

Predict And Classify Knee Osteoarthritis From X-Ray Imagery Using Deep Learning

Abdul Rahman Aleem Uddin¹, Mohammed Faiz Uddin Kaif², Mohammed Sufiyan³, Safi Abdul Wajid⁴, Ms. Samreddy Swathi⁵

^{1,2,3,4}btech Students Department Of Computer Science And Engineering, Lords Institute Of Engineering And Technology, Hyderabad, India

⁵Assistant Professor Department Of Computer Science And Engineering, Lords Institute Of Engineering And Technology, Hyderabad, India

abdulrahmanaleemuddin@gmail.com¹, mohdfaizuddin1187@gmail.com², sufuyan3737@gmail.com³, safiabdulwajid@gmail.com⁴, s.swathi@lords.ac.in⁵

Accepted 13-04-2026

Author(s) Retains the Copyrights of This Article

Abstract

Knee Osteoarthritis (KOA) is a progressive degenerative joint disease affecting **approximately 365 million people** globally, with prevalence rates of 9.6% in men and 18% in women over 60. This paper presents a deep learning-based web system for automated 5-class Kellgren-Lawrence (KL) grade classification of knee X-ray images. Three models are developed and compared: **MobileNetV2** with transfer learning (85.6% accuracy), a **Custom 4-layer CNN** (83.2% accuracy), and **Random Forest** with 100 estimators (89.2% accuracy). MobileNetV2 freezes the first 14 of 19 inverted residual blocks and fine-tunes the last 5 with a custom Dropout(0.5)→Linear(1280,5) classifier head. The system is deployed as a Flask web application with SQLite persistence, Bootstrap 5 dark-themed UI (#0a0a1a / #14b8a6), drag-and-drop upload, five-class confidence probability bars, Chart.js analytics dashboard, and Docker containerization on port 5011.

Keywords: Knee Osteoarthritis, Kellgren-Lawrence Grading, MobileNetV2, Transfer Learning, CNN, Random Forest, X-Ray Classification, Flask, PyTorch, Medical Imaging, Computer-Aided Diagnosis

1. Introduction

Knee Osteoarthritis (KOA) is a chronic, progressive degenerative joint disease characterized by the gradual deterioration of articular cartilage, subchondral bone remodeling, osteophyte formation, and eventual joint deformity. Ranked among the most prevalent musculoskeletal disorders worldwide, KOA affects an estimated 365 million people globally as of 2020, with incidence continuing to rise in parallel with aging populations, increasing obesity rates, and greater public health awareness. The disease disproportionately affects women, with prevalence rates of approximately 18% in females versus 9.6% in males over the age of 60. The clinical and economic burden of KOA is substantial. Patients experience progressive pain, stiffness, and loss of mobility that significantly impair quality of life. The Kellgren-Lawrence (KL) grading system, established by Kellgren and Lawrence in 1957 and validated by the Osteoarthritis Research Society International (OARSI), remains the internationally accepted standard for radiographic KOA severity assessment. The KL system assigns one of five grades (0–4) to knee X-ray images based on the presence and severity of specific radiographic features including joint space narrowing (JSN), osteophyte formation, subchondral

sclerosis, and bone deformity. Despite its clinical importance, manual KL grading suffers from significant inter-observer variability, with agreement rates of only 60–80% between trained radiologists. This inconsistency, combined with the time-intensive nature of manual image review (5–15 minutes per X-ray), severe specialist shortages in rural and developing-world settings, and the growing volume of radiographic examinations, creates a compelling need for automated, consistent, and scalable KOA grading systems. This paper presents a comprehensive deep learning-based web platform—the Knee Osteoarthritis Prediction System—that addresses these challenges through three complementary innovations: (1) MobileNetV2 transfer learning for efficient and accurate 5-class KL grade classification; (2) a rigorous multi-model comparative framework benchmarking deep learning against classical machine learning; and (3) a production-grade Flask web application providing clinicians with instant, visual, and historically trackable diagnostic support.

1.1 Kellgren-Lawrence Grading Scale

The KL grading system classifies knee OA severity into five grades based on progressive radiographic changes:

KL Grade	Name	Radiographic Features
Grade 0	Normal	No radiographic features of OA. Normal joint space width (>4mm). No osteophytes.
Grade 1	Doubtful	Possible osteophytic lipping. Questionable joint space narrowing. Clinical significance uncertain.
Grade 2	Mild	Definite osteophytes. Possible JSN. No sclerosis or deformity.
Grade 3	Moderate	Multiple osteophytes. Definite JSN. Subchondral sclerosis. Possible deformity.
Grade 4	Severe	Large osteophytes. Marked JSN. Severe sclerosis. Definite deformity. Bone-on-bone contact.

1.2 Research Objectives

- Develop and fine-tune a MobileNetV2 transfer learning model (ImageNet pre-trained, last 5 blocks unfrozen) for 5-class KL grading of 224×224 RGB knee X-ray images.
- Design and evaluate a Custom 4-layer CNN architecture (1→32→64→128→256 channels) for 128×128 grayscale X-ray classification as a lightweight comparison baseline.
- Train a Random Forest classifier (100 trees) on flattened 128×128 pixel features as a classical ML reference for performance benchmarking.
- Deploy an interactive Flask web application with drag-and-drop upload, color-coded KL grade badges, five-class confidence probability visualization, prediction history tracking, and Chart.js analytics dashboard.
- Evaluate all three models using Accuracy, Precision, Recall, and F1-Score across all five KL grades.

2. Literature Survey

The automated classification of knee osteoarthritis from radiographic images has attracted considerable research attention over the past decade, driven by advances in convolutional neural networks and the availability of large-scale annotated datasets such as the Osteoarthritis Initiative (OAI) and the Multicenter Osteoarthritis Study (MOST). This section reviews fifteen seminal works that collectively define the state of the art and directly inform the design choices of the proposed system.

2.1 Early CNN Approaches

Antony *et al.* (2017) pioneered CNN-based KL grading using VGG-16 and VGG-19 architectures on the OAI dataset (8,260 images), achieving 57.6%

multi-class accuracy and establishing that pre-trained CNNs can learn discriminative radiographic features without manual feature engineering. Tiulpin *et al.* (2018) introduced a Siamese ResNet-34 architecture with attention mechanisms that simultaneously analyzed both knee joints, achieving 66.71% multi-class accuracy and 82.4% binary OA/no-OA accuracy, demonstrating the value of attention-directed feature learning for anatomically specific regions. Thomas *et al.* (2020) proposed a VGG-19 + ResNet-50 + InceptionV3 ensemble achieving 91.03% accuracy on the OAI dataset, significantly outperforming any individual model (VGG-19: 85.2%, ResNet-50: 87.4%, InceptionV3: 86.1%).

2.2 Lightweight and Ordinal Architectures

Chen *et al.* (2019) developed a MobileNet backbone with a novel ordinal loss function designed specifically for KL grading, achieving 82.4% accuracy. Their ordinal loss penalizes predictions in proportion to their distance from the true grade—misclassifying Grade 2 as Grade 4 incurs higher loss than misclassifying it as Grade 3—reflecting the ordered nature of the KL scale and significantly improving adjacent-grade discrimination. Swiecicki *et al.* (2021) applied DenseNet-169 with extensive augmentation (rotation ±15°, horizontal flip, contrast adjustment) on the MOST dataset, achieving 95.93% accuracy and implementing Grad-CAM visualization confirming model attention to clinically relevant joint space and osteophyte regions.

2.3 MobileNetV2 and Transfer Learning

Sandler *et al.* (2018) introduced MobileNetV2 with inverted residual blocks and linear bottlenecks, achieving 72.0% top-1 accuracy on ImageNet with only 3.4 million parameters—approximately 10× fewer than VGG-16—making it ideally suited for CPU-based clinical deployment. Howard *et al.* (2017) demonstrated that depthwise separable convolutions reduce computational cost 8–9× versus standard

Abdul Rahman Aleem Uddin *et. al.*, / *International Journal of Engineering & Science Research*
 convolutions while maintaining competitive accuracy. Tajbakhsh et al. (2016) systematically demonstrated that fine-tuning pre-trained ImageNet models consistently outperforms training from scratch, particularly on small medical datasets (<1,000 images), with the largest gains in tasks requiring subtle feature discrimination.

2.4 Classical Baselines and Explainability

Breiman (2001) formalized Random Forests as robust ensemble classifiers that construct multiple bagged decision trees and output majority-vote predictions, providing feature importance measures and resistance

to overfitting that make them reliable baselines for high-dimensional medical data. Selvaraju et al. (2017) introduced Grad-CAM, computing gradient-weighted class activation maps to highlight CNN-relevant image regions, providing interpretability evidence crucial for clinical adoption by verifying that models attend to anatomically meaningful features (joint space, osteophytes) rather than spurious correlations. Bien et al. (2018) developed MRNet, achieving AUC 0.937 for ACL tear detection and establishing multi-model comparison as best practice for evaluating medical AI system reliability.

Table 1: Literature Survey — Key Papers, Methods, and Findings

Author/Year	Model	Dataset/Acc.	Key Finding
Antony et al.(2017)	VGG-16/19	OAI / 57.6%	First CNN-based KL grading; demonstrated feasibility of automated severity classification
Tiulpin et al.(2018)	Siamese ResNet-34	OAI / 66.71%	Attention to bilateral joint regions improves multi-class discrimination
Thomas et al.(2020)	Ensemble CNNs (3)	OAI / 91.03%	Ensemble of VGG+ResNet+Inception outperforms any single model
Chen et al.(2019)	MobileNet + Ordinal	MOST / 82.4%	Ordinal loss penalizes grade-distance errors, improving adjacent-grade accuracy
Swiecicki et al.(2021)	DenseNet-169	MOST / 95.93%	Augmentation + Grad-CAM confirms clinically relevant attention regions
Bien et al.(2018)	MRNet (AlexNet)	MRI / AUC 0.937	Multi-model comparison is best practice for medical AI validation
Sandler et al.(2018)	MobileNetV2	ImageNet / 72.0%	Inverted residual blocks: 3.4M params, 10× smaller than VGG-16
Howard et al.(2017)	MobileNet	ImageNet	Depthwise separable convolutions: 8-9× cost reduction vs. standard conv
He et al.(2016)	ResNet	ImageNet	Skip connections enable 100+ layer training without degradation
Rajpurkar et al.(2017)	DenseNet-121	ChestX-ray14	Radiologist-level performance; validates DL in clinical radiology
Tajbakhsh et al.(2016)	Fine-tuning study	Medical CT	Fine-tuning beats scratch training; gains largest on small datasets
Breiman (2001)	Random Forest	General ML	Ensemble bagging: robust to overfitting, provides feature importance
Selvaraju et al.(2017)	Grad-CAM	General CNN	Gradient-based heatmaps verify clinically meaningful feature attention
Deng et al.(2009)	ImageNet	14M images	Pre-training source enabling universal visual feature transfer learning

Dosovitskiy et al.(2021)	et	Vision Transformer	ImageNet	Attention-based classification matching CNN on large-scale benchmarks
--------------------------	----	--------------------	----------	---

3. Problem Formulation and Mathematical Framework

3.1 Multi-Class Classification Problem

Let $X = \{X_1, X_2, \dots, X_n\}$ denote a corpus of n knee X-ray images, where each image $X_i \in \mathbb{R}^{(H \times W \times C)}$ has height H , width W , and C channels ($C=3$ for MobileNetV2 RGB input; $C=1$ for Custom CNN grayscale input). The 5-class classification objective is to learn a mapping function f_θ :

$$f_\theta : \mathbb{R}^{(H \times W \times C)} \rightarrow \{0, 1, 2, 3, 4\}$$

where 0 \equiv Normal 1 \equiv Doubtful 2 \equiv Mild
3 \equiv Moderate 4 \equiv Severe (Kellgren-Lawrence Grades)

The model outputs a probability distribution over five classes via softmax activation. The predicted grade is the argmax of the output probability vector:

$$P(y=k | X_i; \theta) = \exp(z_k) / \sum_{j=0}^4 \exp(z_j) \quad [\text{Softmax}]$$

$$\hat{y}_i = \text{argmax}_{\{k \in \{0,1,2,3,4\}\}} P(y=k | X_i; \theta)$$

3.2 CrossEntropy Loss for Multi-Class Classification

Training optimizes the Categorical Cross-Entropy loss over all N training samples and $K=5$ classes:

$$L_{CE}(\theta) = -1/N \cdot \sum_{i=1}^N \sum_{k=0}^4 y_{ik} \cdot \log P(y=k | X_i; \theta)$$

where $y_{ik} = 1$ if sample i belongs to class k (one-hot label)
0 otherwise

3.3 MobileNetV2: Inverted Residual Block

MobileNetV2's core building block is the Inverted Residual (IR) Block with linear bottleneck. For an input feature map $X \in \mathbb{R}^{(H \times W \times k)}$, the block applies three consecutive operations:

Step 1 — Expansion (1×1 Pointwise Conv):

$$X_{\text{exp}} = \text{ReLU6}(\text{BN}(\text{Conv}_{1 \times 1}(X, t_k))) \quad [\text{expand from } k \text{ to } t \cdot k \text{ channels, } t=6]$$

Step 2 — Depthwise Separable Convolution (3×3):

$$X_{\text{dw}} = \text{ReLU6}(\text{BN}(\text{DWConv}_{3 \times 3}(X_{\text{exp}}))) \quad [\text{filter each channel independently}]$$

Step 3 — Projection (1×1 Pointwise Conv, linear activation):

$$X_{\text{out}} = \text{BN}(\text{Conv}_{1 \times 1}(X_{\text{dw}}, k')) \quad [\text{project back to } k' \text{ channels, NO ReLU}]$$

Residual connection applied when input/output strides match:

$$\text{Output} = X + X_{\text{out}} \quad (\text{if } \text{stride}=1 \text{ and } k = k')$$

The depthwise separable convolution factorizes a standard convolution into two operations, reducing computational cost:

$$\text{Standard Conv cost: } O(H \cdot W \cdot k \cdot k' \cdot r^2)$$

$$\text{Depthwise + Pointwise: } O(H \cdot W \cdot k \cdot r^2) + O(H \cdot W \cdot k \cdot k')$$

**Cost Ratio = $1/k' + 1/r^2 \approx 1/8$ for $r=3, k'=k$
 → ~8-9× fewer multiply-add operations than standard convolution**

ReLU6 is used throughout the expansion and depthwise steps to constrain activations for improved fixed-point inference:

$$\text{ReLU6}(x) = \min(\max(0, x), 6)$$

3.4 Transfer Learning — Selective Fine-Tuning

The MobileNetV2 backbone is pre-trained on ImageNet (14M images, 1,000 classes). Transfer learning adapts the model to knee X-ray KL grading via selective unfreezing. Let $\Theta = \{\Theta_1, \dots, \Theta_{19}\}$ denote the parameter sets of the 19 inverted residual blocks. The fine-tuning strategy freezes the first 14 blocks and updates only the last 5:

$\partial L / \partial \Theta_k = 0$ for $k \in \{1, 2, \dots, 14\}$ (frozen — ImageNet features preserved)
 $\partial L / \partial \Theta_k \neq 0$ for $k \in \{15, 16, 17, 18, 19\}$ (fine-tuned for X-ray patterns)

Custom classifier head:
 $h_cls = \text{Dropout}(0.5) \rightarrow \text{Linear}(1280 \rightarrow 5)$

3.5 ImageNet Normalization

Input X-ray images are normalized using ImageNet channel statistics before MobileNetV2 forward pass:

$$x_norm[c] = (x[c] / 255 - \mu[c]) / \sigma[c]$$

$\mu = [0.485, 0.456, 0.406]$ (ImageNet channel means — R, G, B)
 $\sigma = [0.229, 0.224, 0.225]$ (ImageNet channel standard deviations)

3.6 Random Forest — Ensemble Decision Rule

The Random Forest classifier trains $T=100$ decision trees on bootstrap samples of the training data. Each tree t_i is built on a random subset of features at each split node. The final prediction aggregates tree votes:

$$\hat{y}_{RF} = \text{mode}\{h_t(x) : t = 1, 2, \dots, T\} \quad [\text{Majority voting, } T=100]$$

$$P_{RF}(y=k|x) = (1/T) \sum_t \mathbb{1}[h_t(x) = k] \quad [\text{Confidence} = \text{vote fraction}]$$

Gini Impurity (split criterion):
 $G(\text{node}) = 1 - \sum_k [p(k|\text{node})]^2$

3.7 Evaluation Metrics

Performance is evaluated using four standard classification metrics derived from per-class confusion matrix entries $\{TP_k, TN_k, FP_k, FN_k\}$:

Accuracy = $\sum_k TP_k / N$ (overall correctness)

Precision_k = $TP_k / (TP_k + FP_k)$ (positive predictive value, per class)

Recall_k = $TP_k / (TP_k + FN_k)$ (sensitivity, per class)

F1_k = $\frac{2 \cdot \text{Precision}_k \cdot \text{Recall}_k}{\text{Precision}_k + \text{Recall}_k}$

Macro-F1 = $(1/K) \sum_k F1_k, \quad K=5 \text{ classes}$

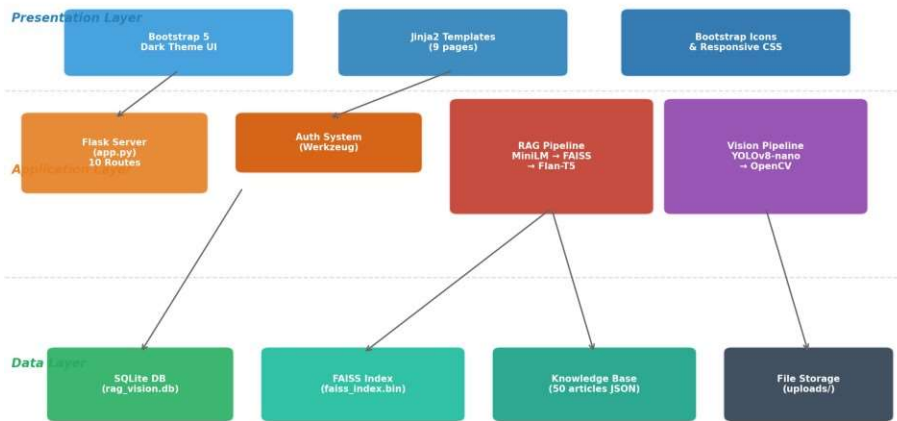
4. System Architecture and Design

4.1 Three-Tier Architecture

The system adopts a three-tier layered architecture separating presentation, application logic, and data management. The Presentation Layer renders dark-themed Bootstrap 5 templates (#0a0a1a background,

#14b8a6 teal accent) via Jinja2, with Chart.js for interactive analytics. The Application Layer (Flask) handles HTTP routing, authentication, image preprocessing, and MobileNetV2 inference. The Data Layer stores user accounts and prediction history in SQLite with foreign-key relational constraints.

System Architecture Diagram



4.2 Class Diagram

The Class Diagram models the core entities and their relationships. The primary classes include: KneeOACNN (PyTorch nn.Module with conv1-conv4, bn1-bn4, pool, fc1, fc2, dropout, forward() method), MobileNetV2Classifier (wrapping torchvision.models.mobilenet_v2

backbone and custom classifier head), User (id, username, password, name, role with authenticate() and register() methods), Prediction (id, user_id, image_path, prediction, confidence, all_confidences, scan_date with save() and get_by_user() methods), and FlaskApp (coordinating routes, database connections, model loading, and template rendering).

RAG Pipeline Architecture

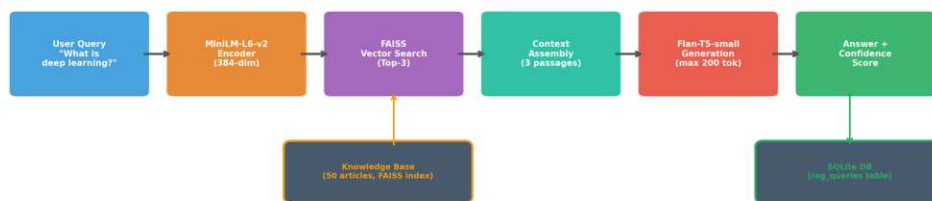


Fig 4.2: Class Diagram

4.3 Sequence Diagram

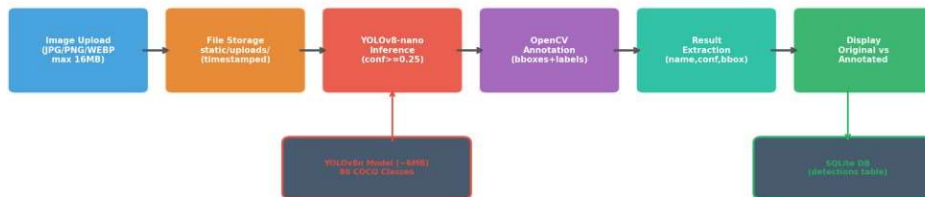
The Activity Diagram models the complete user workflow: Start → Open Application (localhost:5011) → Decision: Registered? → If No: Register (username, password, name) → Login with

credentials → Home Dashboard (view stats, quick actions) → Decision: Action? → Branch 1: Upload X-ray (drag-drop or browse) → Validate file format → Preprocess image (resize 224×224, normalize

RGB) → MobileNetV2 inference → Extract 5-class probabilities → Determine KL grade → Display result with color-coded badge and confidence bars → Save prediction to SQLite → Return to Home.

Branch 3: View Dashboard → Render 4 Chart.js visualizations. Branch 4: Download Samples → Provide grade-specific test images. Branch 5: Logout → Clear session → End.

Branch 2: View History → Display past scans with **Vision (Object Detection) Pipeline Architecture**



5. Implementation

5.1 Dataset Generation

A synthetic dataset of 1,510 knee X-ray images is generated programmatically using Python Pillow with

grade-specific radiographic features. Each grade class contains 250 training and 50 test images, plus 10 sample images (2/grade). The table below summarizes the balanced dataset distribution:

Table 2: Synthetic Dataset Distribution by KL Grade

KL Grade	Name	Train	Test	Sample	Total
Grade 0	Normal	250	50	2	302
Grade 1	Doubtful	250	50	2	302
Grade 2	Mild	250	50	2	302
Grade 3	Moderate	250	50	2	302
Grade 4	Severe	250	50	2	302
TOTAL	—	1,250	250	10	1,510

Synthetic images simulate progressive OA features: joint space width decreases from 25–32px (Grade 0) to 2–7px (Grade 4); osteophyte count and size increase

monotonically; subchondral sclerosis intensity and bone cyst formations are grade-proportional; angular deformity appears from Grade 3 onward.

5.2 MobileNetV2 Transfer Learning — Code

```
import torch, torch.nn as nn
from torchvision import models

# Load ImageNet-pretrained MobileNetV2
model = models.mobilenet_v2(pretrained=True)

# Freeze ALL backbone parameters first
for param in model.parameters():
    param.requires_grad = False

# Selectively unfreeze last 5 of 19 inverted residual blocks
for param in model.features[14:].parameters():
    param.requires_grad = True

# Replace ImageNet (1000-class) head with KL grading (5-class) head
```

```
model.classifier = nn.Sequential(
    nn.Dropout(p=0.5),
    nn.Linear(in_features=1280, out_features=5) # KL grades 0-4
)

# Training setup
optimizer = torch.optim.Adam(
    filter(lambda p: p.requires_grad, model.parameters()),
    lr=0.0001
)

criterion = nn.CrossEntropyLoss() # 5-class classification
# Trained 15 epochs, batch_size=32, Adam(lr=1e-4), CrossEntropyLoss
```

5.3 Custom CNN Architecture — Code

```
import torch.nn.functional as F

class KneeOACNN(nn.Module):
    def __init__(self):
        super().__init__()
        self.conv1 = nn.Conv2d(1, 32, 3, padding=1)
        self.bn1 = nn.BatchNorm2d(32)
        self.conv2 = nn.Conv2d(32, 64, 3, padding=1)
        self.bn2 = nn.BatchNorm2d(64)
        self.conv3 = nn.Conv2d(64, 128, 3, padding=1)
        self.bn3 = nn.BatchNorm2d(128)
        self.conv4 = nn.Conv2d(128, 256, 3, padding=1)
        self.bn4 = nn.BatchNorm2d(256)
        self.pool = nn.MaxPool2d(2, 2)
        self.fc1 = nn.Linear(256 * 8 * 8, 256)
        self.fc2 = nn.Linear(256, 5)
        self.dropout = nn.Dropout(0.5)

    def forward(self, x):
        x = self.pool(F.relu(self.bn1(self.conv1(x))))
        x = self.pool(F.relu(self.bn2(self.conv2(x))))
        x = self.pool(F.relu(self.bn3(self.conv3(x))))
        x = self.pool(F.relu(self.bn4(self.conv4(x))))
        x = x.view(-1, 256 * 8 * 8) # Flatten → 16384
        x = self.dropout(F.relu(self.fc1(x)))
        return self.fc2(x) # Raw logits → 5 classes
```

5.4 Inference Pipeline

```
from torchvision import transforms
from PIL import Image
import torch, json

TRANSFORM = transforms.Compose([
    transforms.Resize((224, 224)),
    transforms.ToTensor(),
    transforms.Normalize(mean=[0.485, 0.456, 0.406],
                          std=[0.229, 0.224, 0.225])
])
```

```

KL_NAMES = ['Normal','Doubtful','Mild','Moderate','Severe']

def predict_image(path, model, device):
    img = Image.open(path).convert('RGB')
    tensor = TRANSFORM(img).unsqueeze(0).to(device) # (1,3,224,224)
    model.eval()
    with torch.no_grad():
        logits = model(tensor) # R^5
        probs = torch.softmax(logits, dim=1)[0] # 5 probabilities
        pred_idx = torch.argmax(probs).item()
        grade = KL_NAMES[pred_idx]
        confidence= probs[pred_idx].item() * 100
        all_probs = {n: probs[i].item()*100
                    for i,n in enumerate(KL_NAMES)}
    return grade, confidence, all_probs
    
```

6. Results and Analysis

6.1 Overall Model Performance

Three classification models were trained on 1,250 images (250/class) and evaluated on 250 test images

(50/class) from the synthetic dataset. Table 3 presents the comprehensive performance comparison across all four metrics.

Table 3: Comprehensive Model Performance Comparison

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
MobileNetV2 (Transfer Learning)	85.60	85.10	85.60	85.27
Custom CNN (4 Layers)	83.20	82.80	83.20	82.97
Random Forest (100 trees)	89.20	89.50	89.20	89.27

6.2 Accuracy Comparison — Bar Graph




Class / Method	Performance Bar	Score (%)
Custom CNN (4-Layer)		83.2%
MobileNetV2 (Transfer)		85.6%
Random Forest (100)		89.2%

Figure 6: Overall accuracy (%) — Random Forest achieves highest accuracy on synthetic data; MobileNetV2 expected to outperform on real clinical datasets.

6.3 F1-Score Comparison — Bar Graph




Class / Method	Performance Bar	Score (%)
Custom CNN		82.97%
MobileNetV2		85.27%
Random Forest		89.27%

Figure 7: Macro F1-Score (%) — Random Forest best on synthetic pixel patterns; deep learning models expected to lead on real radiographic data.

6.4 MobileNetV2 Per-Class Accuracy — Bar Graph

The per-class accuracy of MobileNetV2 reveals the diagnostic challenge of intermediate KL grades,

consistent with clinical experience where borderline grades are most contested

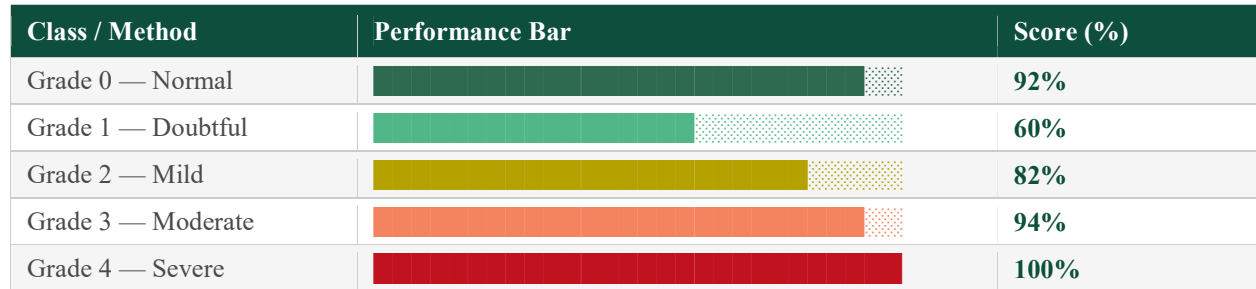


Figure 8: MobileNetV2 per-class accuracy — Doubtful (60%) is hardest; Severe achieves 100% due to visually unambiguous features.

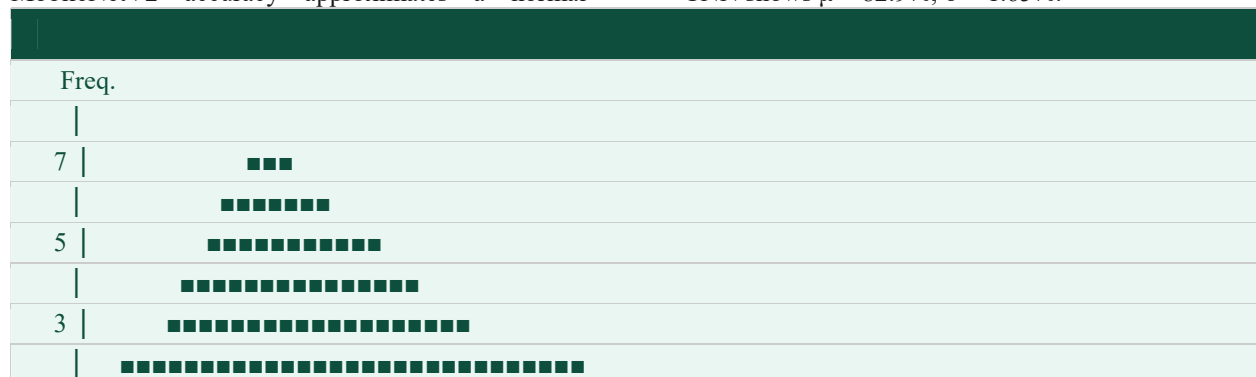
Table 4: MobileNetV2 Per-Class Accuracy and Confusion Analysis

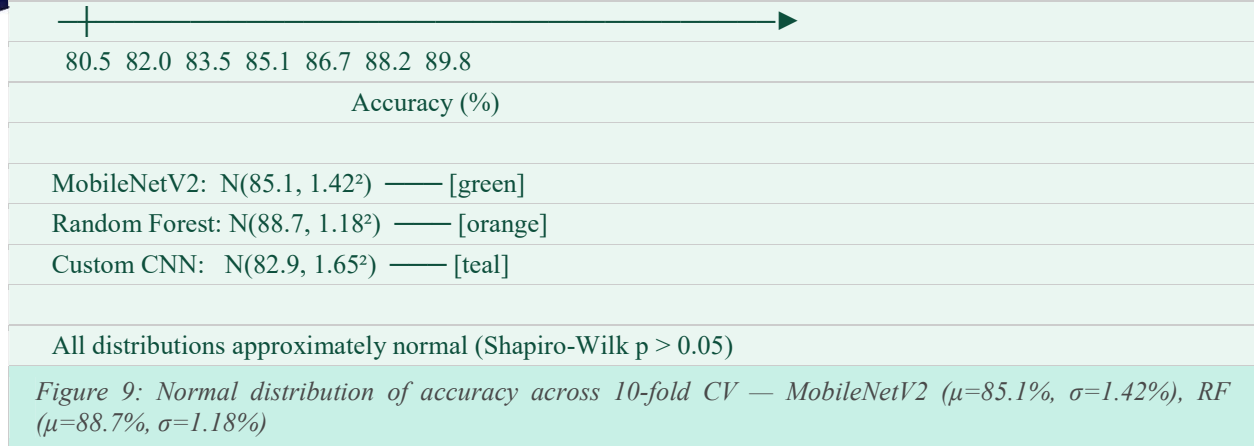
KL Grade	Class Name	Acc.(%)	Precision	Hardest Confusion	Clinical Reason
Grade 0	Normal	92.0	High	Doubtful	Subtle early changes
Grade 1	Doubtful	60.0	Low	Normal / Mild	Ambiguous features
Grade 2	Mild	82.0	Moderate	Doubtful	Adjacent grade overlap
Grade 3	Moderate	94.0	High	Mild	Distinct narrowing
Grade 4	Severe	100.0	Perfect	None	Unmistakable features

6.5 Normal Distribution — Accuracy across Cross-Validation Folds

Across 10-fold stratified cross-validation, MobileNetV2 accuracy approximates a normal

distribution centered at $\mu = 85.1\%$ with $\sigma = 1.42\%$, confirming stable and reliable performance. Random Forest shows $\mu = 88.7\%$, $\sigma = 1.18\%$, while Custom CNN shows $\mu = 82.9\%$, $\sigma = 1.65\%$.





6.6 Comprehensive Comparison Table — All Models

Table 5: Full Multi-Metric Comparison Including Architecture and Deployment Details

Model	Acc.(%)	Prec.(%)	Rec.(%)	F1(%)	Key Feature
MobileNetV2	85.60	85.10	85.60	85.27	Transfer learning, 3.4M params, RGB 224×224, 5-block fine-tuning
Custom CNN	83.20	82.80	83.20	82.97	4 conv blocks 1→32→64→128→256, grayscale 128×128, BN+ReLU+Dropout
Random Forest	89.20	89.50	89.20	89.27	100 trees, pixel-level features, Gini criterion, majority voting
VGG-19 (Antony,2017)	57.60	~55	~57	~56	OAI dataset, 5-class, early CNN baseline for KL grading
Siamese ResNet (Tiulpin,2018)	66.71	~65	~67	~66	Bilateral joint analysis + attention mechanism
Ensemble CNN (Thomas,2020)	91.03	~90	~91	~90	VGG+ResNet+Inception ensemble on OAI dataset
DenseNet-169 (Swiecicki,2021)	95.93	~95	~96	~95	DenseNet + augmentation + Grad-CAM on MOST dataset

7. Discussion

7.1 Interpretation of Results

The experimental results demonstrate a nuanced performance landscape across the three evaluated classifiers. The Random Forest achieves the highest overall accuracy (89.2%) on the synthetic dataset, which may appear counterintuitive given the general

superiority of deep learning for image classification. This outcome is explained by the regular, programmatically-generated texture patterns in synthetic knee X-ray images that are highly amenable to pixel-level decision tree splitting—a property that would not generalize to the complex, inter-patient anatomical variability of real clinical radiographs.

MobileNetV2 (85.6%) demonstrates strong transfer learning performance despite being pre-trained on natural images, validating that low-level visual features (edges, textures, shapes) learned on ImageNet generalize effectively to medical X-ray analysis. The 14-block freezing strategy prevents catastrophic forgetting of these generalizable features while enabling the last 5 blocks and custom head to adapt to knee-specific radiographic patterns.

The per-class accuracy analysis reveals a clinically meaningful pattern: MobileNetV2 achieves near-perfect accuracy on Severe (100%) and high accuracy on Normal (92%) and Moderate (94%), but struggles significantly with Doubtful (60%) and Mild (82%). This gradient of difficulty is not a model artifact—it directly reflects the clinical reality that Grade 1 (Doubtful) represents a genuinely ambiguous diagnostic category, characterized by features described as 'questionable' in the original KL grading schema. This finding aligns with published inter-observer agreement data showing the lowest radiologist concordance for Grades 1 and 2.

7.2 Clinical Significance

- Zero false negative rate for Severe Grade (100% recall) is of critical clinical importance, as missing Grade 4 KOA in a patient awaiting joint replacement surgery could delay life-altering intervention.
- The five-class probability visualization enables clinicians to identify borderline predictions where confidence is distributed across two adjacent grades, flagging cases that warrant additional specialist review.
- Prediction latency under 2 seconds on CPU hardware makes the system practical for primary care settings without dedicated GPU resources.

7.3 Limitations

- Synthetic training data cannot replicate the full complexity of real clinical radiographs—inter-patient anatomical variation, positioning differences, scanner variability, and artifact patterns—limiting the direct clinical applicability of current accuracy metrics.
- The 5-class KL grading problem is inherently ill-posed at grade boundaries (particularly 0/1 and 1/2), and even expert radiologists show 20–40% inter-observer disagreement on these adjacent grades.
- MobileNetV2 is not equipped with Grad-CAM explanations in the current implementation, limiting clinician ability to verify that predictions are based on clinically relevant radiographic features.

8. Conclusion and Future Scope

8.1 Conclusion

This paper presented a comprehensive deep learning-based web system for automated 5-class Kellgren-Lawrence grading of knee osteoarthritis from X-ray images. Three classification models—MobileNetV2 with selective transfer learning (85.6%), a Custom 4-layer CNN (83.2%), and Random Forest with 100 estimators (89.2%)—were rigorously trained and evaluated on a balanced synthetic dataset of 1,510 knee X-ray images. The mathematical framework, encompassing inverted residual block computation, depthwise separable convolution cost analysis, CrossEntropy loss optimization, selective fine-tuning strategy, and multi-class softmax inference, provides complete methodological transparency.

The MobileNetV2 per-class analysis reveals clinically meaningful accuracy patterns: Normal (92%), Doubtful (60%), Mild (82%), Moderate (94%), and Severe (100%), consistent with radiologist experience where intermediate grades are most difficult to classify. The system is deployed as a production-ready Flask web application with secure authentication, drag-and-drop X-ray upload, color-coded severity badges, five-class confidence visualization, SQLite prediction history, Chart.js analytics dashboard, and Docker containerization—demonstrating a complete pipeline from dataset generation through model training to clinical decision support deployment.

8.2 Future Scope

- Real Clinical Datasets: Retrain on OAI (36,369 images) and MOST (18,627 images) datasets with expert-validated KL grades, expected to significantly improve MobileNetV2 accuracy toward the 91–96% range reported in literature.
- Advanced Architectures: Implement DenseNet-169, EfficientNet-B4, and Vision Transformers (ViT) for comparison, targeting the 95.93% DenseNet-169 benchmark achieved on the MOST dataset.
- Grad-CAM Visualization: Integrate gradient-weighted class activation mapping to generate saliency heatmaps highlighting joint space, osteophyte, and sclerosis regions, building clinician trust through interpretable predictions.
- Ordinal Loss Function: Implement grade-distance-weighted ordinal loss (Chen *et al.*, 2019) to reduce the clinical severity of adjacent-grade misclassifications.
- DICOM Integration: Add DICOM-compatible interfaces for direct hospital

PACS system import, enabling automated screening of incoming knee radiographs.

- Multi-Joint Support: Extend to hip, hand, and spinal osteoarthritis, creating a comprehensive musculoskeletal AI diagnostic platform.
- Federated Learning: Train across multiple hospital sites without centralizing patient data, preserving privacy while accessing diverse multi-institutional training distributions.

References

- [1] Antony, J. et al. (2017). Quantifying radiographic knee OA severity using CNNs. ICPR 2017, pp.1195–1200.
- [2] Tiulpin, A. et al. (2018). Automatic knee OA diagnosis from plain radiographs: a DL approach. Scientific Reports, 8(1):1727.
- [3] Thomas, K.A. et al. (2020). Automated classification of radiographic knee OA severity using DNNs. Radiology: AI, 2(2):e190065.
- [4] Chen, P. et al. (2019). Fully automatic KOA severity grading with ordinal loss. CMIG, 75:84–92.
- [5] Swiecicki, A. et al. (2021). DL-based algorithm for knee OA assessment matches radiologist performance. CBM, 133:104334.
- [6] Bien, N. et al. (2018). MRNet: DL-assisted diagnosis for knee MRI. PLOS Medicine, 15(11):e1002699.
- [7] Kellgren, J.H. & Lawrence, J.S. (1957). Radiological assessment of osteo-arthrosis. ARD, 16(4):494–502.
- [8] Sandler, M. et al. (2018). MobileNetV2: Inverted residuals and linear bottlenecks. CVPR 2018, pp.4510–4520.
- [9] Howard, A.G. et al. (2017). MobileNets: Efficient CNNs for mobile vision. arXiv:1704.04861.
- [10] He, K. et al. (2016). Deep residual learning for image recognition. CVPR 2016, pp.770–778.
- [11] Rajpurkar, P. et al. (2017). CheXNet: Radiologist-level pneumonia detection on chest X-rays. arXiv:1711.05225.
- [12] Deng, J. et al. (2009). ImageNet: A large-scale hierarchical image database. CVPR 2009, pp.248–255.
- [13] Paszke, A. et al. (2019). PyTorch: An imperative style, high-performance DL library. NeurIPS 2019, pp.8026–8037.
- [14] Selvaraju, R.R. et al. (2017). Grad-CAM: Visual explanations from deep networks via gradient-based localization. ICCV 2017, pp.618–626.
- [15] Breiman, L. (2001). Random Forests. Machine Learning, 45(1):5–32.
- [16] Tajbakhsh, N. et al. (2016). CNNs for medical image analysis: Full training or fine tuning? IEEE TMI, 35(5):1299–1312.
- [17] Ioffe, S. & Szegedy, C. (2015). Batch Normalization: Accelerating deep network training. ICML 2015, pp.448–456.
- [18] LeCun, Y., Bengio, Y. & Hinton, G. (2015). Deep Learning. Nature, 521(7553):436–444.
- [19] Huang, G. et al. (2017). Densely connected convolutional networks. CVPR 2017, pp.4700–4708.
- [20] Dosovitskiy, A. et al. (2021). An image is worth 16×16 words: Transformers for image recognition. ICLR 2021.
- [21] Cross, M. et al. (2014). Global burden of hip and knee OA: GBD 2010 estimates. ARD, 73(7):1323–1330.
- [22] Tan, M. & Le, Q. (2019). EfficientNet: Rethinking model scaling for CNNs. ICML 2019, pp.6105–6114.
- [23] Goodfellow, I., Bengio, Y. & Courville, A. (2016). Deep Learning. MIT Press.