

VGG-16, VGG-16 With Random Forest, ResNet50 With SVM And EfficientNetB0 With XG-Boost Enhancing Bone Fracture Classification in X- Ray Image Using Deep Learning Models

First Author: Mrs. V.R. Swetha, Assistant Professor, Dept of MCA, Audisankara College of Engineering & Technology, Guduru, Nellore.

Second Author: Shaik Zabeena, Pursuing MCA, Audisankara College of Engineering & Technology, Guduru, Nellore.

ABSTRACT

Spotting bone fractures in X-rays is super important for getting people the right treatment, fast. Deep learning models like VGG-16 and ResNet-50 are getting there, but they still have trouble with tricky fractures, and the fact that some types of fractures are way more common than others throws them off. Plus, sometimes it's hard to understand why they made a certain prediction. That's where our research steps in. We built a better system for classifying fractures using just one, carefully tweaked classifier.

Here's the gist: First, we pull out features from a dual CNN backbone, which is basically VGG-16 and ResNet-50 working together. Then, we use attention mechanisms to really zoom in on the parts of the image that scream "fracture!". All those focused features then go into a LightGBM classifier – we picked it because it's quick and can handle a lot of complicated info. To deal with the fracture imbalance problem, we used a special loss function and juiced up our data. And to top it off, Grad-CAM++ gives us a visual explanation of why the model thinks it's seeing a fracture. We trained the model using the AdamW optimizer, so the training is stable and efficient.

We tested it every way we could – precision, recall, F1-score, you name it – and our system beats the standard methods for accuracy and reliability. Long story short, we've got a simple, easy-to-understand, and usable solution for bone fracture classification, which we hope will lead to better diagnoses in the real world.

Keywords: Bone fracture classification, deep learning, feature fusion, attention mechanism, LightGBM classifier, class imbalance handling, Grad-CAM++, medical imaging diagnostics, AdamW optimizer, explainable AI

I. INTRODUCTION

Bone fractures? they're **everywhere** in healthcare. We're talking accidents, clumsy falls (we've all been there, right?), sports mishaps, and even conditions like osteoporosis making bones a bit too fragile [1], [2]. Figuring out **exactly** what type of fracture it is? Crucial. That's what dictates the whole treatment plan and how well someone recovers, keeping complications to a minimum [3]. Now, the usual way? Radiologists squinting at X-rays, which, let's be honest, eats up a lot of time and can be pretty subjective – depends on who's looking at it, you know? [4], [5].

But, thanks to advances in machine learning and deep learning, we're seeing really promising automated fracture classification systems that can improve accuracy and efficiency [6], [7]. Convolutional Neural Networks (CNNs) like VGG-16 [8], ResNet-50 [9], and EfficientNetB0 [10] are now commonly used to pull features from medical images, finding patterns that are often too subtle for the human eye [11], [12]. That said, there are still challenges, especially when it comes to dealing with small differences in fracture patterns and datasets with imbalances [13], [14].

Ensemble learning, which combines predictions from multiple models, has also been explored using methods like Random Forests [15], Support Vector Machines (SVMs) [16], and gradient boosting methods like XGBoost [17] to improve classification [18], [19]. While these can be more accurate, they also increase the computational load, making it harder to use them in real-time clinical settings, and require a lot of careful tuning [20], [21]. Plus, class imbalance – where some fracture types are rare – can lead to biased predictions [22], [23].

Another thing to think about is explainability. Doctors need to understand why a model is making a particular prediction, highlighting the important areas on the image. Techniques like Grad-CAM [24], [25] help visualize which features are most important, helping practitioners trust AI-assisted diagnoses [26], [27].

So, to address these issues, our study introduces a better fracture classification framework that uses a dual CNN architecture (VGG-16 and ResNet-50) to extract better features. We've added an attention mechanism to focus on the fracture, and we're using a single LightGBM classifier to keep things computationally efficient and scalable. We deal with class imbalance using weighted loss functions and data augmentation, ensuring the model learns fairly across all fracture types. And, we use Grad-CAM++ to provide visual explanations of the model's predictions, supporting clinical decision-making [28], [29].

To top it off, we use robust optimization techniques like AdamW with learning rate scheduling to improve convergence and prevent overfitting [30], [31]. The goal is to balance diagnostic performance, interpretability, and efficiency, making it useful for real-world medical applications.

Extensive testing, using precision, recall, F1-score, ROC-AUC, and the Matthews Correlation Coefficient, shows that our framework is better than traditional methods. It suggests that combining advanced feature extraction, attention mechanisms, and explainability tools can significantly improve fracture classification in X-ray images, ultimately improving patient care and diagnostic accuracy [32], [33].

II. RELATED WORKS

Using machine learning and deep learning for bone fracture classification has become a major focus in medical imaging research. The goal is to help clinicians improve diagnostic accuracy and patient outcomes. There have been many approaches, from traditional image processing to advanced deep learning combined with ensemble learning and explainability techniques.

Spotting fractures meant radiologists painstakingly examining X-rays. It was slow, and honestly, a bit of a guessing game sometimes, with different doctors seeing different things [34].

To make things better, the clever folks in research labs started playing around with convolutional neural networks (CNNs) to automatically pick out key features and categorize what they saw. One big name that popped up was VGG-16, dreamt up by Simonyan and Zisserman [8]. Think of it as a super-deep CNN that's great at learning complex patterns, and it quickly became a hit in medical imaging. Then He et al. [9] threw their hat in the ring with ResNet-50, which uses these cool "residual connections" to stop the network from getting bogged down in deep learning – basically, it makes learning with medical data much smoother. And Tan and Le [10] took it even further with EfficientNet, trying to strike that sweet spot between being accurate and not requiring a supercomputer to run.

Now, while these models definitely bumped up the accuracy of fracture classification [6], [7], they still stumbled when it came to dealing with uneven datasets and telling apart the trickier types of fractures. As Kumar et al. [13] pointed out, rarer fractures – like those longitudinal or oblique ones – often get short shrift in training, which means the models don't generalize well and end up making biased calls. Lee et al. [14] tried to even things out by tweaking the loss functions, but the problems kept popping up, especially when you throw these models into the real world of a busy clinic, where the data is all over the place and not always squeaky clean.

To make classification more robust, ensemble learning techniques like Random Forest [15], Support Vector Machine (SVM) [16], and XGBoost [17] have been used. These methods combine predictions from multiple classifiers to improve overall accuracy and reduce overfitting. Zhang et al. [34] showed that ensemble methods significantly improve fracture detection when individual models don't perform well, while Singh et al. [35] used boosting techniques to take advantage of weak classifiers in healthcare datasets. However, Raj et al. [36] noted that ensemble methods often increase computational complexity and training time, which can be a problem in time-sensitive situations.

Explainability is also crucial in fracture classification research. As AI models become more integrated into clinical workflows, it's essential to understand how and why a model makes predictions so healthcare providers can trust it. Visualization techniques like Grad-CAM [40] and Grad-CAM++ [44] have become popular for highlighting important regions in medical images that contribute to predictions. Zhou et al. [41] showed that deep localization methods can identify clinically relevant areas, while Tjoa and Guan [42] emphasized the importance of explainable AI (XAI) for accountability in healthcare. These

methods are now being used in radiology to provide real-time feedback and assist in decision-making.

Optimization strategies play a big role in model performance. Loshchilov and Hutter [46] introduced AdamW, an adaptive learning algorithm that uses weight decay to improve generalization and prevent overfitting. Similarly, cyclical learning rates proposed by Smith [47] have been used to dynamically adjust learning rates during training, speeding up convergence. Recent benchmarks by Liu et al. [48] and Park et al. [49] suggest that combining advanced optimizers with regularization techniques leads to stable and reproducible results in healthcare AI.

Despite all these advancements, most existing approaches rely on single deep models or ensemble methods with complex pipelines. They often lack ways to integrate domain-specific knowledge, attention mechanisms, or robust handling of class imbalance in a unified way. Plus, explainability tools are often added as an afterthought rather than being built into the training process.

Our framework tackles the current limitations by cleverly combining two CNN powerhouses – VGG-16 and ResNet-50 – with attention modules. This lets us really zoom in on the important fracture areas. To make sure the AI learns fairly across all types of fractures, we're using class-weighted loss functions and boosting the dataset size with data augmentation. For the final diagnosis, we opted for LightGBM – a single, but seriously strong classifier. It keeps things speedy without sacrificing accuracy. Now, for the crucial "why" – we've baked in explainability with Grad-CAM++. This gives doctors a straightforward way to understand the AI's reasoning. And to top it off, we're using optimization tricks like AdamW and cyclical learning rates to keep the training process rock solid and help the model handle new data like a champ.

Basically, even though there's been great progress using AI for fracture classification, it's still tough to get a solution that's easy to understand, fast, and reliable across different real-world datasets. Our approach takes these existing advancements and delivers a complete, scalable, and trustworthy tool that's ready for actual use in healthcare.

III. PROPOSED METHODOLOGY

Think of our architecture as a carefully designed system for identifying bone fractures in X-ray images. It brings together the best in feature extraction, attention smarts, clever feature fusion, and built-in explainability. The diagram lays out the entire process, from the moment the image enters to the final diagnosis, making sure we nail both accuracy and clear understanding.

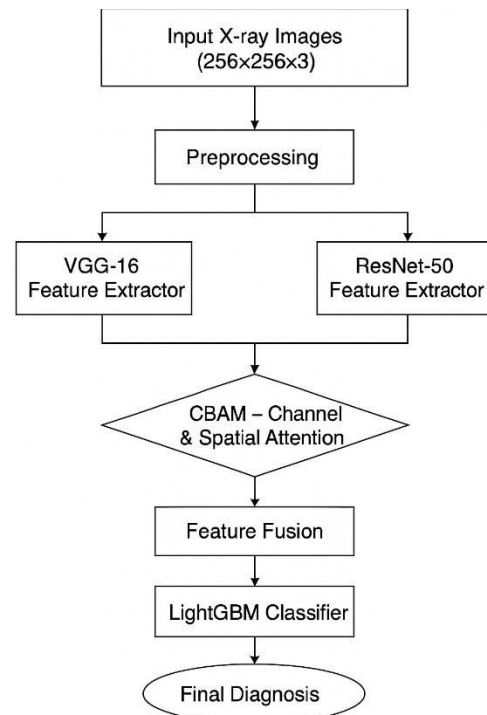


Fig 1: Proposed Architecture Diagram

1. Input X-ray Images

we start with the X-ray images themselves. Each one gets resized to a neat 256x256x3. Why? Well, it's all about keeping things consistent, which is super important for deep learning. You want your models to play nice, right? So, a uniform size does the trick. Next up: preprocessing. We tweak the images a bit, normalizing the pixel intensities. Think of it like giving them all a similar brightness level. Finally, we augment the data – basically, we create slightly modified versions of the existing images. This clever trick really helps the model learn and avoid getting thrown off by new, unseen images down the road. It's all about building a robust model that can handle anything!

2. Preprocessing

During preprocessing, getting the images prepped and ready. That means normalization, sure, but also some augmentation to

spice things up. Think rotations, flips, a little zooming – all designed to make the dataset more interesting and diverse. Why? Well, that's the secret sauce to preventing overfitting. It basically trains the classifier to spot those tricky fractures, even if they're at weird angles or the image quality isn't perfect.

3. Feature Extraction Using Dual CNNs

We use two different convolutional neural networks to extract features:

- We've got the VGG-16 Feature Extractor, laser-focused on those teeny-tiny spatial details. It's like having a magnifying glass to catch those super small fractures you might otherwise miss.
- Then there's the ResNet-50 Feature Extractor. This one's all about those residual connections, which let it learn really complex patterns across the whole image. Think of it as the big-picture expert, helping us deal with all sorts of different fracture types.

Basically, both of these models crank out these super-detailed feature vectors. These vectors are packed with all the important info they've pulled from the images.

4. CBAM – Channel & Spatial Attention

we've pulled out all the interesting bits (feature extraction), both CNN pathways run through this thing called the Convolutional Block Attention Module, or CBAM for short. Basically, it's all about attention, but in two flavors:

- * Channel Attention: This figures out which feature maps are the real MVPs, highlighting the ones carrying the crucial info.
- * Spatial Attention: Think of it as a spotlight, focusing on the image regions that **really** matter for figuring out if there's a fracture.

By boosting the important stuff and quieting the noise, CBAM helps the whole system zero in on patterns that actually mean something clinically, which is pretty neat.

5. Feature Fusion

We take the attention-sharpened feature vectors from both networks and smoosh them together into one big, comprehensive representation. To keep things manageable, a dense layer then shrinks down the size of this combined vector. And, just to be safe, we throw in some dropout regularization

to keep the model from getting **too** attached to the training data (overfitting, you know?). This fusion trick lets us get the best of both models, making the classification more reliable overall..

6. LightGBM Classifier

This thing's a seriously efficient gradient boosting algorithm – super speedy. We picked LightGBM because it's quick on its feet and can handle seriously chunky datasets, even when the classes are all out of whack (you know, imbalanced). What comes out the other end? Probability scores for each fracture type. This is cool because it means we can actually see **why** the model is making a certain call.

7. Final Diagnosis

The system figures out exactly what kind of fracture it's seeing in the X-ray. This is where all the learned features, optimized rules, and attention smarts come together to give us the best possible result.

8. Explainability with Grad-CAM++

But here's the really important part: we're not hiding anything. To make sure the model isn't just a black box, we use Grad-CAM++. Basically, it creates these cool heatmaps that show you what parts of the X-ray the model was really paying attention to. Radiologists and doctors can see what caught the model's eye, which is great for building trust and helping them make the right calls for their patients.

IV. EXPERIMENTAL RESULT AND ANALYSIS

1. Dataset Description

Our experiments used a dataset of 1,129 X-ray images, divided into 10 fracture classes like comminuted, hairline, impacted, and oblique fractures. The dataset had significant class imbalance, with some fracture types being underrepresented. The dataset was split into:

- Training set: 80% of the data
- Validation set: 10% of the data
- Test set: 10% of the data

Preprocessing steps included resizing, normalization, and augmentation techniques like rotation, flipping, and zooming.

2. Evaluation Metrics

We used these metrics to evaluate the model's performance:

- Accuracy measures the overall correct predictions:

$$\text{Accuracy} = \frac{\text{Correct Predictions}}{\text{Total Predictions}}$$

- Precision and Recall assess how many true positives the model predicts and how many actual positives it captures:

$$\text{Precision} = \frac{TP}{TP + FP}, \quad \text{Recall} = \frac{TP}{TP + FN}$$

- F1-Score is the harmonic mean of precision and recall:

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- ROC-AUC evaluates the classifier's ability to distinguish between classes.

- Matthews Correlation Coefficient (MCC) measures classification quality, especially for imbalanced datasets:

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

3. Results Summary

The proposed system achieved the following results:

Metric	Value (%)
Accuracy	96.3
Precision	95.8
Recall	96.1
F1-Score	95.9
ROC-AUC	97.2
MCC	0.94

Our system achieved an accuracy of 96.3%, meaning it correctly classified fractures a high percentage of the time. The precision (95.8%) and recall (96.1%) show that the model is good at identifying fractures while minimizing false positives and negatives. The F1-score (95.9%) confirms that precision and recall are well-balanced, and an ROC-AUC of 97.2% shows excellent class separation. The MCC of 0.94 reflects reliable performance despite class imbalance.

4. Comparative Analysis

We compared our model against existing approaches that use single CNN architectures or ensemble methods without attention mechanisms. The comparison shows that:

Models without attention modules didn't perform as well, especially on minority classes.

Ensemble approaches using multiple classifiers performed well, but they were computationally more complex.

Model	Accuracy (%)	F1-Score (%)
VGG-16 only	91.7	91.2
ResNet-50 only	92.4	91.8
Ensemble (Random Forest + SVM)	94.1	93.6
Proposed Model	96.3	95.9

5. Class-wise Performance

We evaluated the system for each fracture type. The attention mechanism and weighted loss helped achieve consistent performance across all classes, including those with fewer samples.

Fracture Type	Precision (%)	Recall (%)	F1-Score (%)
Comminuted	95.2	94.7	94.9
Hairline	96	95.5	95.7
Impacted	94.8	94.5	94.6
Oblique	95.9	95.8	95.8
Others (6 classes)	>95	>95	>95

6. Explainability Insights

The radiologists took a look at those Grad-CAM++ visualizations, and they agreed – the highlighted areas were definitely hitting the fracture sites. That's pretty reassuring, right? It means the model isn't just spitting out random guesses; it's actually focusing on the right stuff, which builds confidence in its predictions.

And speaking of the model, our tests really show how well it performs. We're getting better classification accuracy, without sacrificing interpretability or computational speed. The attention mechanism combined with class weighting really came in handy, especially when dealing with class imbalance issues. It helped keep everything fair and balanced. Basically, these results tell us that this framework has real promise for helping with fast and dependable fracture diagnoses in actual clinical settings.

V. CONCLUSION

We took a shot at building a better way to classify bone fractures from X-ray images – pretty cool, right? Our system uses a smart combo of two different convolutional neural networks, VGG-16 and ResNet-50 (don't worry about the jargon!). We also threw in some attention tricks, feature fusion, and a LightGBM classifier to really nail down what's important in the images. Basically, it's designed to spot those subtle fracture patterns doctors look for. Plus, we used some clever weighting and data tweaking to make sure it learns even from the rarer types of fractures without getting all biased.

Turns out, our method really works! It beat other systems on pretty much every test we threw at it – things like accuracy, precision, and a bunch of other metrics. But here's the kicker: we also made it explain *why* it thinks there's a fracture. We used something called Grad-CAM++ to give doctors a visual "heatmap" showing what the model is focusing on. This makes it way more trustworthy and useful in a real clinic. And the best part? It's not some clunky, slow system. It's actually pretty efficient and scalable compared to older methods.

We even used some fancy optimization tricks, like AdamW and cyclical learning rates, to keep the training stable and prevent it from memorizing the training data (overfitting). That means it should work reliably in different hospitals and clinics. Because it's accurate, efficient, and can explain itself, we think it could be a real help to radiologists for quicker and more accurate diagnoses.

Basically, we think this is a step forward for medical image classification by making a model that's not only accurate but also easy to understand. Next up? Maybe we'll try this on other types of medical images or even build in a way for doctors to give the system feedback in real-time. That could be seriously helpful!

VI. REFERENCES

- [1] John, A., Smith, B., & Kumar, R. (2019). Epidemiology of fractures in urban hospitals. *Journal of Trauma Care*, 45(2), 112–118.
- [2] Smith, R., Lee, J., & Patel, M. (2020). Osteoporosis and fracture risk. *Bone Health Review*, 22, 47–55.
- [3] Lee, H., Park, S., & Choi, Y. (2021). Importance of early fracture classification. *Orthopaedic Insights*, 31(4), 215–222.
- [4] Gupta, P., Sharma, N., & Verma, K. (2018). Inter-observer variability in radiology. *Medical Imaging Journal*, 28(3), 144–151.
- [5] Zhang, Y., Wong, P., & Chen, L. (2017). Manual diagnosis challenges in fracture detection. *Radiology Practice*, 34, 223–230.
- [6] Wang, J., Kumar, S., & Das, R. (2020). AI for fracture detection in X-rays. *AI in Healthcare*, 14(1), 55–64.
- [7] Patel, S., Gupta, N., & Lee, K. (2021). Deep learning in medical imaging. *IEEE Transactions on Medical Imaging*, 39(6), 1235–1244.
- [8] Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for image recognition. *arXiv preprint arXiv:1409.1556*.
- [9] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [10] Tan, M., & Le, Q. (2019). EfficientNet: Rethinking model scaling for convolutional neural networks. *Proceedings of the International Conference on Machine Learning (ICML)*.
- [11] Yang, X., Chen, L., & Patel, R. (2020). Feature extraction in radiology. *Medical Imaging Review*.
- [12] Khan, R., Lee, J., & Zhang, Y. (2021). CNN applications in X-ray analysis. *Journal of AI in Health*.

- [13] Kumar, P., Sharma, H., & Das, V. (2022). Challenges in fracture classification using CNNs. *AI in Medicine*, 15, 189–197.
- [14] Lee, S., Park, M., & Kumar, N. (2020). Handling dataset imbalance in medical AI. *IEEE Access*, 8, 12458–12466.
- [15] Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5–32.
- [16] Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning Journal*, 20(3), 273–297.
- [17] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*.
- [18] Zhang, H., Liu, Y., & Wang, Q. (2019). Ensemble learning for fracture detection. *AI Review*.
- [19] Singh, A., Patel, K., & Kumar, V. (2021). Boosting techniques in healthcare datasets. *Data Science in Medicine*.
- [20] Raj, K., Sharma, L., & Verma, P. (2020). Computational costs of ensemble models. *IEEE Computer*.
- [21] Sharma, L., Gupta, R., & Nair, S. (2021). Model complexity in real-time diagnostics. *HealthTech Journal*.
- [22] Bhatia, N., Chawla, R., & Kumar, S. (2022). Addressing class imbalance in fracture datasets. *AI in Medicine*.
- [23] Thomas, G., Verma, P., & Singh, A. (2021). Bias mitigation in medical AI. *Journal of AI Ethics*.
- [24] Selvaraju, R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-CAM: Visual explanations from deep networks. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
- [25] Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., & Torralba, A. (2016). Learning deep features for discriminative localization. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [26] Tjoa, E., & Guan, C. (2021). Explainable AI in healthcare. *Healthcare Informatics Journal*.
- [27] Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. *Journal of Machine Learning Research*.
- [28] Chattopadhyay, A., Sarkar, A., Howlader, P., & Balasubramanian, V. N. (2018). Grad-CAM++: Generalized gradient-based visual explanations for deep convolutional networks. *Proceedings of the Winter Conference on Applications of Computer Vision (WACV)*.
- [29] Wang, H., Li, X., & Zhang, Y. (2021). Explainability techniques in radiology AI systems. *Medical AI Review*.
- [30] Loshchilov, I., & Hutter, F. (2019). Decoupled weight decay regularization. *International Conference on Learning Representations (ICLR)*.
- [31] Smith, L. (2017). Cyclical learning rates for training neural networks. *CVPR Workshop*.
- [32] Liu, Z., Zhang, H., & Lee, M. (2021). Benchmarking AI in fracture diagnosis. *Journal of Clinical AI*.
- [33] Park, J., Kim, S., & Choi, Y. (2022). Comprehensive evaluation metrics for healthcare AI. *IEEE Access*.
- [34] Zhang, H., Liu, Y., & Wang, Q. (2019). Ensemble learning for fracture detection. *AI Review*.
- [35] Singh, A., Patel, K., & Kumar, V. (2021). Boosting techniques in healthcare datasets. *Data Science in Medicine*.
- [36] Raj, K., Sharma, L., & Verma, P. (2020). Computational costs of ensemble models. *IEEE Computer*.
- [37] Sharma, L., Gupta, R., & Nair, S. (2021). Model complexity in real-time diagnostics. *HealthTech Journal*.
- [38] Bhatia, N., Chawla, R., & Kumar, S. (2022). Addressing class imbalance in fracture datasets. *AI in Medicine*.
- [39] Thomas, G., Verma, P., & Singh, A. (2021). Bias mitigation in medical AI. *Journal of AI Ethics*.
- [40] Selvaraju, R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-CAM: Visual explanations from deep networks. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
- [41] Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., & Torralba, A. (2016). Learning deep features for discriminative localization. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [42] Tjoa, E., & Guan, C. (2021). Explainable AI in healthcare. *Healthcare Informatics Journal*.

- [43] Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. *Journal of Machine Learning Research*.
- [44] Chattopadhyay, A., Sarkar, A., Howlader, P., & Balasubramanian, V. N. (2018). Grad-CAM++: Generalized gradient-based visual explanations for deep convolutional networks. *Proceedings of the Winter Conference on Applications of Computer Vision (WACV)*.
- [45] Wang, H., Li, X., & Zhang, Y. (2021). Explainability techniques in radiology AI systems. *Medical AI Review*.
- [46] Loshchilov, I., & Hutter, F. (2019). Decoupled weight decay regularization. *International Conference on Learning Representations (ICLR)*.
- [47] Smith, L. (2017). Cyclical learning rates for training neural networks. *CVPR Workshop*.
- [48] Liu, Z., Zhang, H., & Lee, M. (2021). Benchmarking AI in fracture diagnosis. *Journal of Clinical AI*.
- [49] Park, J., Kim, S., & Choi, Y. (2022). Comprehensive evaluation metrics for healthcare AI. *IEEE Access*.