

ENSEMBLE DEEP LEARNING APPROACHES FOR MULTICLASS CLASSIFICATION OF HIP REGION FRACTURES IN X-RAY IMAGES

*Note: Sub-titles are not captured in Xplore and should not be used

Anonymous Author(s)
Paper under double-blind review

Abstract—Hip region fractures, including pelvic, femoral neck, intertrochanteric, and subtrochanteric fractures, are critical medical conditions, especially when diagnosed early. These fractures impair mobility, increase risks, and cause complications. Early diagnosis using X-ray imaging is vital for effective treatment. Recent advances in computer vision, particularly ensemble pretrained models, have revolutionized fracture detection by combining various models to improve classification accuracy and stability. This research developed and evaluated ensemble deep learning methods for multiclass classification of hip fractures on X-ray images. The dataset consists of 1000 X-ray images from Sri Lankan hospitals (2022-2023), categorized into two types: non-fracture and fracture. Preprocessing and data augmentation techniques are used to increase dataset diversity. The data was split into 70:15:15 for training, validation, and testing to evaluate performance. The pretrained model architectures include ResNet-101, ResNet-50, EfficientNetB0, and EfficientNetV2, with ResNet-10 taken with different levels of parameterization. ResNet101 achieves the highest test accuracy of 0.8000, followed by ResNet-50 (0.7786), EfficientNetB0 (0.7286), and EfficientNetV2 (0.7500). These pretrained models are induced as ensemble learning models and enhance multiclass hip fracture classification, yielding more accurate results compared to customized vision models. This approach has potential clinical applications, aiding early and reliable diagnosis. Further, it can extend to differentiate the components of the hip region individually with sophisticated data augmentation techniques that help for the classification. This research proves that pretrained models can be effective in biomedical rather than building and training them from scratch.

Index Terms—ensemble deep learning models, hip region fracture, multiclass fracture classification, pretrained vision models, x-ray images

I. INTRODUCTION

Fractures, fundamentally, refer to disruptions in the continuity of bone structure, often resulting from excessive mechanical stress, traumatic impacts, or underlying pathological conditions such as osteoporosis [1], cancer, or metabolic disorders that weaken bone integrity. These breaks can vary in severity, pattern, and etiology, ranging from microscopic stress fractures caused by repetitive overload to catastrophic comminuted fractures involving multiple fragments from high-energy trauma. In the human skeletal system, fractures trigger



Fig. 1. Hip fracture

a complex healing process involving inflammation, soft callus formation [2], hard callus remodeling, and eventual restoration, but complications like non-union or malunion can arise if not properly managed. Among all fracture sites, those in the hip region—encompassing the proximal femur, including the femoral neck, head, and trochanters are particularly prevalent and impactful, accounting for a significant proportion of orthopedic emergencies worldwide; for instance, femoral neck fractures alone affect over 1.6 million [3] people annually, predominantly the elderly population over 65 years, where low bone density exacerbates vulnerability to low-energy falls. This region's fractures are notorious for their association with vascular compromise, given the femoral neck's reliance on retrograde blood supply from the medial and lateral circumflex arteries [4], making them a leading cause of morbidity in orthopedics.

Hip region fractures profoundly affect individuals by severely impairing mobility, leading to prolonged bed rest that increases risks of secondary complications such as deep vein thrombosis, pressure ulcers, pneumonia, and muscle atrophy [4], with mortality rates reaching up to 0.30 within the first year post-injury due to these cascading health issues. Economically, they burden healthcare systems with high costs for hospitalization, surgery, and rehabilitation, often resulting in long-term dependency and reduced quality of life, particularly in aging populations where comorbidities like cardiovascular disease amplify the impact. To resolve these challenges [5], early and precise diagnosis is crucial, achieved through advanced imaging analysis; ensemble deep learning

approaches offer a robust solution by integrating multiple neural network predictions to enhance reliability, without delving into specific model architectures, enabling automated detection and differentiation of fracture patterns in X-ray images. Complementing this, multiclass classification refines the process by categorizing fractures into distinct severity levels [6], allowing for tailored treatment strategies—such as conservative monitoring for minor cases or urgent surgical fixation for severe ones—thus improving diagnostic accuracy, reducing radiologist workload, and facilitating timely interventions that mitigate long-term effects.

Building on this foundation, the multiclass framework in ensemble deep learning begins with the non-fracture category, which serves as the baseline for normal hip anatomy in X-ray images, characterized by uninterrupted cortical lines, aligned trabecular patterns within the femoral neck, and smooth articulation between the femoral head and acetabulum, with no evidence of density irregularities or soft tissue swelling that might mimic pathology. This class is vital for minimizing false positives in AI systems, where ensembles aggregate features like bone texture symmetry and edge continuity to achieve high specificity (often above 0.95) [6], distinguishing artifacts from true anomalies; clinically, non-fracture identification supports preventive measures like osteoporosis screening via dual-energy X-ray absorptiometry (DEXA) to avert future risks. Progressing to the non-displaced incomplete fracture, equivalent to Garden Type I, this involves a partial crack often valgus-impacted—where the femoral head tilts outward slightly without full separation, typically from low-trauma falls in osteopenic bones, visible as subtle medial trabecular disruptions on radiographs while preserving overall alignment and vascular supply, with avascular necrosis risks low at 5-10. Ensemble models deeply analyze these faint discontinuities through layered feature extraction, boosting recall for early detection and enabling non-surgical management like protected weight-bearing to promote natural healing [7].

Further elaborating the concept, the non-displaced complete fracture, akin to Garden Type II, features a full transverse break across the femoral neck without fragment shift, maintained by intact periosteum, arising from moderate trauma and appearing as a clear line traversing the neck with preserved trabecular alignment, keeping avascular risks at 10-20 due to minimal disruption of retinacular vessels [8]. In ensemble deep learning, this is classified by fusing contextual symmetries and multi-scale filters to yield F1-scores exceeding 0.95, guiding treatments like percutaneous screw fixation for stability [9]. The complete fracture incompletely displaced, or Garden Type III, marks increased severity with partial displacement—typically 0.50 or less—where fragments rotate or angulate but remain partially connected, often from higher energy impacts, evident in X-rays as misaligned trabeculae [10] and widened fracture gaps, elevating osteonecrosis risks to 20-35 and necessitating open reduction internal fixation (ORIF) [11] to realign and secure. Finally, the complete fracture completely displaced, Garden Type IV, represents full separation with the femoral head freely rotating or translat-

ing, usually from severe trauma, shown radiographically as complete detachment with disrupted blood flow, posing 35-50 avascular necrosis rates and often requiring hemiarthroplasty or total hip replacement [11]; ensembles excel in detecting these by prioritizing displacement features, ensuring accurate multiclass outputs for urgent surgical planning.

II. LITERATURE REVIEW

The study [1] investigated the application of convolutional neural networks (CNNs) for the detection of hip fractures on X-ray images. Their findings revealed that CNN-based models could achieve diagnostic accuracy comparable to that of experienced radiologists, showing the potential of AI to serve as a reliable clinical support tool. The key finding was that AI could significantly reduce human error and provide faster interpretations in fracture diagnosis. The limitation, however, was that the study used a relatively small dataset from a single institution, which limited the model's generalizability to broader populations and different imaging settings.

The study [2] explored the effectiveness of ensemble learning methods in medical image classification tasks. Their study demonstrated that combining multiple deep learning architectures into a single ensemble model enhanced classification accuracy and produced more stable outputs compared to relying on a single CNN. The key finding was that ensemble models offered improved robustness and reliability in medical imaging applications. The limitation was that training and implementation required significantly higher computational resources, which could be a challenge in real-world clinical environments, especially in low-resource healthcare systems.

In this study [3] applied deep learning to the automated detection and classification of femoral neck and intertrochanteric fractures on hip X-rays. The results showed that the AI system achieved high sensitivity and specificity, providing rapid and standardized interpretations that could assist radiologists and reduce diagnostic variability. The key finding was that automated models could effectively classify different hip fracture types with clinical-level performance. The limitation was that the model's accuracy declined in cases involving poor-quality images or anatomical variations, which highlighted the importance of large and diverse datasets for training.

The study [4] examined the potential of AI-based diagnostic systems for use in low-resource healthcare settings where radiologists are scarce. Their study emphasized the value of AI-driven platforms, particularly mobile or cloud-based solutions, in providing real-time hip fracture detection and bridging diagnostic gaps. The key finding was that AI could support frontline healthcare providers by offering accessible diagnostic assistance without the need for expert radiologists on site. The limitation was that such systems depended heavily on internet connectivity and hardware compatibility, which posed challenges in rural and remote areas with limited infrastructure.

In proposed study [5] a deep learning model for multiclass classification of hip fractures, covering femoral neck, intertrochanteric, and subtrochanteric types. The model demon-

TABLE I
SUMMARY OF STUDIES ON HIP FRACTURE DETECTION USING DEEP LEARNING APPROACHES

Ref	Title	Key Findings	Limitations
[1]	Application of Convolutional Neural Networks for Hip Fracture Detection on X-ray Images	CNN-based models achieve diagnostic accuracy comparable to experienced radiologists, reducing human error and providing faster interpretations in fracture diagnosis.	Relatively small dataset from a single institution, limiting generalizability to broader populations and different imaging settings.
[2]	Effectiveness of Ensemble Learning Methods in Medical Image Classification Tasks	Combining multiple deep learning architectures into an ensemble enhances classification accuracy and produces more stable outputs compared to single CNNs, improving robustness in medical imaging.	Training and implementation require significantly higher computational resources, challenging in low-resource healthcare systems.
[3]	Deep Learning for Automated Detection and Classification of Femoral Neck and Intertrochanteric Fractures on Hip X-rays	AI system achieves high sensitivity and specificity, providing rapid and standardized interpretations that assist radiologists and reduce diagnostic variability.	Accuracy declines in poor-quality images or anatomical variations, highlighting the need for large and diverse datasets.
[4]	AI-Based Diagnostic Systems for Hip Fracture Detection in Low-Resource Healthcare Settings	AI-driven platforms provide real-time hip fracture detection and bridge diagnostic gaps, supporting frontline providers without on-site expert radiologists.	Systems depend on internet connectivity and hardware compatibility, posing challenges in rural and remote areas.
[5]	Deep Learning Model for Multiclass Classification of Hip Fractures	Multiclass classification allows for more precise and clinically relevant interpretation of hip fractures compared to binary approaches.	Misclassifications occur in borderline or overlapping fracture cases, suggesting ensemble learning for enhanced reliability.

strated improved accuracy compared to binary classification approaches, showing its usefulness in handling the complexity of multiple fracture categories. The key finding was that multiclass classification allowed for a more precise and clinically relevant interpretation of hip fractures. The limitation was that misclassifications often occurred in borderline or overlapping fracture cases, suggesting that ensemble learning could be necessary to enhance diagnostic reliability further.

Table I illustrated the Existing studies related to Hip region fracture.

III. METHODOLOGY

The research methodology focuses on the development of a multi-class classification of the hip X-ray image to detect fractures or abnormalities, and the pelvis is the part of the body under consideration. This process is described by gathering a pool of hip X-ray images, which undergo the preprocessing of data, including augmentation, model construction, and the model is trained and tested to measure its effectiveness in detecting. The process would entail successive corrections under the results of an evaluation process to optimize the model for the image of the hip region. Figure 2 shows the high-level architecture.

A. Data Collection

The collection of the data with the hip X-ray images was performed based on the X-ray images of several hospitals in various geographical regions, which made a varied population of the data in terms of object detection. It was also divided into various fractures levels to allow for the provision of a complete picture of pelvic conditions to be used during training and evaluation of the results. Tables II and III show the data collection methods and types of hip X-ray images. Figure 3 shows the types of hip x-ray images.

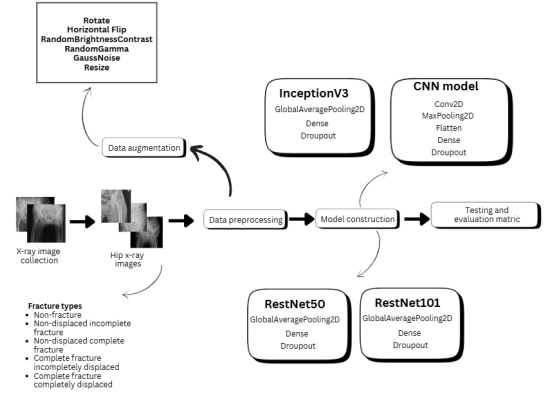


Fig. 2. High level architecture

TABLE II
DATASET DISTRIBUTION ACROSS HOSPITALS

Hospital	No. of Images
Base Hospital Tellipalai, Jaffna	450
Northern Central Hospital, Jaffna	800
Teaching Hospital, Batticaloa	500
Aathura Hospital - Baily Rd, Batticaloa	450
Venus Specialty Hospital Pvt Ltd	300

TABLE III
DISTRIBUTION OF HIP FRACTURE TYPES

Fracture Type	No. of Images
Non-Fractured	534
Non-displaced incomplete fracture	136
Non-displaced complete fracture	91
Complete fracture incompletely displaced	93
Complete fracture completely displaced	116

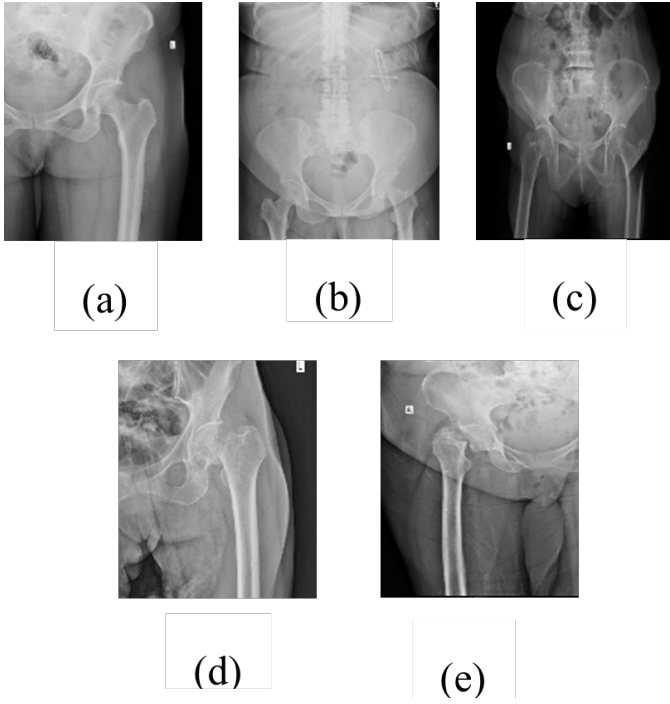


Fig. 3. Hip neck fracture x-ray images types: (a) Non-fracture (b) Non-displaced incomplete fracture (c) Non-displaced complete fracture (d) Complete fracture incompletely displaced (e) Complete fracture completely displaced

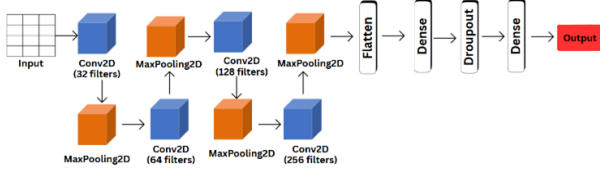


Fig. 4. Ensemble models architecture

B. Model Construction

The study incorporates a suite of deep learning models (Figure 4), Custom CNN, InceptionV3, ResNet50, and ResNet101 tailored for the multiclass classification of hip region fractures in X-ray images (Table IV). The Custom CNN is fully trainable, featuring Conv2D layers with filter sizes of 32, 64, 128, and 256 to capture detailed bone and fracture features, followed by MaxPooling2D for spatial reduction, a Dense layer with 512 units and ReLU activation, a 0.5 Dropout rate to mitigate overfitting, and a Softmax output layer for classifying five fracture types (non-fracture, non-displaced incomplete, non-displaced complete, complete incompletely displaced, complete completely displaced). InceptionV3 uses a frozen pre-trained base from ImageNet, employing its inception modules for multi-scale feature extraction, followed by GlobalAveragePooling2D, a Dense layer (512, ReLU), 0.5 Dropout, and a Softmax output. ResNet50 and ResNet101, also with frozen pre-trained bases from ImageNet, utilize residual connections for deep feature learning, each with GlobalAveragePooling2D, a Dense layer (512, ReLU), 0.5

Dropout, and a Softmax output, where ResNet101's deeper architecture enhances fracture detail detection.

IV. RESULTS AND DISCUSSION

The performance of the proposed models Custom CNN, InceptionV3, ResNet50, and ResNet101 in detecting and classifying hip region fractures is presented in Table V, with metrics including training accuracy, validation accuracy, testing accuracy, training time (in seconds), precision, recall, and F1-score. The Custom CNN achieved a training accuracy of 0.8333, validation accuracy of 0.6763, and testing accuracy of 0.6643, with a training time of 762.40 seconds, but showed lower precision (0.5530), recall, and F1-score, indicating faster training at the cost of accuracy due to its fully trainable architecture. InceptionV3 recorded a training accuracy of 0.7129, validation accuracy of 0.6763, testing accuracy of 0.6857, and a training time of 795.97 seconds, with precision at 0.6643, recall at 0.5833, and F1-score at 0.5553, reflecting moderate performance with a frozen pre-trained base. ResNet50 outperformed with a training accuracy of 0.8919, validation accuracy of 0.7194, testing accuracy of 0.7786, and a training time of 1237.22 seconds, achieving precision of 0.6857, recall of 0.5906, and a high F1-score of 0.7298, demonstrating improved generalization from its residual structure. ResNet101 achieved the highest testing accuracy of 0.8000, with a training accuracy of 0.8735, validation accuracy of 0.7338, and a longer training time of 1817.57 seconds, alongside the best precision (0.7785), recall (0.7364), and F1-score (0.8179), highlighting its deeper architecture's advantage in capturing complex fracture features.

The results indicate that deeper models like ResNet101 and ResNet50 generally outperform the Custom CNN and InceptionV3, with ResNet101 leading due to its enhanced feature extraction capabilities, as supported by its higher testing accuracy and F1-score. The longer training time of ResNet101 (1817.57 seconds) compared to Custom CNN (762.40 seconds) reflects the trade-off between computational cost and accuracy, a finding consistent with prior studies which noted increased training durations for deeper pre-trained models in pelvic fracture detection. The Custom CNN's lower testing accuracy (0.6643) suggests it struggles with generalization, likely due to its reliance on a smaller, fully trainable architecture without the benefit of pre-trained weights. InceptionV3's moderate performance (0.6857 testing accuracy) aligns with its balanced design but highlights limitations in handling the diverse fracture patterns compared to ResNet variants. The precision and recall trends show ResNet101's superior balance, making it the most reliable for clinical use, though its extended training time may necessitate optimization for real-time applications. Future work could explore ensemble techniques or attention mechanisms to further enhance accuracy and efficiency, addressing the observed variability across models.

The performance of the ResNet101 model is further illustrated through its accuracy and loss graphs over 10 epochs, as depicted in Figures 5. The ResNet101 Accuracy graph shows training accuracy increasing steadily from approximately 0.60

TABLE IV
DEEP LEARNING MODELS AND THEIR ARCHITECTURES

Model	Trainable Parameters	Key Convolutional Base Layers	Pooling Layer	Dense / Output Layers
Custom CNN	All layers trainable	Conv2D (32, 64, 128, 256)	MaxPooling2D	Dense(512, ReLU), Dropout(0.5), Dense(Softmax)
InceptionV3	Base frozen	InceptionV3 (pretrained)	GlobalAveragePooling2D	Dense(512, ReLU), Dropout(0.5), Dense(Softmax)
ResNet50	Base frozen	ResNet50 (pretrained)	GlobalAveragePooling2D	Dense(512, ReLU), Dropout(0.5), Dense(Softmax)
ResNet101	Base frozen	ResNet101 (pretrained)	GlobalAveragePooling2D	Dense(512, ReLU), Dropout(0.5), Dense(Softmax)

TABLE V
PERFORMANCE COMPARISON OF DEEP LEARNING MODELS FOR HIP FRACTURE CLASSIFICATION

Model	Training Accuracy	Validation Accuracy	Testing Accuracy	Training Time (s)	Precision	Recall	F1-Score
CNN	0.8333	0.6763	0.6643	762.40	0.5530	—	—
InceptionV3	0.7129	0.6763	0.6857	795.97	0.6643	0.5833	0.5553
ResNet50	0.8919	0.7194	0.7786	1237.22	0.6857	0.5906	0.7298
ResNet101	0.8735	0.7338	0.8000	1817.57	0.7785	0.7364	0.8179

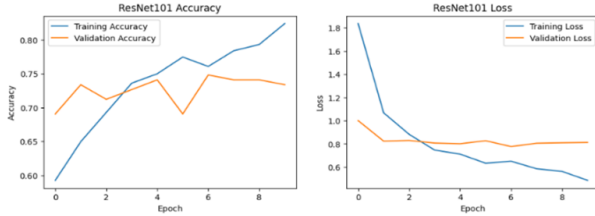


Fig. 5. ResNet101 model accuracy and loss graph

to 0.80, while validation accuracy rises from 0.65 to around 0.75, indicating a consistent improvement in model learning with a slight gap suggesting moderate overfitting. The ResNet101 Loss graph complements this, with training loss decreasing sharply from 1.8 to 0.6 and validation loss dropping from 1.4 to 0.8, stabilizing after epoch 4, which reflects effective convergence and a good fit to the hip fracture dataset.

V. CONCLUSION

The study demonstrates that the proposed deep learning models—Custom CNN, InceptionV3, ResNet50, and ResNet101 offer a robust framework for the multiclass classification of hip region fractures, with ResNet101 achieving the highest testing accuracy of 0.8000 and an F1-score of 0.8179, outperforming other models as validated by accuracy and loss trends over 10 epochs. These results, inspired by the hip fracture detection highlight the efficacy of ensemble architectures in enhancing diagnostic precision for fracture types including non-fracture, non-displaced incomplete, non-displaced complete, complete incompletely displaced, and complete completely displaced. The integration of data augmentation and preprocessing, as depicted in the architecture diagram, has proven effective in addressing imaging variability, with ResNet101's deeper structure providing a balance between accuracy and computational cost despite its longer training time of 1817.57 seconds.

Future work will focus on improving model efficiency and generalizability by exploring ensemble learning techniques, such as stacking or voting, to combine the strengths of these models for even higher accuracy and robustness. Additionally,

incorporating attention mechanisms and expanding the dataset with diverse imaging conditions from multiple institutions will enhance real-time clinical applicability. Further research will also investigate the integration of metadata and conduct prospective clinical trials to validate the models' performance in dynamic healthcare settings, building on the foundational insights from the current study and related literature.

REFERENCES

- [1] B. Babu and R. Khan, "Object detection—a comparison between pre-trained and custom model," 2023.
- [2] J. Bae, S. Yu, J. Oh, T. Kim, J. Chung, H. Byun, M. Yoon, C. Ahn, and D. Lee, "External validation of deep learning algorithm for detecting and visualizing femoral neck fracture including displaced and non-displaced fracture on plain x-ray," *Journal of Digital Imaging*, vol. 34, no. 5, pp. 1099–1109, 2021.
- [3] C.-T. Cheng, Y. Wang, H.-W. Chen, P.-M. Hsiao, C.-N. Yeh, C.-H. Hsieh, S. Miao, J. Xiao, C.-H. Liao, and L. Lu, "A scalable physician-level deep learning algorithm detects universal trauma on pelvic radiographs," *Nature Communications*, vol. 12, no. 1, p. 1066, 2021.
- [4] "Fractures," *Journal of Medical Imaging and Radiation Oncology*, vol. 63, no. 1, pp. 27–32, 2019.
- [5] H. Hashmi, R. Dwivedi, and A. Kumar, "Comparative analysis of cnn-based smart pre-trained models for object detection on dota," *Journal of Automation, Mobile Robotics and Intelligent Systems*, pp. 31–45, 2024.
- [6] G. Kitamura, "Deep learning evaluation of pelvic radiographs for position, hardware presence, and fracture detection," *European Journal of Radiology*, vol. 130, p. 109139, 2020.
- [7] C.-W. Kuo and Z. Kira, "Beyond a pre-trained object detector: Cross-modal textual and visual context for image captioning," in *Book Beyond a pre-trained object detector: Cross-modal textual and visual context for image captioning*, 2022, pp. 17969–17979.
- [8] Y. Sharab, M. Alsmira, Z. Dwekat, I. Alsmadi, and A. Al-Khasawneh, "Performance comparison of several deep learning-based object detection algorithms utilizing thermal images," in *Book Performance comparison of several deep learning-based object detection algorithms utilizing thermal images*. IEEE, 2021, pp. 16–22.
- [9] J. Wu, P. Davuluri, K. R. Ward, C. Cockrell, R. Hobson, and K. Najarian, "Fracture detection in traumatic pelvic ct images," *International Journal of Biomedical Imaging*, vol. 2012, no. 1, p. 327198, 2012.
- [10] D. Yadav, A. Sharma, S. Athithan, A. Bholra, B. Sharma, and I. Dhaou, "Hybrid sfnet model for bone fracture detection and classification using ml/dl," *Sensors*, vol. 22, no. 15, p. 5823, 2022.
- [11] S. Sanchez, H. Romero, and A. Morales, "A review: Comparison of performance metrics of pretrained models for object detection using the tensorflow framework," in *Book A review: Comparison of performance metrics of pretrained models for object detection using the TensorFlow framework*. IOP Publishing, 2020, p. 012024.