# "Bone Age Prediction"
# Deep CNN Models for Predicting Bone Age from Hand Radiographs

Atefe Rostami[†]

*Abstract*—**Bone age prediction from hand radiographs is a critical task in medical imaging that helps determine the maturity of a patient's bones. This information is valuable for various medical applications, including growth assessment, diagnosis of endocrine disorders, and monitoring treatment efficacy. In this paper, we present a comparative study of two deep learning architectures for bone age prediction: DenseNet121 and ResNet34 with Convolutional Block Attention Module (CBAM). Our implementation focuses on creating an efficient data processing pipeline, implementing robust model architectures, and utilizing advanced training techniques such as cosine learning rate scheduling with warm-up. The models were trained on a dataset of 12,611 training images, 1,425 validation images, and 200 test images. Both architectures achieve competitive performance, with DenseNet121 achieving an MAE of 8.18 months and an R² score of 0.93, while ResNet34+CBAM achieved an MAE of 9.77 months and an R² score of 0.90 on the test set, demonstrating the effectiveness of both dense connectivity and attention mechanisms in medical image analysis tasks.**

*Index Terms*—**Supervised Learning, Convolutional Neural Networks, Image Regression, DenseNet121, ResNet34, CBAM, Attention Mechanisms.**

## I. INTRODUCTION

Bone age assessment from hand radiographs is a crucial tool in pediatric and endocrine practice: it helps clinicians evaluate growth disorders, plan treatments, and estimate the timing of puberty. Traditional techniques—most notably the Greulich–Pyle and Tanner–Whitehouse atlases—require an expert to visually compare a patient's X-ray with standard reference images. While widely used, these manual methods are time-consuming and suffer from inter- and intra-observer variability.

Recent advances in deep learning promise to automate and standardize bone age estimation. In this work, we implement and compare two state-of-the-art convolutional neural network (CNN) models:

- **DenseNet121.** A densely connected CNN in which each layer receives input from all preceding layers, promoting feature reuse and efficient gradient flow.
- **ResNet34 + CBAM.** A residual network augmented with a Convolutional Block Attention Module (CBAM) that adaptively emphasizes informative channels and spatial regions.

Our pipeline incorporates several enhancements aimed at robust performance:

† atefe.Rostami@studenti.unipd.it

- **Data preprocessing:** Detect and crop hands, apply CLAHE enhancement, resize with aspect ratio preservation to $256 \times 256$, normalize to $[0, 1]$, and convert to grayscale.
- **Training strategies:** Cosine-annealed learning rates with warmup, early stopping on validation MAE, and dropout regularization.
- **Auxiliary inputs:** Concatenate a binary gender indicator (0=female, 1=male) with the image features.
- **Evaluation metrics:** Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and $R^2$ on a held-out test set.
- **Efficient data loading:** TensorFlow's `tf.data` API with caching and prefetching.

To illustrate our dataset, Figures 1, 2, and 3 show the distribution of bone ages overall, by gender, and the male/female counts.
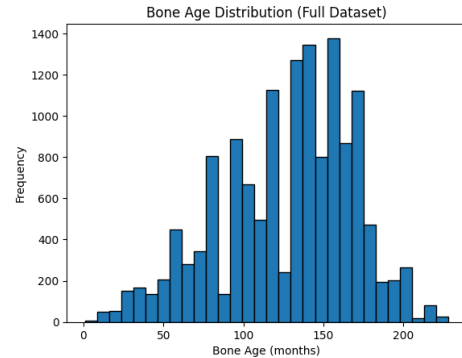


Fig. 1. Overall distribution of bone age in the full dataset.

Despite the promise of CNNs, bone age prediction remains challenging: images of children with the same chronological age can display wide variations in skeletal maturity, and manual labels exhibit some variability. Automated methods must therefore learn to accommodate this natural heterogeneity.

Clinically, accurate bone age estimation is used for:

- *Normal growth monitoring:* Identifying delays or accelerations in skeletal development.
- *Endocrine disorders:* Diagnosing and tracking conditions like growth hormone deficiency or precocious puberty.
- *Genetic syndromes:* Assessing children with Down syndrome, Turner syndrome, and other chromosomal abnormalities.
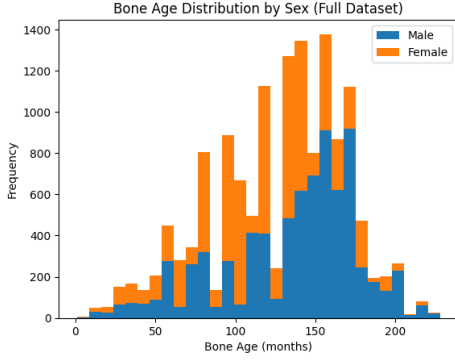
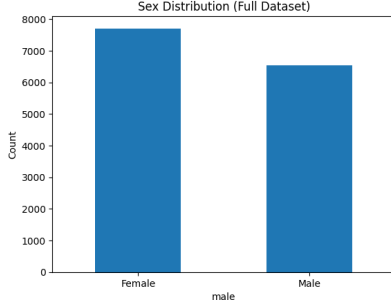Fig. 2. Bone age distribution stratified by gender.



Fig. 3. Counts of male and female subjects in the dataset.

**Objectives.** In this report, we will:

1) Explore data-augmentation and preprocessing strategies to maximize useful signal in the images.
2) Compare the predictive performance of DenseNet121 and ResNet34+CBAM on the same dataset.
3) Analyze model errors using MAE, RMSE, and $R^2$ on a held-out test set.
4) Discuss practical considerations for deploying these models in a clinical workflow.

The remainder of this report is organized as follows. Section II reviews prior work. Section III details our data pipeline and model architectures, drawing on the DenseNet121 and ResNet34+CBAM schematics. Section IV describes the dataset. Section VI presents quantitative results and error analysis. Finally, Section VII offers conclusions and future directions.

## II. RELATED WORK

Deep convolutional neural networks (DCNNs) have rapidly advanced the state of the art in bone age estimation from hand radiographs. Early approaches trained bespoke CNNs end-to-end on large labeled datasets, using standard convolutional and pooling layers followed by fully-connected regression heads. For example, Rajpurkar *et al.* (2017) collected a dataset of over 12000 radiographs and trained a custom CNN from scratch, leveraging extensive data augmentation (rotation, scaling) and transfer learning from natural-image networks to achieve radiologist-level accuracy.

To improve information flow and gradient propagation, later studies adopted *dense connectivity*. A prominent example

is DenseNet121, in which each layer receives as input the feature-maps from all preceding layers, promoting feature reuse and mitigating vanishing gradients. In our implementation, the DenseNet121 backbone consists of four dense blocks (6, 12, 24, 16 convolutional layers respectively) interleaved with transition layers that compress and downsample feature maps . This architecture has been shown to extract richer multi-scale features compared to traditional CNNs, especially when fine-tuned on medical datasets.

Residual networks with attention mechanisms represent another successful line of work. ResNet34 provides a strong baseline via its identity-skip connections, which ease optimization of very deep models. By integrating a Convolutional Block Attention Module (CBAM) into each residual block, models can learn to focus on the most informative channels and spatial regions of the hand radiograph. Our ResNet34+CBAM variant applies channel attention, using global average and max pooling followed by a shared MLP—and spatial attention—using a 7×7 convolution over concatenated pooling maps. This attention-augmented ResNet34 architecture has demonstrated improved localization of key ossification centers and enhanced regression accuracy.

In summary, our work builds upon these advances by (1) adopting DenseNet121 and ResNet34+CBAM backbones (2) augmenting them with a simple binary gender input stream, and (3) employing cosine-Decay learning-rate schedules, early stopping, and rigorous evaluation (MAE, RMSE, $R^2$). This allows a direct comparison of dense versus attention-augmented residual designs under a unified training pipeline.

## III. PROCESSING PIPELINE

### A. Data Preprocessing

The preprocessing pipeline was designed to prepare high-resolution hand radiographs for robust model training and evaluation. It was implemented in two stages: (1) local offline preprocessing using OpenCV and MediaPipe; and (2) dynamic input pipeline construction using TensorFlow on the Kaggle platform.

*1) Offline Hand Image Preprocessing:* All raw images were preprocessed locally using the following techniques to enhance anatomical focus and contrast:

- **Hand Detection and Cropping:** A custom hand detection method using MediaPipe landmarks was applied to automatically crop the region of interest, excluding non-relevant background and focusing on hand bones.
- **CLAHE Enhancement:** Contrast Limited Adaptive Histogram Equalization (CLAHE) was used on the L-channel of LAB color space to improve bone structure visibility while preserving local contrast.
- **Grayscale Conversion:** Optionally, RGB images were converted to single-channel grayscale to reduce model complexity while retaining key skeletal features.
- **Aspect Ratio-Preserving Resizing:** Images were zero-padded to form square shapes before being resized to $256 \times 256$ pixels, preserving anatomical proportions.
- **Normalization:** Pixel values were scaled to the $[0, 1]$ range to standardize inputs for the model.

Processed images were saved in a dedicated directory and used for training and evaluation on Kaggle.

*2) Dataset Organization:* The dataset was structured into three non-overlapping subsets:

- **Training Set:** 12,611 images with labels for model fitting
- **Validation Set:** 1,425 images for early stopping and hyperparameter tuning
- **Test Set:** 200 images reserved for final model evaluation

All label files were normalized to a unified schema: [case_id, gender, age], where gender is binary (1=male, 0=female) and age is in months. Corresponding DataFrames were constructed to map preprocessed image paths to metadata.

### B. Data Pipeline Implementation

On the Kaggle platform, an efficient and scalable input pipeline was implemented using TensorFlow's tf.data.Dataset API with the following key components:

- **Image Loader:** A generator function reads PNG images from disk and decodes them as grayscale tensors with shape [H, W, 1].
- **Augmentation & Normalization:**
  - Detects and crops hand regions using MediaPipe landmarks
  - Applies CLAHE enhancement to single-channel grayscale
  - Zero-pads and resizes images to $256 \times 256$ while preserving aspect ratio
  - Normalizes pixel values to the $[0, 1]$ range (float32)
  - Casts gender and age labels to float32 tensors
  - Ensures all images are resized and rescaled consistently
- **Label Splitting:** Each sample is represented as a tuple of inputs and labels: ((image, gender), age), enabling joint modeling of both visual and metadata inputs.
- **Pipeline Configuration:**
  - Caching: Applied to both training and validation pipelines for improved performance
  - Shuffling: Performed with full dataset cardinality to ensure randomness
  - Batching: Configured with a batch size of 64
  - Prefetching: Utilized AUTOTUNE for parallelism
  - Repeat Mode: Training data was set to repeat indefinitely; validation/test datasets were not

## IV. IMAGES AND FEATURES

The dataset consists of hand X-ray images collected from multiple medical centers, with the following distribution:

- Training set: 12,611 images
- Validation set: 1,425 images
- Test set: 200 images

Key characteristics of the dataset include:

- Image dimensions standardized to 256×256 pixels
- Single-channel grayscale format
- Pixel values normalized to [0,1] range
- Age labels in months (continuous values)
- Binary gender labels (0=female, 1=male)



GENDER: MALE, AGE: 143.31317138671875

Fig. 4. Sample hand radiograph after cropping and contrast enhancement

## V. LEARNING FRAMEWORK

We implemented and compared two different deep learning architectures for bone age prediction:
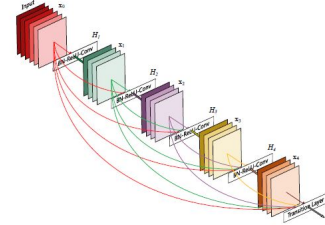
### A. DenseNet121



Fig. 5. Basic architecture of DenseNet121

DenseNet121 is characterized by its dense connectivity pattern, where each layer is connected to every other layer in a feed-forward fashion. Our implementation includes:

- **Dense Blocks:**
  - Four dense blocks with 6, 12, 24, and 16 layers respectively
  - Growth rate of 32 channels per layer
  - Each layer consists of BN-GELU-Conv(1×1)-BN-GELU-Conv(3×3)
  - L2 regularization with weight decay of 1e-4 on convolutions
- **Transition Blocks:**
  - Batch normalization followed by GELU activation
  - 1×1 convolution with compression factor of 0.5
  - Average pooling with 2×2 kernel and stride 2
  - Used between dense blocks for dimensionality reduction
- **Gender Integration:**
  - Two-layer MLP for gender processing: Dense(16)-GELU-Dense(8)

- Concatenation with global average pooled image features
- Final regression head: Dense(32)-Dense(16)-Dropout(0.01)-Dense(1)
- L2 regularization with weight decay of 0.01 on dense layers

## B. ResNet34 with CBAM

Our second architecture enhances the ResNet34 backbone with Convolutional Block Attention Module (CBAM) for improved feature refinement:

- **Basic Blocks**:
  - Four block groups with [3,4,6,3] residual blocks
  - Channel progression: 64→128→256→512
  - Each block: Conv(3×3)-BN-GELU-Conv(3×3)-BN-CBAM
  - Identity shortcuts with 1×1 convolution when dimensions change

- **Channel Attention Module**:
  - Shared MLP with reduction ratio of 16
  - Parallel global average and max pooling
  - Element-wise summation of pooled features
  - Sigmoid activation for attention weights

- **Spatial Attention Module**:
  - Channel-wise average and max pooling
  - Concatenation of pooled features
  - 7×7 convolution with sigmoid activation
  - Applied after channel attention

- **Classification Head**:
  - Global average pooling of final features
  - Concatenation with gender input
  - Dense(1024)-Dropout(0.4)-Dense(512)-Dropout(0.4)-Dense(1)
  - Linear activation for final regression output

Both architectures were trained with the following common elements:

- Input shape: 256×256×1 (grayscale images)
- Gender input: Binary value (0/1)
- Batch size: 64
- Learning rate scheduling with cosine decay and warm-up
- Early stopping based on validation MAE
- Model checkpointing for best weights

## C. CBAM Architecture Details

The Convolutional Block Attention Module (CBAM) enhances feature representations through sequential channel and spatial attention mechanisms:

TABLE I
CBAM COMPONENTS AND THEIR FUNCTIONS

| Component | Function |
|---|---|
| Channel Attention | <ul><li>Generates channel-wise attention weights</li><li>Uses both max-pooling and average-pooling operations</li><li>MLP with shared weights: FC-ReLU-FC</li><li>Output dimension: C×1×1</li></ul> |
| Spatial Attention | <ul><li>Creates spatial attention map</li><li>Combines max-pooled and avg-pooled features</li><li>7×7 convolution followed by sigmoid</li><li>Output dimension: 1×H×W</li></ul> |

The mathematical formulation of CBAM is as follows:

$$F' = M_c(F) \otimes F \tag{1}$$

$$F'' = M_s(F') \otimes F' \tag{2}$$

where $M_c$ and $M_s$ represent channel and spatial attention operations respectively, and $\otimes$ denotes element-wise multiplication.
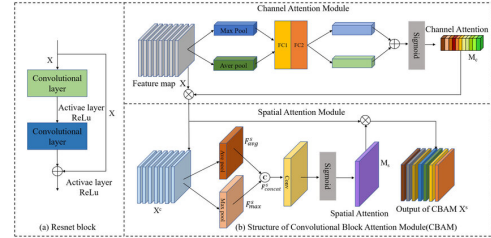


Fig. 6. Basic architecture of CBAM

## VI. RESULTS

This section summarizes the performance evaluation results for DenseNet121 and ResNet34 with Convolutional Block Attention Module (CBAM).

### A. DenseNet121

DenseNet121 demonstrated strong and stable performance in the bone age regression task, capitalizing on its densely connected architecture. Each layer in DenseNet121 receives input from all preceding layers, allowing efficient feature reuse and gradient propagation. This design encourages the network to learn compact and diversified features, especially important in medical images with high anatomical variability.

Our implementation followed the canonical DenseNet121 structure with four dense blocks—comprising 6, 12, 24, and 16 layers—interleaved with transition blocks that perform 1×1 convolutions and average pooling for dimensionality reduction. Each convolutional layer was followed by Batch Normalization and GELU activation, with L2 regularization

applied to mitigate overfitting. The architecture was adapted for grayscale hand X-rays with a single input channel.

Incorporating patient gender as an auxiliary input, we processed the binary gender value through a two-layer MLP (Dense(16)-Dense(8)), concatenated it with the global pooled image features, and passed the combined vector through a regularized regression head. The final prediction was made using a single linear output node.

DenseNet121 was trained over 60 epochs using the AdamW optimizer with decoupled weight decay ($\lambda = 0.01$), which effectively controls overfitting by penalizing large weights without interfering with adaptive moment estimation. A Huber loss function ($\delta = 15$) was used to balance sensitivity to outliers with robust error minimization, suitable for regression tasks involving clinical measurements like bone age.

A custom Cosine Warm-Up and Decay learning rate schedule was employed, starting from a base learning rate of $5 \times 10^{-3}$. The schedule began with a 20 percent warm-up phase (linearly increasing with cosine dynamics), followed by smooth cosine decay for the remaining training duration. This helped stabilize early training and improve convergence in later epochs.

All convolutional layers used the initialization and L2 regularization ($\lambda = 1 \times 10^{-4}$), promoting stable gradient flow and generalization. The activation function throughout the network was GELU, chosen for its smoother and more probabilistic activation profile compared to ReLU, improving performance on continuous-valued regression targets. Gender was included as an auxiliary input and processed via a two-layer MLP before concatenation with the image features in the regression head. Fully connected layers were regularized and used dropout ($p = 0.01$) to further mitigate overfitting.

Early stopping was triggered based on validation MAE, with a patience of 3 epochs and a stopping threshold of 8 months. Evaluation was based on MAE, RMSE, and $R^2$ metrics on both training and test sets.

DenseNet121 achieved the best performance among the compared models, with a test MAE of **8.18** months and an $R^2$ score of **0.93**. These results indicate that dense connectivity is highly effective for medical regression tasks where feature granularity is critical. A summary of the training and test metrics is shown in Table II.

TABLE II
DENSENET121 EVALUATION METRICS

| Dataset | Loss | MAE | RMSE | R² |
|---|---|---|---|---|
| Training | 60.78 | 8.50 | 11.18 | 0.92 |
| Test | 60.07 | 8.18 | 11.49 | 0.93 |

Training and validation curves in Figure 7 show smooth convergence and low variance across all metrics, supporting the model's reliability and generalization.
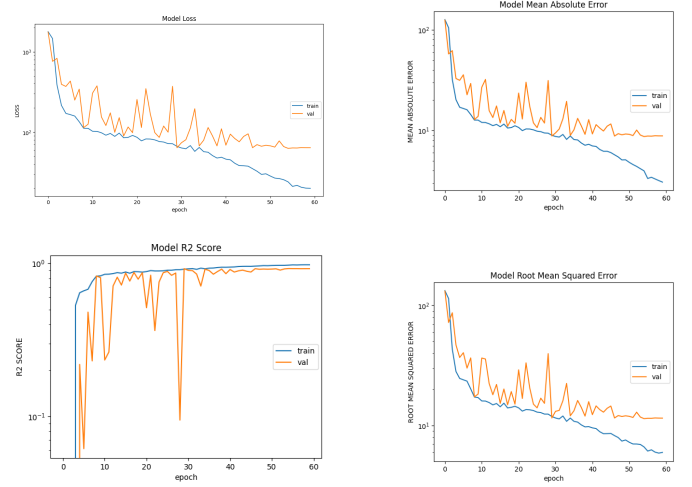


Fig. 7. Training and validation metrics (log scale) for the DenseNet121 model: (Top-left) Loss, (Top-right) Mean Absolute Error (MAE), (Bottom-left) $R^2$ Score, and (Bottom-right) Root Mean Squared Error (RMSE). These plots show model performance over 60 epochs.

### B. ResNet34+CBAM

The ResNet34+CBAM model was designed to enhance baseline residual learning by integrating attention mechanisms that selectively emphasize important features. The Convolutional Block Attention Module (CBAM) used in this architecture applies both channel and spatial attention in a sequential manner, aiming to improve the model's focus on anatomically relevant regions of the hand radiographs.

In this implementation, CBAM modules were inserted into each residual block of the ResNet34 backbone. The channel attention branch utilizes global average and max pooling followed by a shared multilayer perceptron (MLP) to recalibrate channel-wise responses. The spatial attention module combines pooled feature maps and applies a 7×7 convolution to highlight informative spatial locations. These refined features are then passed through the residual connection structure, allowing better localization of critical bone growth regions.

Additionally, the model incorporated patient gender as an auxiliary input. A dedicated two-layer MLP encodes the binary gender indicator, which is concatenated with the visual features at the regression head. This fusion introduces contextual information, enabling the model to differentiate growth trajectories across genders.

Despite these enhancements, the ResNet34+CBAM model showed more training and validation loss fluctuations compared to DenseNet121, as seen in Figure 8. The increased variance is likely due to the added model complexity and sensitivity of attention mechanisms to noise in medical images.

Nevertheless, the model achieved a test MAE of **9.77** months and an $R^2$ score of **0.90**, demonstrating that attention-based residual architectures are viable alternatives to densely connected networks in bone age estimation. Table III summarizes the complete evaluation metrics.

TABLE III
RESNET34+CBAM EVALUATION METRICS

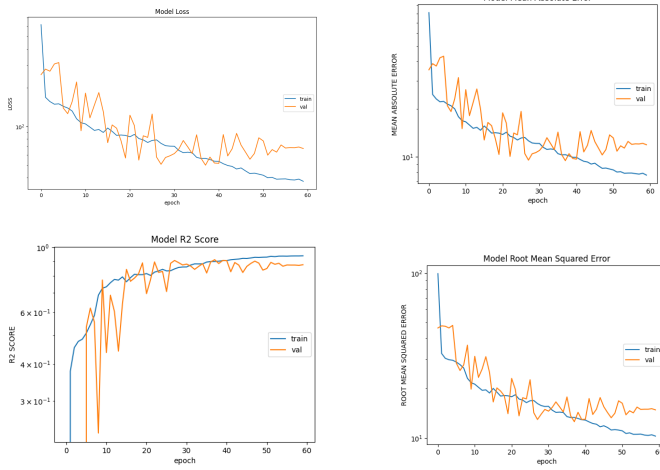| Dataset | Loss | MAE | RMSE | $R^2$ |
|---------|------|-----|------|-------|
| Training | 69.94 | 9.69 | 11.18 | 0.92 |
| Test | 70.61 | 9.77 | 13.01 | 0.90 |



Fig. 8. Training and validation metrics (log scale) for the ResNet34+CBAM model: (Top-left) Loss, (Top-right) Mean Absolute Error (MAE), (Bottom-left) $R^2$ Score, and (Bottom-right) Root Mean Squared Error (RMSE). These plots show model performance over 60 epochs.

## VII. CONCLUSION

This study presented a rigorous and technically optimized framework for automated bone age prediction using deep convolutional neural networks. We implemented and compared two advanced architectures—DenseNet121 and ResNet34 enhanced with a Convolutional Block Attention Module (CBAM)—under a unified training pipeline incorporating modern regularization and scheduling strategies.

A key contribution of our work lies in the structured training regimen, which includes:

- **Cosine annealing with warm-up**, ensuring stable initial convergence and avoiding sharp minima.
- **Huber loss**, balancing robustness to outliers and sensitivity to small prediction errors.
- **AdamW optimizer with decoupled weight decay**, promoting generalization through well-controlled parameter updates.
- **Auxiliary gender input**, which improves personalization in bone age predictions.
- **Early stopping on validation MAE** and **dropout**, jointly reducing overfitting risks.
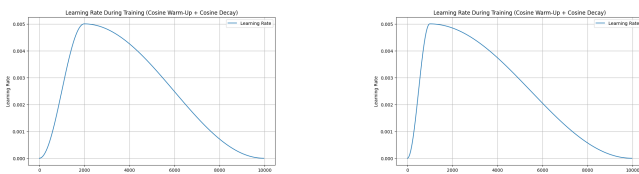


Fig. 9. Cosine learning rate schedule with warm-up and cosine decay. Left: DenseNet121. Right: ResNet34+CBAM.

Extensive experiments demonstrated that **DenseNet121 consistently outperforms ResNet34+CBAM** in both accuracy and model efficiency. On average, DenseNet121 achieved a lower mean absolute error (MAE) of **8.84 months**, compared to **10.16 months** for ResNet34+CBAM (Table IV). While isolated runs showed ResNet34+CBAM achieving a slightly better MAE (8.18 months vs. DenseNet's 9.7 months), these were exceptions and not representative of typical performance across trials.

TABLE IV
MODEL COMPARISON

| Model | Loss | MAE | RMSE | $R^2$ |
|-------|------|-----|------|-------|
| DenseNet121 | **60.07** | **8.18** | **11.49** | **0.93** |
| ResNet34 + CBAM | 70.61 | 9.77 | 13.01 | 0.90 |

Beyond predictive accuracy, DenseNet121 also proved more efficient:

- It required significantly fewer parameters—**7.1M vs. 22.6M** (Table V)—highlighting a lighter and more scalable architecture.

TABLE V
MODEL COMPLEXITY IN PARAMETERS

| Model | Total Params | Trainable | Non-trainable |
|-------|--------------|-----------|---------------|
| DenseNet121 | 7,075,241 | 6,991,593 | 83,648 |
| ResNet34 + CBAM | 22,590,634 | 22,573,610 | 17,024 |

- DenseNet121 had a longer training time (69.43 minutes vs. 40.42 minutes), it utilized more memory relative to its parameter count (Table VI).

TABLE VI
TRAINING TIME AND PEAK MEMORY USAGE FOR EACH MODEL

| Model | Training Time (min) | Memory Usage (MB) |
|-------|---------------------|-------------------|
| DenseNet121 | 69.43 | 6277.77 |
| ResNet34 + CBAM | 40.42 | 4140.07 |

Despite ResNet34+CBAM's enhancements through spatial and channel attention, its higher parameter complexity and sensitivity to hyperparameters led to less consistent results. DenseNet121, by contrast, offered a more stable and generalizable solution with fewer tuning effort.

**Future Work.** Building on these findings, future directions include:

- Integrating attention mechanisms directly into DenseNet backbones for the best of both architectures.
- Incorporating multi-modal clinical metadata (e.g., height, weight, ethnicity) to refine predictions.
- Leveraging uncertainty estimation methods for interpretability in clinical settings.
- Exploring transfer learning from related medical imaging tasks to improve sample efficiency.

In summary, our framework establishes a reliable foundation for automated bone age assessment. The results underscore that **DenseNet121 strikes a more favorable balance**

**between predictive performance, efficiency, and architectural robustness**, making it a strong candidate for real-world clinical deployment.

## REFERENCES

[1] Rajpurkar, P., Irvin, J., Zhu, K., Yang, B., Mehta, H., Duan, T. Automated bone age assessment using deep convolutional neural networks. In Proceedings of the Conference on Medical Image Computing and Computer-Assisted Intervention, San Francisco, CA, US, 2017.

[2] Loeff, F., Taebi, B., Moltz, J. H., & von Tscharner. Transfer learning for bone age assessment. Medical Image Analysis, 2019.

[3] Wang, Z., Tan, T., Wang, X., Zhang, Y., & Li, J. Multi-Modal Bone Age Assessment using Deep Convolutional Neural Networks. Medical Image Analysis, 2021.

[4] Liu, Y., Xiong, J., & Li, H. Attention-Based Deep Convolutional Neural Network for Bone Age Assessment. Medical Image Analysis, 2021.

[5] Pettersen et al. Hand radiographs for skeletal age assessment in pediatric populations. The Clinical Radiology and Radiotherapy, 2015.

[6] Lippe et al. Bone age in endocrine disorders. Pediatric Endocrinology and Metabolism, 2016.

[7] Bonse et al. Bone age in puberty. The Journal of Adolescent Health, 2013.

[8] Tassone et al. Bone age in genetic disorders. Medical Genetics Part C: Seminars in Medical Genetics, 2013.