

Undergraduate Research Opportunities Programme in Science (UROPS)  
Project Report

# **Detecting bone fractures with attention based CNN and VLM**

By  
Wong Song Ming

Department of Statistics and Data Science  
Faculty of Science  
National University of Singapore

2025/2026

Undergraduate Research Opportunities Programme in Science (UROPS)  
Project Report

# **Detecting bone fractures with attention based CNN and VLM**

By  
Wong Song Ming

Department of Statistics and Data Science  
Faculty of Science  
National University of Singapore

2025/2026

Project ID: UXXXXX

Advisor: Professor Swapnil Mishra

Deliverables:

Report: 1 Volume

## Abstract

Elbow fractures are a common type of bone fracture among children, and can lead to serious complications if not quickly diagnosed and treated. Deep learning models such as convolutional neural networks have proven effective at detecting elbow fractures, matching the performance of trained radiologists. However, the performance of these models is hampered by a lack of training data. In this project, we developed a framework to train attention based CNN models to detect elbow fractures. The framework was tested with 3 CNNs: ConvNeXt, EfficientNetV2 and MobileNetV3. In addition, we finetuned 2 medical vision language models, MedGemma and LLaVA-Rad, and evaluated their performance against the CNNs. The CNN models used were modified with the addition of convolutional block attention modules and pretrained on the MURA dataset before training on the target dataset. Subsequently, feature fusion was used to train an ensemble model. The vision language models were finetuned on the same training dataset using Low Rank Adaptation. The attention based CNN models showed a significant improvement over the original models, ranging from 1.5% to 3.1% improvement in accuracy while the ensemble model achieved the highest recall of 94.3% and second highest precision of 91.9%. Medgemma had poor performance with an accuracy of only 88.1%, below that of the attention based CNN models. However LLaVA-Rad had an accuracy of 91.5% which surpassed 2 of the CNN models and had a high recall of 93.8%, demonstrating high potential for further work.

## **Acknowledgements**

I would like to thank Professor Swapnil Mishra for giving me the opportunity to take on this research project and providing me with the guidance I needed to complete the project. I would also like to thank my family and friends for supporting and encouraging me.

# Contents

<b>Abstract</b>	<b>i</b>
<b>Acknowledgements</b>	<b>ii</b>
<b>List of Figures</b>	<b>iv</b>
<b>List of Tables</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Related Work</b>	<b>3</b>
<b>3 Materials and Methods</b>	<b>5</b>
3.1 Dataset . . . . .	5
3.2 Data Processing . . . . .	6
3.3 CNN Models . . . . .	7
3.4 Vision Language Models . . . . .	12
<b>4 Evaluation</b>	<b>14</b>
4.1 Results . . . . .	14
4.2 Discussion . . . . .	19
<b>5 Conclusion</b>	<b>21</b>
<b>References</b>	<b>22</b>

# List of Figures

3.1	VGG16 architecture   Credits to Kenneth Leung . . . . .	7
3.2	Model architecture . . . . .	10
4.1	Confusion matrices for CNN models . . . . .	16
4.2	Confusion matrices for CNN models without pretraining and attention . . . . .	17
4.3	Confusion matrices for VLM models . . . . .	18

# List of Tables

3.1	Training set . . . . .	5
3.2	Pretraining set . . . . .	5
4.1	CNN results . . . . .	15
4.2	CNN results without pretraining and attention . . . . .	15
4.3	VLM results . . . . .	15

# Chapter 1

## Introduction

Bone fractures are a significant health concern globally with 178 million new cases and 455 million prevalent cases of acute or long term symptoms in 2019 alone. (Wu et al., 2021) Elbow fractures in particular are one of the most common types of bone fractures among children, accounting for over 10% of paediatric fractures and up to 5% of fractures among adults. Accurate and timely diagnosis of elbow fractures is important to avoid potential complications such as joint stiffness, malunion and vascular injury. (Waseem et al., 2025) However, interpretation of elbow X-rays is a difficult task. Some parts of the anatomy of the elbow are obscured depending on the perspective of the X-ray (anteroposterior, lateral or oblique) and non displaced elbow fractures are difficult to spot on X-rays (McGinley et al., 2006). Diagnosis of paediatric elbow fractures is especially difficult due to the differences in anatomy and injury patterns from the adult elbow such as the presence of ossification centres (Iyer et al., 2012).

Recent advances in machine learning and deep learning have allowed for the use of automated systems to quickly detect fractures from X-rays. These often rely on convolutional neural networks (CNN) which are able to automatically learn complex features and patterns from images and have been shown to have excellent performance in fracture detection, on par with or exceeding trained radiologists (Jung et al., 2024). The use of these automated systems could help support doctors in accurately diagnosing elbow fractures, especially in emergency settings and for paediatric elbow fractures which may be initially reviewed by radiologists who do not specialise in paediatrics (Lindsey et al., 2018; Taves et al., 2017).

One challenge faced by CNN models as well as other deep learning models is the amount of available data which is often limited in medical domains. To address this, deep learning models can be pretrained on other datasets which are related to the target dataset. In the case of diagnosing elbow fractures, this can include pretraining on X-rays of other bones such as the shoulder, or even X-rays of other body parts such as chest X-rays (Alammar et al., 2023).



Other techniques to improve the performance of deep learning models include the use of attention mechanisms which allow the model to focus on important parts of an image and ensembling techniques to combine the output of multiple models.

Other than CNN models, another area that shows potential is the use of vision language foundation models. Foundation models are pretrained on huge multimodal datasets which contain both image and text and cover many different domains. This allows the models to perform a wide variety of tasks with or without task specific finetuning. Foundation models developed for general purposes may not perform well in the medical domain due to the lack of medical data in their training sets. However, there have also been foundation models which have been specifically trained for medical domains to perform tasks such as generating medical reports or segmenting images (Khan et al., 2025). Although foundation models are much larger and more computationally expensive than CNNs, the same model could potentially be used for many different medical tasks including those which combine image and text data which is a significant advantage over a CNN which is more narrow.

This project aims to improve on the generic CNN architecture by developing a framework to train attention based CNN models to accurately detect elbow fractures from X-rays. We also aim to finetune pretrained medical vision language foundation models on the same dataset to compare against the CNN models and evaluate their performance.

# Chapter 2

## Related Work

Luo et al. (2021) proposed a multiview deep learning method to classify elbow fractures, evaluated using a dataset of 1964 elbow radiographs which were divided into non fracture, ulnar fracture and radial fracture classes. In the single view stage, 2 Vgg16 models were trained on images in either the frontal view or the lateral view. The trained model weights were then reused in the frontal and lateral view modules of the multiview model. After feature extraction, the multiview model was split into 2 single view branches and one which combined features from both views. If both views were available for the same patient, the output of the merge branch was taken as the prediction. Otherwise, only the relevant branch would be used. The training process used knowledge-guided curriculum learning which involved feeding easier samples before harder samples as quantified by radiologists. The frontal view and lateral view models had a balanced accuracy of 57% and 80.7% and a binary accuracy of 73.2% and 89.5% respectively. The multiview model had a balanced accuracy of 86.4% and binary accuracy of 91%.

Malik et al. (2022) proposed a 2 phase method to identify complex fractures using a subset of 16984 images from the MURA dataset. In the first phase, the images were preprocessed and converted into RGB format. In the second phase, 2 pretrained deep learning models (Darknet-53 and Xception) were used to extract features from the images. These were combined with features extracted using histogram of oriented gradient and local binary pattern. The resulting feature vector then underwent feature selection using principal component analysis followed by the whale optimisation algorithm. Lastly, the 1049 selected features were used to train an SVM, K-NN and neural network classifier. Under 10-fold cross validation, the SVM classifier achieved 91.4% accuracy and 0.82 kappa score, the K-NN classifier achieved 97.1% accuracy and 0.94 kappa score and the neural network achieved 86.5% accuracy and 0.73 kappa score.

Ahmed and Hawezi (2023) investigated the use of traditional machine learning techniques to classify bone fractures. A dataset of 270 x-rays of the lower leg was used for testing. First, the images were converted to greyscale and processed using a gaussian filter for denoising as well as adaptive histogram equalisation and canny edge detection to improve contrast. Next, grey level co-occurrence matrix (GLCM) was used to extract 5 texture properties (energy, correlation, dissimilarity, homogeneity, contrast) over 4 distances and 7 angles for a total of 140 features per image. The GLCM features were then used to train 5 machine learning classifiers (naive bayes, decision tree, K-NN, random forest, SVM). Naive bayes had an accuracy of 64.2%, decision tree had an accuracy of 80.3%, K-NN had an accuracy of 83.9%, random forest had an accuracy of 85.7% and SVM had an accuracy of 92.9%.

Alzubaidi et al. (2024) developed a trustworthy deep learning framework using the MURA dataset to detect shoulder abnormalities. The shoulder X-rays were used as the target dataset and X-rays of the remaining 6 body parts were used for pretraining. A total of 7 individual deep learning models were pretrained on MURA followed by finetuning on the shoulder X-rays. The features of the individual models were then extracted at the fully connected layers and combined through feature fusion and used to train several classifiers such as logistic regression and SVM. The best classifier achieved an accuracy of 99.2% with pretraining and 78.5% without pretraining which was a significant improvement over the individual models which ranged from 72.4% to 77.6% accuracy and 65.7% to 72.6% accuracy respectively. This surpassed the performance of 3 surgeons invited to classify the dataset who had an average accuracy of 79.1%. The individual deep learning models were further validated using activation visualisation and locally interpretable model-independent explanations for interpretability of outputs.

Tahir et al. (2024) proposed an ensemble of 4 CNN models (MobileNetV2, Vgg16, InceptionV3, ResNet50) to detect bone fractures. The ensemble consisted of a logistic regression model trained on the output probabilities of the individual models. The dataset used was the subset of 6542 humerus radiographs from the MURA dataset. The images were processed with histogram equalisation to improve contrast followed by data augmentation using random rotations, horizontal flipping and random scaling. On the validation set, the 4 individual models achieved an accuracy of 88%, 82.2%, 81% and 86% respectively while the ensemble model had an accuracy of 93%.

Alam et al. (2025) introduced a model called MobLG-Net to extract features for training machine learning classifiers. The method was tested on a publicly available dataset on Kaggle consisting of 9463 X-rays of various body parts. MobLG-Net consists of the pretrained input layer of MobileNet combined with a custom sequential model with convolutional layers, pooling, dropout and fully connected layers. The features extracted by MobLG-Net were then used to train several classifiers such as light gradient boosting machine and logistic regression. The classifiers trained using this approach had an average accuracy of 97.6% to 98.5%.

# Chapter 3

## Materials and Methods

### 3.1 Dataset

The main dataset used in this project was collected from a local hospital and consists of 4369 elbow radiographs in either the lateral or anterior posterior view. Each radiograph was named according to the type of injury present (supracondylar fracture, dislocation, contusion etc). Using this information, the data was grouped into either injury or normal. The injury set includes radiographs with untreated fractures and/or dislocations while the normal set consists of all other radiographs. The data was further divided into training, validation, and test sets with a 80/10/10 split as shown in Table 3.1.

In addition, the MURA dataset from Stanford Machine Learning Group was used for pretraining the CNN models (Rajpurkar et al., 2017). This dataset consists of 40561 radiographs of 7 different bodyparts: elbow, finger, forearm, hand, humerus, shoulder and wrist. Each image is labelled as either positive (abnormality) or negative (no abnormality). The division of the data into training and validation sets is shown in Table 3.2. Although MURA also has elbow images, the criteria for abnormality used in MURA differs from the criteria for injury in the main dataset. For example, a contusion would be classified as normal in the main dataset but would be classified as abnormal in the MURA dataset. Because of this mismatch, the elbow X-rays in MURA could not be combined with the main dataset for the training stage.

Set	Normal	Injury	Both
Training	1957	1538	3495
Validation	251	186	437
Test	244	193	437
Total	2452	1917	4369

Table 3.1: Training set

Set	Normal	Injury	Both
Training	21935	14870	36805
Validation	1533	1667	3200
Total	23468	16537	40005

Table 3.2: Pretraining set

## 3.2 Data Processing

To improve the quality of the radiographs, all of the radiographs underwent preprocessing. Firstly, the images were resized according to the model used. For the CNN models, the image size used for pretraining and training was 384 x 384 and 584 x 584 pixels respectively. As the models are fully convolutional, they could accept any image size without modifications to the architecture. For MedGemma and LLaVA-Rad, the images were resized to 896 x 896 and 518 x 518 respectively according to the requirements for each model.

Next, the radiographs underwent Contrast Limited Adaptive Histogram Equalisation (CLAHE) using OpenCV. Adaptive Histogram Equalisation is an image processing technique which calculates histograms of the pixel intensities in different parts of an image. The pixel intensities are then redistributed using the histograms to improve the contrast of the image (Reza, 2004). The implementation in OpenCV splits the image into a non overlapping grid and applies the technique subject to a clip limit which limits the amount of redistribution to reduce noise. The specific settings used were a clip limit of 2.5 and a grid size of 4. CLAHE is useful for detecting fractures as it increases the contrast between a fracture and the surrounding bone which improves its visibility.

Next, the radiographs were sharpened further using unsharp masking. First, each image was processed with a gaussian filter to produce a blurred image. This blurred image was then subtracted from the original image to produce a mask which is multiplied by a factor of 2 and added back to the original image. This process sharpens the image and further increases contrast at the edges. The kernel used for the gaussian filter was 5 x 5 for the MURA dataset and 7 x 7 for the main dataset.

Lastly, the data was augmented with a combination of geometric transformations as follows: Random rotation between  $-60^\circ$  and  $60^\circ$ , 75% probability of horizontal flip, scaling by a factor between 0.85 and 1.1 and translation along the x and y axis of between 0.02 and 0.12. Data augmentation artificially expands the training set for supervised learning which can reduce overfitting and improve generalisation to new data. Each image in the training set was augmented twice for the target dataset and once for the MURA dataset for a total of 10485 images and 73610 images respectively. Brightness and contrast transformations are also commonly used in image augmentation but were found to negatively impact performance and were thus excluded.

## 3.3 CNN Models

### 3.3.1 Architecture

Convolutional Neural Networks (CNNs) are a type of neural network which use convolutional filters or kernels to perform convolutions on images and produce intermediate feature maps. These feature maps are then passed onto the following convolutional layers. At the final layer, the feature maps are pooled and flattened to produce a feature vector which is processed by a fully connected neural network to generate an output. An example of a CNN architecture is shown below.

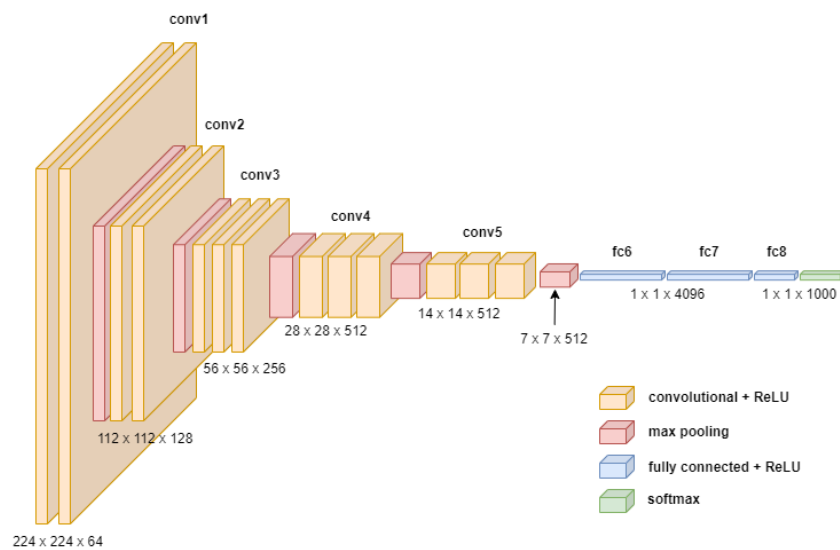


Figure 3.1: VGG16 architecture | Credits to Kenneth Leung

For this project, 3 CNN models were used: ConvNeXt, EfficientNetV2 and MobileNetV3. ConvNeXt is a modernised ResNet or Residual Network which is a deep neural network architecture that is widely used in tasks involving medical imaging due to its high performance (Liu et al., 2022; Xu et al., 2023). EfficientNetV2 is another family of CNNs which are designed for efficiency and faster training and have been successful in tasks such as identifying breast cancer from histopathology slides (Tan and Le, 2021; Hayat et al., 2024). MobileNetV3 is a lightweight CNN architecture which can be run on mobile devices due to its small size and low memory consumption and is useful for real time medical applications (Howard et al., 2019).

After training the individual models, an ensemble model was created using feature fusion for the final classification. The framework used for the CNN models is outlined in the following pages.

- 1. Transfer Learning

Transfer learning is a machine learning technique where a machine learning model is trained on one task and is then reused for another related task instead of starting again from scratch. This technique not only saves on training time but can also improve performance on the second task, especially if the amount of available training data is limited.

In this case, the CNN models were initialised with weights which had been pretrained on ImageNet 1k. ImageNet 1k is a large dataset with over a million RGB radiographs divided into 1000 categories such as plants, animals and furniture and is commonly used as a source of transfer learning for image classification tasks. The weights for the models were obtained from Pytorch Hub which is a publicly available model repository.

- 2. Greyscale vs RGB

ImageNet is a colour dataset with 3 channel (RGB) radiographs. Hence, the models cannot directly accept greyscale radiographs which only have 1 channel. This issue can be solved by duplicating the single channel twice to create a 3 channel image. However, this process is inefficient as it increases the memory cost of storing the image as well as the computational cost of processing the image through the model. To solve this, the first convolutional layer of each model was modified by summing the per channel weights of each kernel. This allows the models to directly accept greyscale radiographs without affecting the output of the first layer.

- 3. Convolutional Block Attention Module

Convolutional Block Attention Module (CBAM) is an attention mechanism for CNN models which can be incorporated into an existing architecture and trained together with the rest of the model (Woo et al., 2018). CBAM takes in an intermediate feature map and infers a channel attention map and a spatial attention map. The intermediate feature map is then refined using element wise multiplication with the 2 attention maps. CBAM has been shown to improve the performance of CNNs on different computer vision tasks while incurring little additional computational cost.

ConvNeXt was modified with 4 CBAM modules placed after each of the 4 ConvNeXt blocks. EfficientNet was modified with a single CBAM module between the FusedMBConv modules and the MBConv modules, as well as a second CBAM module after the final convolutional layer. MobileNet was modified with a single CBAM module after the final convolutional layer.

- 4. Ensembling

Ensembling is a commonly used technique in machine learning which combines the output of several different base models to get a single model which outperforms the individual ones. This can be done through several methods such as a simple majority vote or training a model on the outputs of the individual models.

In this case, feature fusion was used to combine the outputs of the 3 individual CNN models. This is done by concatenating the flattened feature vectors of each model before they are sent to the fully connected layer. By default, Convnext has 768 output channels while EfficientNet has 1280 output channels and MobileNet has 960 output channels. To ensure each model has equal input to the ensemble, the final convolutional layer of EfficientNet was modified to have 768 output channels and MobileNet had an additional convolutional layer added to reduce the channels from 960 to 768.

The concatenated feature vector with 2304 features is then normalised and passed to a dropout layer which randomly sets 30% of the features to 0 to mitigate overfitting. The feature vector is then used as the input for a small neural network to generate the final prediction.



Input Image (584 x 584 x 1)
Conv2d (146 x 146 x 96)
Layernorm (146 x 146 x 96)
ConvNeXt block (73 x 73 x 192)
CBAM (73 x 73 x 192)
Downsample (73 x 73 x 192)
ConvNeXt block (36 x 36 x 384)
CBAM (36 x 36 x 384)
Downsample (36 x 36 x 384)
ConvNeXt block (18 x 18 x 768)
CBAM (18 x 18 x 768)
Downsample (18 x 18 x 768)
ConvNeXt block (18 x 18 x 768)
CBAM (18 x 18 x 768)
Adaptive Average Pooling (1 x 1 x 768)
Layernorm (1 x 1 x 768)
Flatten
Linear

(a) Convnext architecture

Input Image (584 x 584 x 1)
Conv2d (292 x 292 x 24)
Batchnorm (292 x 292 x 24)
FusedMBConv (3x) (292 x 292 x 24)
FusedMBConv (5x) (146 x 146 x 48)
FusedMBConv (5x) (73 x 73 x 80)
CBAM (73 x 73 x 80)
MBConv (7x) (37 x 37 x 160)
MBConv (14x) (37 x 37 x 176)
MBConv (18x) (19 x 19 x 384)
MBConv (5x) (19 x 19 x 512)
Conv2d (19 x 19 x 768)
Batchnorm (19 x 19 x 768)
CBAM (19 x 19 x 768)
Adaptive Average Pooling (1 x 1 x 768)
Flatten
Linear

(b) EfficientNet architecture

Input Image (584 x 584 x 1)
Conv2d (292 x 292 x 16)
Batchnorm (292 x 292 x 16)
Inverted Residual (292 x 292 x 16)
Inverted Residual (2x) (146 x 146 x 24)
Inverted Residual (3x) (73 x 73 x 40)
Inverted Residual (4x) (37 x 37 x 80)
Inverted Residual (2x) (37 x 37 x 112)
Inverted Residual (3x) (19 x 19 x 160)
Conv2d (19 x 19 x 960)
Batchnorm (19 x 19 x 960)
Conv2d (19 x 19 x 768)
Batchnorm (19 x 19 x 768)
CBAM (19 x 19 x 768)
Adaptive Average Pooling (1 x 1 x 768)
Flatten
Linear

(c) MobileNet architecture

Figure 3.2: Model architecture

### 3.3.2 CNN Training

As mentioned previously, ImageNet is a dataset of natural objects and is very different from the greyscale radiographs in our target dataset. Hence, it is not an ideal source of transfer learning as the features learnt from ImageNet may not be useful. To address this issue, the MURA dataset was used for a pretraining step before training on the target dataset.

For pretraining, the Stochastic Gradient Descent (SGD) optimiser was used together with a cosine annealing learning rate scheduler with warm restarts. The training hyperparameters were batch size of 32, max epochs of 35, initial learning rate of  $1.2e-3$ ,  $8e-4$  and  $1e-4$  for ConvNeXt, EfficientNet and MobileNet respectively, momentum of 0.6 and weight decay of  $1e-3$ . The cosine annealing scheduler was set to  $T_0$  of 5 and  $T_{mult}$  of 2 which means the learning rate is decayed over 5 epochs in the first cycle, 10 epochs in the second cycle and 20 in the third.

Following pretraining, the pretrained weights of the 3 models were used for training on the target dataset. The hyperparameters for training were largely the same as pretraining. The differences were max epochs of 28, momentum of 0.8 and 0.7 for ConvNeXt and EfficientNet respectively and  $T_0$  of 4. Lastly, for the ensemble model, the different hyperparameters were max epochs of 12, momentum of 0.8 and initial learning rate of  $2e-4$  for the fully connected layers and  $5e-5$  for the backbone.

## 3.4 Vision Language Models

### 3.4.1 Models

Large language models (LLM) are a type of generative deep learning model which use transformers to learn statistical patterns in natural language and predict the next word in a sequence by training on massive text datasets. LLMs break down input text into tokens which are represented by numerical vectors called embeddings. The self attention mechanism of the transformers generates query, key and value vectors from the embeddings using trainable weights. The attention score is then computed from the dot product of the query and key vectors and determines which tokens are more important. This process allows LLMs to capture the relationship between words and generate coherent outputs (Raiaan et al, 2024).

Vision language models (VLM) are multimodal generative models which expand the capabilities of an LLM by combining it with a vision encoder which is usually a vision transformer model. The vision encoder processes image inputs into an embedding vector which has the same dimensions as the embedding generated by the LLM from the text input. The embeddings are then fused into a single embedding which combines both visual and textual information. This embedding is then passed on to the remainder of the LLM to generate a text output. This allows VLMs to perform tasks involving image and/or text such as image captioning, question answering and image classification (Zhang et al, 2024).

Foundation models are large models which have been pretrained on a vast amount of data, usually through self supervised learning. Foundation models are versatile and can be used for a wide variety of tasks or as a base for further finetuning. In the case of VLMs, the training data includes unlabelled text and images or labelled text-image pairs (Awais et al, 2025). However, the datasets used for general foundation models usually do not have much medical data which impacts the performance of these models on medical tasks. To address this issue, foundation VLMs have been specifically trained on medical datasets for use in tasks such as medical report generation and diagnosis. For this project, the 2 VLMs used are MedGemma 4B and LLaVA-Rad due to their high performance and relatively small size.

MedGemma is a collection of 2 publicly available foundation models created by google and trained for a wide range of medical related tasks such as radiology report generation and diagnosis (Selligren et al., 2025). MedGemma is based on the Gemma 3 VLM and has a 4 billion parameter and 27 billion parameter variant. Both use the SigLIP image encoder and have been trained on a large variety of medical data such as chest X-rays, histopathology slides and medical question answer pairs. LLaVA-Rad is another publicly available foundation model created by Microsoft and focused on analysing chest X-rays and generating reports (Zambrano Chaves et al., 2025). LLaVA-Rad uses the BiomedCLIP image encoder and the Vicuna LLM and is trained on a large dataset of chest X-rays with radiology reports.

### 3.4.2 VLM Training

Finetuning the parameters of an entire VLM is extremely computationally and memory intensive due to the large size of the model and also carries the risk of catastrophic forgetting which is a phenomenon where a model loses performance on old tasks when finetuned for a new task (Goodfellow et al., 2015). To mitigate these issues, parameter efficient finetuning (PEFT) strategies are often used instead (Bafghi et al., 2024). Rather than finetuning the entire model, only a small number of additional parameters are trained.

The specific PEFT strategy used here is a technique called Low Rank Adaptation (LoRA) which involves freezing the parameters of a pretrained model and injecting additional low rank matrices into the targeted layers of the model (Hu et al., 2021). The updates of the original weight matrices from training can be represented by the product of a pair of low rank matrices which are much smaller than the original when the rank is small. During training, the original matrices are frozen and only the parameters of the low rank matrices are optimised which greatly decreases the memory and time required. For inference, the product of the low rank matrices is added to the original weight matrices to obtain the finetuned weights without significantly increasing inference time. This process can be easily reversed by subtraction to obtain the original weights allowing the model to swap between LoRA weights to perform different tasks.

As the target dataset consists solely of images, a prompt was added to each image to instruct the model to analyse the radiograph and classify it as either A: 'No injury' or B: 'Fracture or dislocation'. The correct answer is either 'No injury' or 'Fracture or dislocation' depending on the class of the image. For both MedGemma and LLaVA-Rad, the pretrained weights, tokeniser and image processor were downloaded from Huggingface and used as a starting point for finetuning. The checkpoint used for LLaVA-Rad was from the alignment step which only tuned the image encoder. This is because the final finetuning step used LoRA instead of directly training the image encoder and LLM. To further reduce memory usage, the existing weights of the model were compressed with 4-bit quantisation using the bitsandbytes library (Dettmers et al., 2023).

The TRL and PEFT libraries were used in a custom training script for finetuning MedGemma while LLaVA-Rad was finetuned by modifying the provided script from the LLaVA-Rad repository. For both models, the weights of the linear layers of the model were targeted for finetuning. The Adam-W adaptive optimiser was used instead of SGD as with the CNN models to increase the convergence speed. For MedGemma, the hyperparameters used were batch size of 32 divided into 8 accumulation steps, max epochs of 5, learning rate of  $8e-5$  with linear scheduling, weight decay of  $1e-3$ , max grad norm of 0.03, LoRA rank of 16 and alpha of 32. For LLaVA-Rad, the different hyperparameters were max epochs of 4, learning rate of  $1e-4$ , LoRA rank of 32 and alpha of 64.

# Chapter 4

## Evaluation

### 4.1 Results

#### 4.1.1 Evaluation Metrics

The models were evaluated using accuracy, precision, recall, F1 score and AUROC (for the CNN models only). The first 4 metrics are calculated based on the true positive (TP), true negative (TN), false positive (FP) and false negative (FN) values at a particular prediction threshold. AUROC is calculated from a plot of the sensitivity (recall) of a model against specificity at different prediction thresholds. The equations are as follows:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (4.1)$$

$$Precision = \frac{TP}{TP + FP} \quad (4.2)$$

$$Recall = \frac{TP}{TP + FN} \quad (4.3)$$

$$F1\ score = \frac{2 * TP}{2 * TP + FP + FN} \quad (4.4)$$

$$Specificity = \frac{TN}{TN + FP} \quad (4.5)$$

### 4.1.2 Results

The 4 CNN models were evaluated at a balanced threshold of 0.5 for all metrics except AUROC. Convnext achieved an accuracy of 93.4%, precision of 92.7%, recall of 92.2%, F1 score of 92.5% and AUROC of 0.976. EfficientNet achieved an accuracy of 90.1%, precision of 90.3%, recall of 87.0%, F1 score of 88.7% and AUROC of 0.971. MobileNet achieved an accuracy of 90.6%, precision of 86.5%, recall of 93.2%, F1 score of 89.8% and AUROC of 0.953. The Ensemble model achieved an accuracy of 93.8%, precision of 91.9%, recall of 94.3%, F1 score of 93.1% and AUROC of 0.978.

Model	Accuracy	Precision	Recall	F1 score	AUC-ROC
ConvNeXt	0.934	0.927	0.922	0.925	0.976
EfficientNet	0.901	0.903	0.870	0.887	0.971
MobileNet	0.906	0.865	0.932	0.898	0.953
Ensemble	0.938	0.919	0.943	0.931	0.978

Table 4.1: CNN results

Without pretraining or any changes to the architecture from the publicly available version, Convnext achieved an accuracy of 90.6%, precision of 87.3%, recall of 92.2%, F1 score of 89.7% and AUROC of 0.945. EfficientNet achieved an accuracy of 87.0%, precision of 83.7%, recall of 87.6%, F1 score of 85.6% and AUROC of 0.941. MobileNet achieved an accuracy of 89.2%, precision of 86.1%, recall of 90.2%, F1 score of 88.1% and AUROC of 0.935.

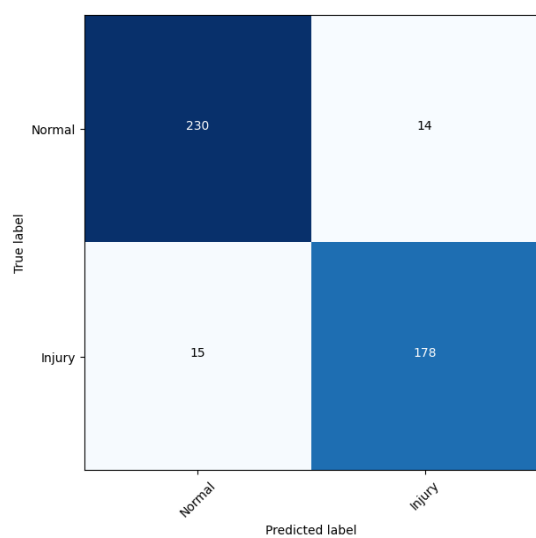
Model	Accuracy	Precision	Recall	F1 score	AUC-ROC
ConvNeXt	0.906	0.873	0.922	0.897	0.945
EfficientNet	0.870	0.837	0.876	0.856	0.941
MobileNet	0.892	0.861	0.902	0.881	0.935

Table 4.2: CNN results without pretraining and attention

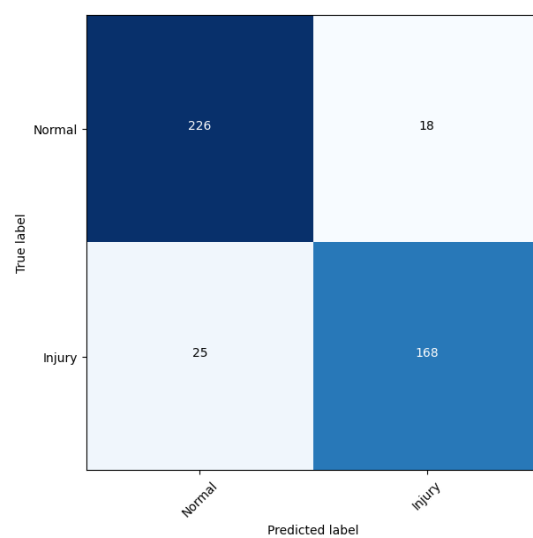
For MedGemma and LLaVA-Rad, the predictions are obtained from the generated text output and does not have a threshold like the CNN models. Hence, AUC-ROC could not be calculated. MedGemma achieved an accuracy of 88.1%, precision of 86.2%, recall of 87.0% and F1 score of 86.6%. LLaVA-Rad achieved an accuracy of 91.5%, precision of 87.9%, recall of 93.8% and F1 score of 90.7%.

Model	Accuracy	Precision	Recall	F1 score	AUC-ROC
MedGemma	0.881	0.862	0.870	0.866	-
LLaVA-Rad	0.915	0.879	0.938	0.907	-

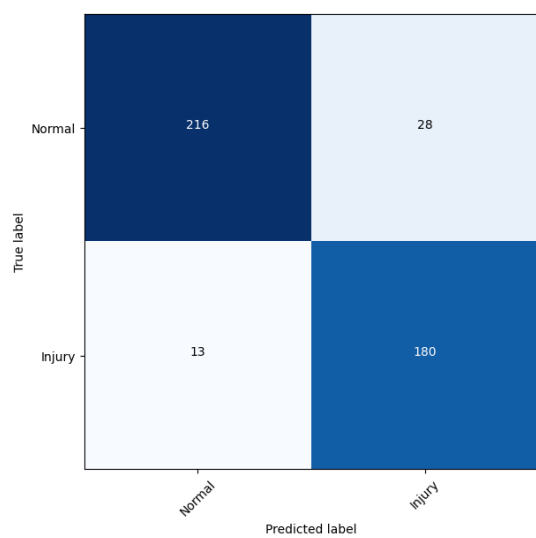
Table 4.3: VLM results



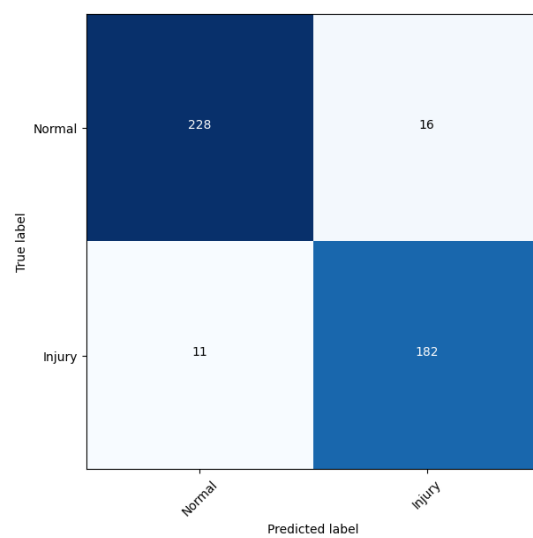
(a) Convnext



(b) EfficientNet

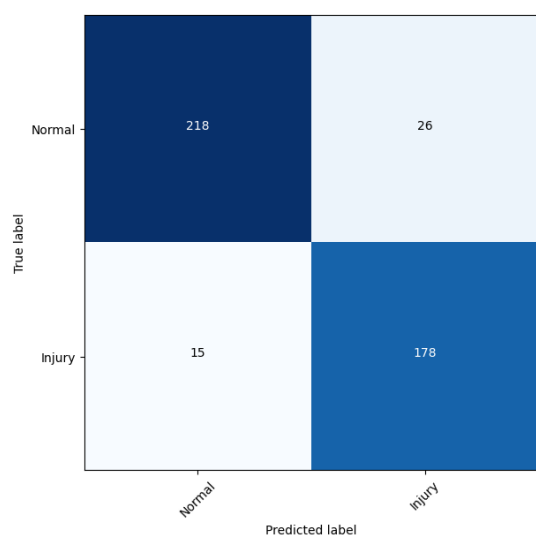


(c) MobileNet

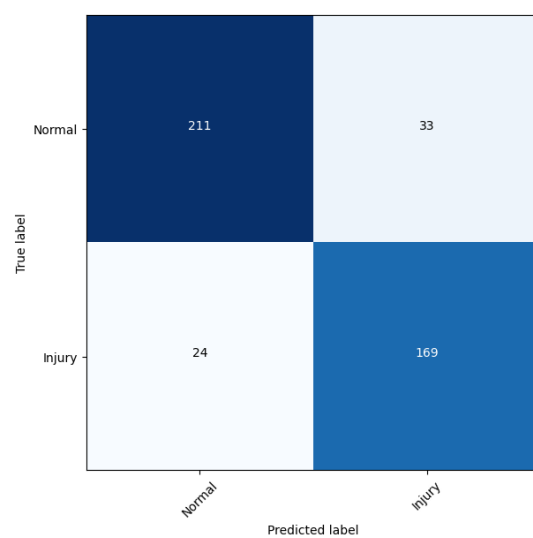


(d) Ensemble

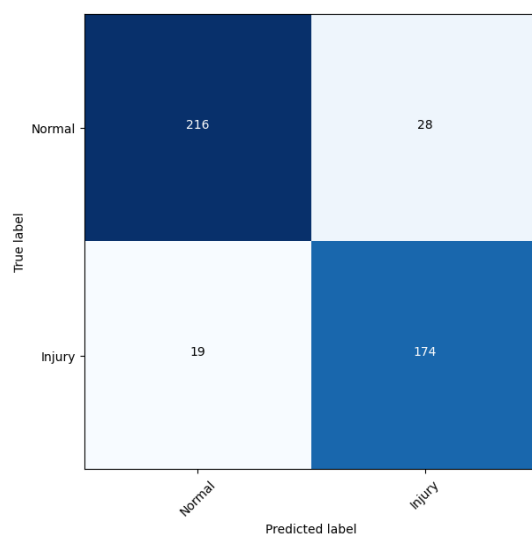
Figure 4.1: Confusion matrices for CNN models



(a) Convnext



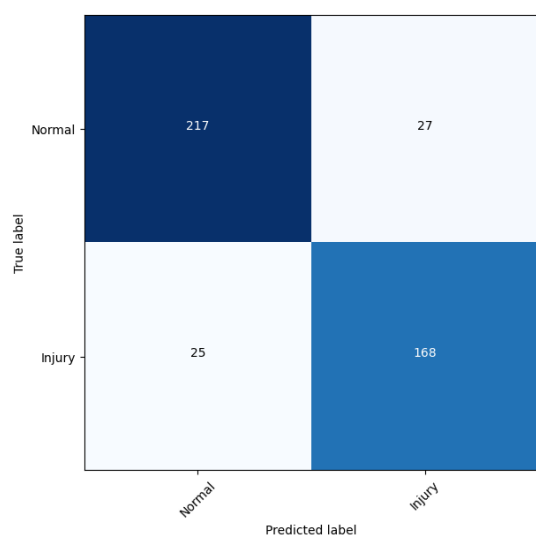
(b) EfficientNet



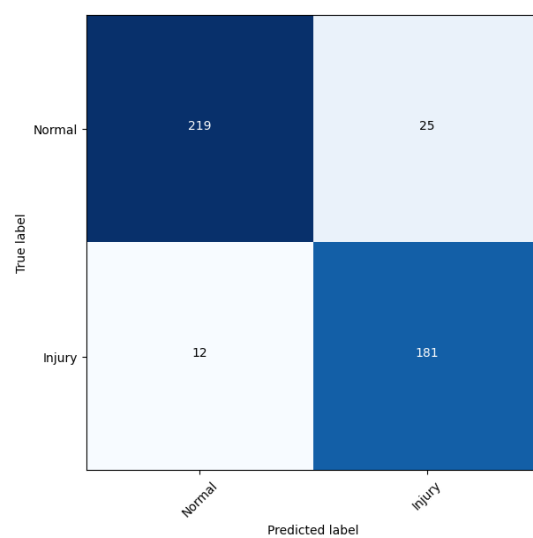
(c) MobileNet

Figure 4.2: Confusion matrices for CNN models without pretraining and attention





(a) Medgemma



(b) LLaVA-Rad

Figure 4.3: Confusion matrices for VLM models

## 4.2 Discussion

From the results, the combination of pretraining on a related dataset and the changes made to the architectures of the CNN models had a significant impact on performance. Convnext had a 2.8% increase in accuracy and 0.031 increase in AUROC, EfficientNet had a 3.1% increase in accuracy and 0.03 increase in AUROC, MobileNet had a 1.4% increase in accuracy and 0.018 increase in AUROC. This indicates that pretraining on the MURA dataset helped to reduce the mismatch issue between ImageNet and the target dataset. The addition of convolutional block attention modules also improved the ability of the models to focus on relevant parts of the image thus improving performance. The improvement was larger for Convnext and EfficientNet than MobileNet which may be due to the larger size of the first 2 models.

In Alzubaidi et al., pretraining yielded improvements of between 2% and 9% accuracy over using just ImageNet which is much larger than the improvement seen here. This may be because of differences in the target dataset. In Alzubaidi et al, the target dataset was the subset of humerus X-rays in MURA which comes from the same source as the rest of the dataset. In contrast, the elbow X-rays used here come from a different source and have different characteristics from the X-rays in the MURA dataset which reduces the impact of the pretraining step. A possible solution could be to supplement the pretraining set with additional X-rays from other publicly available online sources as well as from local healthcare institutions. This could also include unlabelled X-rays for an additional self-supervised learning step before the pretraining.

The ensemble model had a 0.4% increase in accuracy and 0.002 increase in AUROC compared to Convnext which is the best performing individual model. This is a relatively small difference in performance. However, the ensemble model had a higher recall than all 3 of the individual models with 2.1% higher recall than Convnext, 7.3% higher recall than EfficientNet and 1.1% higher recall than MobileNet. The model also had a higher precision than EfficientNet and MobileNet although it had a 0.8% lower precision than Convnext.

The ensemble models in Alzubaidi et al. and Tahir et al. showed an improvement of over 20% and 2% accuracy respectively over their best performing individual model. This is much larger than the improvement seen in our ensemble model. The comparatively small improvement suggests that the features collected from the 3 individual models may be too similar to each other resulting in a smaller benefit from the feature fusion. This might be mitigated by using a larger number of smaller CNN models or by training each model on a slightly different subset of the training data.

For the VLM models, MedGemma generally had a worse performance than any of the individual CNN models, only beating EfficientNet without pretraining. In contrast, LLaVA-Rad had a better accuracy and F1 score than EfficientNet and MobileNet, only losing

to ConvNeXt and the Ensemble model. In addition, LLaVA-Rad had the second highest recall out of all the models. In the case of automated fracture detection, the cost of a false negative is significantly higher than the cost of a false positive.

As mentioned in the introduction, an undiagnosed elbow fracture could lead to further complications for a patient whereas a false diagnosis wastes the time of the radiologist which is a less serious problem. Hence, a high recall is arguably more important than a high precision. From this standpoint, the Ensemble was the best model followed by LLaVA-Rad. In addition because of the computational requirements, the VLMs were not pretrained on the MURA dataset which might have increased their performance significantly. Hence, although the performance of the VLMs are not the best, LLaVA-Rad in particular shows potential for further development.

# Chapter 5

## Conclusion

In this project, we developed a framework for the training of attention based CNN models to detect elbow fractures from X-rays. The framework was validated by training 3 attention based CNN models as well as an ensemble model using feature fusion. The results showed that the framework achieved a significant improvement over the baseline models. In addition, we finetuned 2 VLMs, MedGemma and LLaVA-Rad, on the same dataset and found that LLaVA-Rad had competitive performance compared to the attention based CNNs. Lastly, we suggested possible ways to improve on the performance of the CNNs and VLMs.

# References

- Ahmed, K. D., & Hawezi, R. (2023). Detection of bone fracture based on machine learning techniques. *Measurement: Sensors*, 27. <https://doi.org/10.1016/j.measen.2023.100723>
- Alam, A., Al-Shamayleh, A. S., Thalji, N., Raza, A., Barajas, E., Thompson, E. B., Diez, I. D. L. T., & Ashraf, I. (2025). Novel transfer learning based bone fracture detection using radiographic images. *BMC Medical Imaging*, 25. <https://doi.org/10.1186/s12880-024-01546-4>
- Alammar, Z., Alzubaidi, L., Zhang, J., Li, Y., Lafta, W., & Gu, Y. (2023). Deep transfer learning with enhanced feature fusion for detection of abnormalities in x-ray images. *Cancers*, 15. <https://doi.org/10.3390/cancers15154007>
- Alzubaidi, L., Salhi, A., A Fadhel, M., Bai, J., Hollman, F., Italia, K., Pareyon, R., Albahri, A. S., Ouyang, C., Santamaría, J., Cutbush, K., Gupta, A., Abbosh, A., & Gu, Y. (2024). Trustworthy deep learning framework for the detection of abnormalities in x-ray shoulder images. *PLOS One*, 19. <https://doi.org/10.1371/journal.pone.0299545>
- Awais, M., Naseer, M., Khan, S., Anwer, R. M., Cholakkal, H., Shah, M., Y, M.-H., & K, F. S. (2025). Foundational models defining a new era in vision: A survey and outlook. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47, 2245–2264. <https://doi.org/10.1109/TPAMI.2024.3506283>
- Bafghi, R. A., Harilal, N., Monteleoni, C., & Raissi, M. (2024). Parameter efficient fine-tuning of self-supervised vits without catastrophic forgetting. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 3679–3684. <https://doi.org/10.1109/CVPRW63382.2024.00371>
- Dettmers, T., Pagnoni, A., Holtzman, A., & Zettlemoyer, L. (2023). Qlora: Efficient finetuning of quantized llms. <https://doi.org/10.48550/arXiv.2305.14314>
- Goodfellow, I. J., Mirza, M., Da, X., Courville, A., & Bengio, Y. (2015). An empirical investigation of catastrophic forgetting in gradient-based neural networks. <https://doi.org/10.48550/arXiv.1312.6211>
- Hayat, M., Ahmad, N., Nasir, A., & Zeeshan, A. T. (2024). Hybrid deep learning efficientnetv2 and vision transformer (effnetv2-vit) model for breast cancer histopathological image classification. *IEEE Access*, 12, 184119–184131. <https://doi.org/10.1109/ACCESS.2024.3503413>

- Howard, A., Sandler, M., Chen, B., Wang, W., Chen, L.-C., Tan, M., Chu, G., Vasudevan, V., Zhu, Y., Pang, R., Adam, H., & Le, Q. (2019). Searching for mobilenetv3. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. <https://doi.org/10.1109/ICCV.2019.00140>
- Hu, E. J., Shen, Y., Wallis, P., Zeyuan, A.-Z., Li, Y., Wang, S., Wang, L., & Chen, W. (2021). Lora: Low-rank adaptation of large language models. <https://doi.org/doi.org/10.48550/arXiv.2106.09685>
- Iyer, R. S., Thapa, M. M., Khanna, P. C., & Chew, F. S. (2012). Pediatric bone imaging: Imaging elbow trauma in children—a review of acute and chronic injuries. *American Journal of Roentgenology*, 198. <https://doi.org/10.2214/AJR.10.7314>
- Jung, J., Dai, J., Liu, B., & Wu, Q. (2024). Artificial intelligence in fracture detection with different image modalities and data types: A systematic review and meta-analysis. *PLOS Digital Health*, 3(1). <https://doi.org/10.1371/journal.pdig.0000438>
- Khan, W., Leem, S., See, K. B., Wong, J. K., Zhang, S., & Fang, R. (2025). A comprehensive survey of foundation models in medicine. *IEEE Reviews in Biomedical Engineering*. <https://doi.org/10.1109/RBME.2025.3531360>
- Lindsey, R., Daluiski, A., Chopra, S., Lachapelle, A., Mozer, M., Sicular, S., Hanel, D., Gardner, M., Gupta, A., Hotchkiss, R., & Potter, H. (2018). Deep neural network improves fracture detection by clinicians. *Proceedings of the National Academy of Sciences of the United States of America*, 115. <https://doi.org/10.1073/pnas.1806905115>
- Luo, J., Kitamura, G., Arefan, D., Doganay, E., Panigrahy, A., & Wu, S. (2021). Knowledge-guided multiview deep curriculum learning for elbow fracture classification. *Machine learning in medical imaging*. [https://doi.org/10.1007/978-3-030-87589-3\\_57](https://doi.org/10.1007/978-3-030-87589-3_57)
- Malik, S., Amin, J., Sharif, M., Yasmin, M., Kadry, S., & Anjum, S. (2022). Fractured elbow classification using hand-crafted and deep feature fusion and selection based on whale optimization approach. *Mathematics*, 10. <https://doi.org/10.3390/math10183291>
- McGinley, J. C., Roach, N., Hopgood, B. C., & Kozin, S. H. (2006). Nondisplaced elbow fractures: A commonly occurring and difficult diagnosis. *The American Journal of Emergency Medicine*, 24, 560–566. <https://doi.org/10.1016/j.ajem.2006.01.010>
- Raiaan, M. A. K., Mukta, S. H., Fatema, K., Fahad, N. M., Sakib, S., & Mim, M. M. J. (2024). A review on large language models: Architectures, applications, taxonomies, open issues and challenges. *IEEE Access*, 12, 26839–26874. <https://doi.org/10.1109/ACCESS.2024.3365742>
- Rajpurkar, P., Irvin, J., Bagul, A., Ding, D., Duan, T., Mehta, H., Yang, B., Zhu, K., Laird, D., Ball, R. L., Langlotz, C., Shpanskaya, K., Lungren, M. P., & Ng, A. Y. (2017). Mura: Large dataset for abnormality detection in musculoskeletal radiographs. <https://doi.org/10.48550/arXiv.1712.06957>
- Reza, A. M. (2004). Realization of the contrast limited adaptive histogram equalization (clahe) for real-time image enhancement. *The Journal of VLSI Signal Processing-Systems for*

- Signal, Image, and Video Technology*, 38, 38–44. <https://doi.org/10.1023/B:VLSI.0000028532.53893.82>
- Sellergren, A., Kazemzadeh, S., Jaroensri, T., Kiraly, A., Traverse, M., Kohlberger, T., Xu, S., Jamil, F., Hughes, C., Lau, C., Chen, J., Mahvar, F., Yatziv, L., Chen, T., Sterling, B., Baby, S., Baby, S. M., Lai, J., Schmidgall, S., & ... Yang, L. (2025). Medgemma technical report. <https://doi.org/doi.org/10.48550/arXiv.2507.05201>
- Tahir, A., Saadia, A., Khan, K., Gul, A., Qahmash, A., & Akram, R. N. (2024). Enhancing diagnosis: Ensemble deep-learning model for fracture detection using x-ray images. *Clinical Radiology*, 79, e1394–e1402. <https://doi.org/10.1016/j.crad.2024.08.006>
- Tan, M., & Le, Q. V. (2021). Efficientnetv2: Smaller models and faster training. *International Conference on Machine Learning*. <https://api.semanticscholar.org/CorpusID:232478903>
- Taves, J., Skitch, S., & Valani, R. (2017). Determining the clinical significance of errors in pediatric radiograph interpretation between emergency physicians and radiologists. *Canadian Journal of Emergency Medicine*, 20, 420–424. <https://doi.org/10.1017/cem.2017.34>
- Waseem, M., Saeed, W., & Launico, M. V. (2025, July). *Elbow fractures overview*. StatPearls Publishing. <https://www.ncbi.nlm.nih.gov/books/NBK441976/>
- Woo, S., Park, J., Lee, J.-L., & Kweon, I. S. (2018). Cbam: Convolutional block attention module. *Computer Vision – ECCV 2018: 15th European Conference, Munich, Germany, September 8–14, 2018, Proceedings, Part VII*, 3–19. [https://doi.org/10.1007/978-3-030-01234-2\\_1](https://doi.org/10.1007/978-3-030-01234-2_1)
- Wu, A. M., Bisignano, C., James, S., Abady, G. G., Abedi, A., Abu-Gharbieh, E., Alhassan, R. K., Alipour, V., Arabloo, J., Asaad, M., Asmare, V. N., Awedew, A. F., Banach, M., Banerjee, S. K., Bijani, A., Birhanu, T. T. M., Bolla, S. R., Cámara, L. A., Chang, J. C., & ... Vos, T. (2021). Global, regional, and national burden of bone fractures in 204 countries and territories, 1990-2019: A systematic analysis from the global burden of disease study 2019. *The Lancet Healthy Longevity*, 2, e580–e592. [https://doi.org/10.1016/S2666-7568\(21\)00172-0](https://doi.org/10.1016/S2666-7568(21)00172-0)
- Xu, W., Fu, Y.-L., & Zhu, D. (2023). Resnet and its application to medical image processing: Research progress and challenges. *Computer Methods and Programs in Biomedicine*, 240. <https://doi.org/doi.org/10.1016/j.cmpb.2023.107660>
- Zambrano Chaves, J. M., Huang, S.-C., Xu, Y., Xu, H., Usuyama, N., Zhang, S., Wang, F., Xie, Y., Khademi, M., Yang, Z., Awadalla, H., Gong, J., Hu, H., Yang, J., Li, C., Gao, J., Gu, Y., Wong, C., Wei, M., & ... Poon, H. (2025). A clinically accessible small multimodal radiology model and evaluation metric for chest x-ray findings. *Nature Communications*, 16(3108). <https://doi.org/10.1038/s41467-025-58344-x>

- Zhang, J., Huang, J., Jin, S., & Lu, S. (2024). Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46, 5625–5644. <https://doi.org/10.1109/TPAMI.2024.3369699>
- Zhuang, L., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., & Xie, S. (2022). A convnet for the 2020s. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 11966–11976. <https://doi.org/10.1109/CVPR52688.2022.01167>