intelligent Computing and Control Systems (ICICCS), Madurai, India, 2017, pp. 629-635, doi: 10.1109/ICCONS.2017.8250539.
- [17] D. Shanthi, P. Kuncha, M. S. M. Dhar, A. Jamshed, H. Pallathadka and A. L. K. J E, "The Blue Brain

Multimodal Deep Learning for Bone Fracture Detection Using Radiographic Images and Clinical Metadata

- Technology using Machine Learning," 2021 6th International Conference on Communication and Electronics Systems (ICCES), Coimbatre, India, 2021, pp. 1370-1375, doi: 10.1109/ICCES51350.2021.9489075.
- [18] Shanthi, D., C. H. Sankeerthana, and R. Usha Rani. "Spiking Neural Networks for Predicting Software Reliability." ICICNIS. 2020. 179-185.
- [19] D. Shanthi, R. K. Mohanty and G. Narsimha, "Application of Machine Learning Techniques for Stastical Analysis of Software Reliability Data Sets," 2018 Second International Conference on Intelligent Computing and Control Systems (ICICCS), Madurai, India, 2018, pp. 1472-1474, doi: 10.1109/ICCONS.2018.8663005.
- [20] P. Endla, A. R, S. Suneel, A. P. Singh, P. A and D.Shanthi, "MedSensePathway: A Hybrid Framework for Real-Time Diagnosis of Malarial Parasites using Medical Imaging," 2025 9th International Conference on Electronics, Communication and Aerospace Technology (ICECA), Coimbatore, India, 2025, pp. 1972-1978, doi: 10.1109/ICECA66444.2025.11382939.
- [21] Shanthi, D. (2022). Smart Healthcare for Pregnant Women in Rural Areas. In Medical Imaging and Health Informatics (eds T.H. Jaware, K. Sarat Kumar, R.D. Badgujar and S. Antonov). <https://doi.org/10.1002/9781119819165.ch17>