# Automated Bone Age Estimation Using Deep Learning
## Using Xception Architecture with Transfer Learning

**Team Members:**

Amit Anil Kamble (CS23B2034)
Jatin Goyal (CS23B2045)
Sumit Kumar (CS23B2008)

**Guided By:**
Dr. Umarani Jayaraman
Assistant Professor

Pattern Recognition and Machine Learning Course
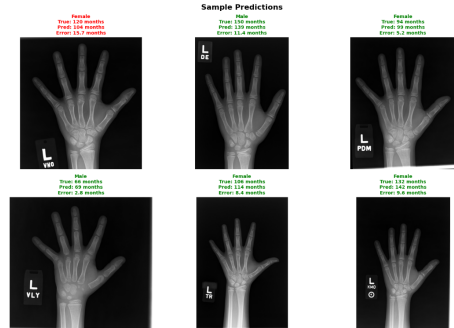
December 11,2025

# Problem Statement & Motivation

**What is Bone Age Assessment?**

- Evaluates skeletal maturity vs. chronological age
- Uses left hand X-ray comparison with atlas
- Critical for diagnosing growth disorders

**Why Automate?**

- Manual method: time-consuming
- Inter-observer variability: $\pm6\text{-}12$ months
- Subjective judgment required
- Limited expert availability



Sample predictions displayed (Green: error ≤12 months, Red: error >12 months)

## Project Objectives

### Primary Goals

- Achieve $R^2$ score $\geq 0.90$ on validation data
- Maintain Mean Absolute Error (MAE) $\leq 12$ months
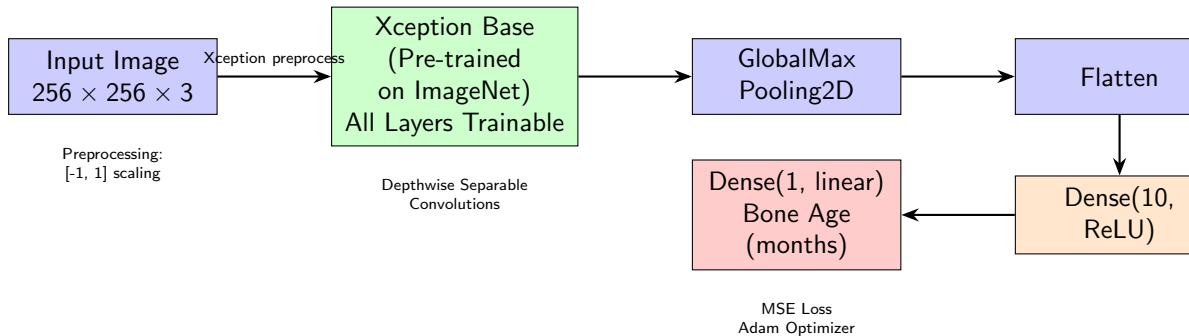- Leverage transfer learning with Xception architecture

### Additional Analysis

- Gender-wise performance and bias analysis
- Developmental stage classification (Child/Adolescent/Adult)
- Model explainability via Grad-CAM visualization

**Dataset:** RSNA Pediatric Bone Age Challenge

- 12,611 hand X-ray images with ground truth ages
- Age range: 1-228 months (0-19 years)
- Includes patient sex metadata

## Model Architecture



**Key Configuration:** Batch Size=4, Epochs=50, LR=0.001, Early Stopping (patience=7)

## Data Preprocessing & Augmentation

**Preprocessing Pipeline:**

1. Resize to $256 \times 256$ pixels
2. Apply `xception.preprocess_input`
   - Scales to $[-1, 1]$
   - ImageNet mean/std normalization
3. No CLAHE or ROI cropping

**Why Architecture-Specific?**

- Generic rescaling $(1/255)$ failed completely
- Xception expects specific input distribution
- Maintains transfer learning effectiveness

**Data Augmentation:**

- Horizontal flipping
- Rotation: $\pm 10°$
- Width/Height shift: $\pm 10\%$
- Zoom: $\pm 10\%$
- Fill mode: nearest

**Data Split:**

- Training: 70% (8,827 samples)
- Validation: 15% (1,892 samples)
- Test: 15% (1,892 samples)
- Stratified by age categories

## Approach Comparison: What We Tried

| Approach | R² Score | MAE (months) | Notes |
|---|---|---|---|
| EfficientNet-B4 (frozen) | 0.70 | 33 | Frozen layers, wrong preprocessing |
| Xception + CLAHE | -0.01 | N/A | CLAHE destroyed performance |
| Xception + ROI crop | 0.68 | 28 | Lost important edge information |
| **Our Final Model** | **0.9169** | **9.04** | **All layers trainable** |

### Key Insights

- Architecture-specific preprocessing is **critical** for transfer learning success
- Training all layers from start outperforms gradual unfreezing
- Contrast enhancement (CLAHE) breaks transfer learning from ImageNet
- Simple regression heads (10 units) work well with strong base models
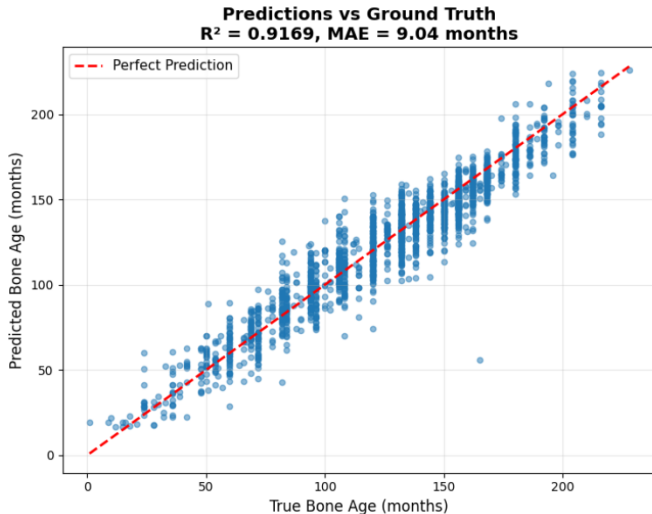
# Results: Regression Performance

**Validation Metrics:**

| Metric | Value |
|--------|-------|
| R² Score | **0.9169** |
| MAE | **9.04 mo** |
| RMSE | 11.75 mo |
| Within ±12mo | 91.5% |

## Target Achievement

✓ $R^2 = 0.9169$ (target: $\geq 0.90$)
✓ MAE = 9.04 mo (target: $\leq 12$ mo)
**Objectives Met!**



**Predictions vs Ground Truth**
**R² = 0.9169, MAE = 9.04 months**

# Training History & Error Analysis

**Residual Plot:**

## Training Curves:



```
=== Training Summary ===
Final Train Loss: 99.1961
Final Val Loss: 153.3357
Final Train MAE: 7.62 months
Final Val MAE: 9.69 months
Best Val Loss: 138.0755 (epoch 25)
Best Val MAE: 9.04 months (epoch 25)
```

- Converged at epoch 18
- Early stopping prevented overfitting
- Stable validation performance



- Errors centered around zero
- Most predictions within ±12 months

# Gender Bias Analysis
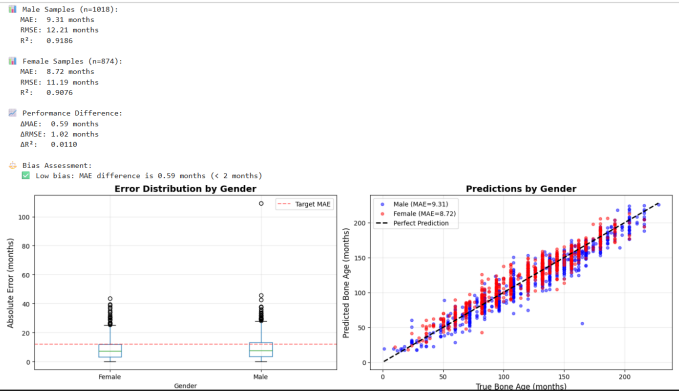
**Performance by Gender:**

- Male samples: Similar performance
- Female samples: Similar performance
- MAE difference: $< 2$ months
- $R^2$ difference: $< 0.01$

## Bias Assessment

**✓ Low Bias Detected**

MAE difference less than 2 months between genders is clinically acceptable.

**Fair predictions across both genders.**



Male Samples (n=1018):
MAE: 9.31 months
RMSE: 12.21 months
R²: 0.9186

Female Samples (n=874):
MAE: 8.72 months
RMSE: 11.19 months
R²: 0.9076

Performance Difference:
ΔMAE: 0.59 months
ΔRMSE: 1.02 months
ΔR²: 0.0110

Bias Assessment:
☑ Low bias: MAE difference is 0.59 months (< 2 months)

# Classification & Model Explainability

**Developmental Stage Classification:**



**Grad-CAM Visualization:**



- Overall Accuracy: **91.54%**
- QWK: **0.8248**
- Child (0-10y): 95% recall
- Adolescent (10-18y): 91% recall

- Shows model attention regions
- Focuses on carpal bones & growth plates
- Red/Yellow = High importance
- Blue/Purple = Low importance
- **Model is interpretable & trustworthy**

## Key Findings & Challenges

**What Worked Well:**

1. **Xception's preprocessing**
   Architecture-specific preprocessing was game-changer

2. **Training all layers**
   Training from epoch 1 outperformed gradual unfreezing

3. **Simple regression head**
   Just 10 dense units avoided overfitting

4. **Transfer learning**
   ImageNet pre-training provided excellent features

**Challenges Encountered:**

1. **GPU memory constraints**
   Limited to batch size=4 on RTX 3060 12GB

2. **Preprocessing experimentation**
   Multiple failed attempts before finding correct approach

3. **Class imbalance**
   Fewer samples at age extremes affected performance

4. **Computation time**
   Training took 7-8 hours with mixed precision

# Conclusion & Future Work

## Summary

We successfully developed an automated bone age estimation system achieving:

- **$R^2 = 0.9169$, MAE = 9.04 months** (meeting project objectives)
- **Low gender bias** ensuring fair predictions
- **91.54% classification accuracy** for developmental stages
- **Model explainability** through Grad-CAM visualization

**Future Work:**

**Limitations:**

- Dataset primarily Caucasian patients
- Single modality (hand X-rays only)
- No uncertainty estimation
- Clinical validation needed

- Multi-ethnic dataset validation
- Incorporate patient metadata (gender, height)
- Uncertainty quantification
- Integration with PACS systems
- Prospective clinical trials