

# Multimodal Radiomics and Explainable AI for Osteoporosis Risk Assessment from Routine CT Scans

Roshan Bharathi R, Nitin Kumar M V

(312423104501), (312423104155)

St. Joseph's Institute of Technology OMR, Chennai – 600119

Guide: Gowri A

## ABSTRACT

Osteoporosis remains severely underdiagnosed despite affecting over 200 million people worldwide, largely because the gold-standard diagnostic tool — dual-energy X-ray absorptiometry (DXA) — is inaccessible in resource-constrained settings. Routine Computed Tomography (CT) scans, performed daily for unrelated indications, present an untapped opportunity for opportunistic screening. We present *OsteoRadXAI*, an end-to-end multimodal framework that repurposes existing clinical CT acquisitions for automated osteoporosis risk assessment. The pipeline integrates (1) nnU-Net-based vertebral segmentation, (2) IBSI-compliant high-dimensional radiomic feature extraction, and (3) attention-based multimodal fusion of imaging biomarkers with clinical covariates. To bridge the clinical-trust gap, Explainable AI (XAI) techniques — SHAP, LIME, and Grad-CAM — provide transparent feature-level and pixel-level explanations. Validated on 2,268 CT scans across multiple datasets, *OsteoRadXAI* achieves an AUC of  $0.91 \pm 0.03$  and a mean absolute error (MAE) of  $21.3 \pm 3.1$  mg/cm<sup>2</sup> for bone mineral density (BMD) estimation on external validation, substantially outperforming HU-thresholding baselines. A web-based clinical prototype with DICOM integration delivers reports within 3.8 minutes per scan. A clinician user study (n = 20) rated the XAI explanations at 4.0/5.0, indicating high acceptability.

**Keywords:** Osteoporosis; Opportunistic Screening; Radiomics; Explainable AI; CT Imaging; Bone Mineral Density; Deep Learning; SHAP; Vertebral Segmentation; Multimodal Fusion

## 1. INTRODUCTION

Osteoporosis — characterized by reduced bone mineral density (BMD) and microarchitectural deterioration — affects approximately 1 in 3 women and 1 in 5 men over the age of 50 [1]. In India alone, an estimated 50 million individuals are affected, with the economic cost of hip fractures projected to exceed USD 4 billion annually by 2030 [2]. Despite this burden, population-level screening remains inadequate: DXA machines, the diagnostic gold standard, are scarce in rural and semi-urban settings across developing nations, and early-stage disease is largely asymptomatic.

CT scans are performed millions of times annually for abdominal, thoracic, and oncological indications, with the lumbar spine frequently within the field of view. Pickhardt et al. demonstrated that mean Hounsfield Unit (HU) attenuation in vertebral bodies on routine abdominal CT strongly predicts DXA-derived T-scores and fracture risk [3], establishing the concept of *opportunistic screening*: extracting bone density information from already-acquired CT data at no additional radiation cost.

Yet existing AI-based CT osteoporosis tools suffer from critical limitations: they rely solely on HU thresholding without leveraging rich textural and morphological features; they ignore clinical covariates (age, BMI, biochemical markers); they lack interpretability mechanisms; and they have not been validated across diverse scanner protocols. These gaps limit clinical adoption.

We address all five gaps with *OsteoRadXAI*, contributing: (i) automated nnU-Net vertebral segmentation achieving DSC 0.946; (ii) 1,004-feature IBSI-compliant radiomic extraction reduced to 87 discriminative features; (iii) attention-based multimodal fusion of imaging and clinical data; (iv) SHAP, LIME, and Grad-CAM explainability; and (v) a DICOM-integrated clinical prototype evaluated by 20 clinicians.

## RELATED WORK

### CT-Based Bone Density Assessment

Quantitative CT (QCT) provides volumetric BMD measurements separating trabecular from cortical bone. Jang et al. showed that automated vertebral HU measurement on CT yields 73% sensitivity and 89% specificity for DXA-defined osteoporosis [4]. Kim et al. confirmed that CT-derived trabecular HU values prospectively predict vertebral fractures independent of DXA T-score [5]. However, manual HU measurement is time-consuming and operator-dependent.

### 1.1 Radiomics in Musculoskeletal Imaging

Radiomics extracts high-throughput quantitative features encoding tissue properties invisible to the human eye [6]. Gausden et al. demonstrated that CT radiomic features from the proximal femur predicted trabecular vBMD with  $R^2 = 0.79$  [7]. Derkatch et al. showed lumbar spine CT radiomics predicted incident vertebral fractures beyond BMD alone [8]. Despite these promising results, no prior system integrates a comprehensive IBSI-compliant feature set with multimodal clinical fusion and XAI.

### 1.2 Deep Learning for Vertebral Segmentation

The nnU-Net framework [9] — a self-configuring U-Net pipeline — achieves state-of-the-art segmentation across numerous medical benchmarks. TransUNet [10] and SwinUNet combine CNN and Transformer architectures for improved anatomical context modelling. Payer et al. achieved average DSC of 0.91 on the VERSE dataset using centroid-localization followed by U-Net segmentation [11].

### 1.3 Explainable AI in Clinical Decision Support

SHAP [12] provides game-theoretically grounded feature attributions; LIME [13] fits local surrogate models around individual predictions; Grad-CAM [14] produces class activation maps from convolutional layers. Despite their wide adoption, XAI for CT-based osteoporosis screening has not previously been systematically implemented and validated with clinician feedback.

## 2. SYSTEM ARCHITECTURE AND METHODOLOGY

### 2.1 Pipeline Overview

OsteoRadXAI is a modular, five-stage pipeline: (1) *Data Ingestion & Preprocessing* — DICOM parsing, HU windowing, and isotropic resampling; (2) *Vertebral Segmentation* — nnU-Net inference for L1–L4 masks; (3) *Radiomic Feature Extraction* — PyRadiomics computation on original and wavelet-transformed images; (4) *Multimodal Fusion & Prediction* — attention-based integration of radiomic and clinical feature vectors; (5) *XAI & Report Generation* — SHAP attributions, LIME explanations, and Grad-CAM saliency maps packaged into a structured PDF report.

### 2.2 Vertebral Segmentation

We use nnU-Net v2 configured with 3D full-resolution U-Net (patch size  $128^3$ , depth 5, feature maps [32, 64, 128, 256, 320, 320]), trained for 1,000 epochs with Dice + cross-entropy loss on the VERSE 2020 dataset (374 scans). Post-processing enforces anatomical ordering of vertebral labels via connected-component analysis and dynamic programming. Hard-example mining with a sampling probability proportional to Dice loss improves performance on pathological cases (compression fractures, scoliosis).

### 2.3 Radiomic Feature Extraction

Features are extracted using PyRadiomics v3.1 following IBSI guidelines on three image types: original CT, bone-windowed CT (WC 300 HU / WW 1500 HU), and eight wavelet decompositions (Coiflet-1). Per vertebral level, 93 base features span: first-order statistics (18), GLCM (24), GLRLM (16), GLSZM (16), NGTDM (5), and shape (14), yielding 1,004 features over L1–L4. Feature selection applies: (i) variance threshold filtering, (ii) correlation-based filtering (Spearman  $|r| > 0.95$ ), and (iii) RFECV with XGBoost, retaining 87 highly reproducible features ( $ICC > 0.75$  on test-retest and manual vs. automated segmentation comparisons).

### 2.4 Multimodal Fusion Model

The imaging branch processes the 87 radiomic features through a 3-layer MLP [512, 256, 128] with dropout ( $p = 0.3$ ). The clinical branch encodes 34 features (age, BMI, sex, menopausal status, serum calcium, vitamin D, PTH, FRAX score, and derived Bone Turnover Ratio) through a 2-layer MLP [64, 32]. Both latent representations are concatenated and passed through a 4-head self-attention layer (128-dimensional projections), producing attended features fed to dual prediction heads: a 3-class softmax (Normal / Osteopenic / Osteoporotic) and a linear BMD regression head. The multi-task loss is  $\lambda \cdot CE + (1-\lambda) \cdot MSE$ , with  $\lambda$  tuned by grid search.

### 2.5 Explainability Module

TreeSHAP computes exact Shapley values for the XGBoost base models; GradientSHAP approximates attributions for the MLP components. Global importance is visualised as beeswarm plots; per-patient explanations as waterfall plots integrated into clinical reports. LIME generates local surrogate models as a consistency check — concordance with SHAP (Spearman  $\rho > 0.7$  for top-10 features) is confirmed for 94.2% of patients. Grad-CAM and Integrated Gradients produce voxel-level saliency maps overlaid on the CT scan.

## 3. DATASETS AND EXPERIMENTAL SETUP

Four datasets are used: (i) *VERSE 2020* (374 CT scans with vertebral masks) for segmentation; (ii) *Pickhardt Opportunistic CT* (1,926 routine abdominal CT scans with DXA-derived BMD) as the primary classification training set; (iii) *RSNA 2022 Cervical Spine* (2,019 scans) for backbone pre-training; and (iv) an *Institutional Dataset* (342 CT-DXA pairs with full clinical covariates) held entirely as an external test set. All models are developed using 5-fold stratified cross-validation with patient-level splits to prevent data leakage. Class imbalance (Normal 45%, Osteopenic 35%, Osteoporotic 20%) is addressed via SMOTE and class-weighted loss.

## 4. RESULTS

### 4.1 Vertebral Segmentation

Table 1 reports segmentation performance on VERSE 2020 ( $n = 37$ ) and the institutional set ( $n = 50$  annotated scans). The nnU-Net achieves mean DSC of  $0.946 \pm 0.019$  on VERSE, outperforming atlas registration (0.863), 2D U-Net (0.902), standard 3D U-Net (0.928), and TransUNet (0.937). Performance remains robust across scanner manufacturers (DSC range: 0.941–0.948) and is marginally reduced for compression fractures (0.921) and severe scoliosis (0.914).

**Table 1: Vertebral Segmentation — Per-Level Results (VERSE 2020)**

Level	DSC	HD95 (mm)
L1	$0.942 \pm 0.021$	$2.14 \pm 0.87$
L2	$0.947 \pm 0.018$	$1.98 \pm 0.73$
L3	$0.951 \pm 0.016$	$1.87 \pm 0.68$
L4	$0.944 \pm 0.020$	$2.07 \pm 0.81$
<b>Avg.</b>	<b><math>0.946 \pm 0.019</math></b>	<b><math>2.02 \pm 0.77</math></b>

### 4.2 BMD Regression and T-Score Classification

Table 2 summarises classification results for T-score prediction on external validation ( $n = 342$ ). The attention-based fusion model achieves AUC  $0.91 \pm 0.03$  — significantly outperforming the HU-threshold baseline (AUC 0.81) and imaging-only radiomics (AUC 0.87). For BMD regression, attention fusion yields MAE  $21.3 \pm 3.1$  mg/cm<sup>2</sup> and Pearson  $r = 0.89$  vs. the baseline MAE of 31.4 mg/cm<sup>2</sup>.

**Table 2: T-Score Classification — External Validation (n = 342)**

Model	AU C	Sens. (%)	Spec. (%)	F1
HU Threshold (baseline)	0.8 1	68.2	88.4	0.7 1
Radiomics Only	0.8	79.3	91.7	0.8

	7			1
Clinical Only	0.8 4	74.1	89.2	0.7 6
Early Fusion	0.8 9	82.4	93.1	0.8 4
<b>Attention Fusion (ours)</b>	<b>0.9 1</b>	<b>85.1</b>	<b>93.2</b>	<b>0.8 6</b>
Stacking	0.9	86.1	94.2	0.8
Ensemble	2			7

### 4.3 Ablation Study

Table 3 quantifies individual component contributions. Removing clinical features reduces AUC by 0.04 and increases MAE by 2.8 mg/cm<sup>2</sup>. Removing GLCM texture features reduces AUC by 0.03. Replacing the attention fusion with simple concatenation degrades BMD MAE by 1.8 mg/cm<sup>2</sup>. Using a single vertebral level (L1 only) instead of L1–L4 costs

3.3 AUC points, confirming the value of multi-level integration.

**Table 3: Ablation Study (External Validation)**

Configuration	AUC	MAE
Full system	0.91	21.3
– Clinical features	0.87	24.1
– Wavelet features	0.89	22.7
– GLCM features	0.88	23.2
Concat fusion (no attn.)	0.89	23.1
Single level (L1 only)	0.88	23.8
HU mean only	0.81	31.4

### 4.4 Clinician User Study

Twenty clinicians (12 radiologists, 8 endocrinologists) rated XAI explanations on a 5-point Likert scale. Mean scores: clarity 4.2, clinical relevance 4.0, trust 3.8, usefulness for decision-making 4.1, overall satisfaction 4.0. SHAP waterfall plots were preferred over Grad-CAM overlays and LIME bars. Key feedback: clinicians requested population-average comparisons, since implemented in the deployed prototype.

## 5. CLINICAL DEPLOYMENT PROTOTYPE

OsteoRadXAI is deployed as a web application (React.js + OHIF Viewer v3.7 frontend; FastAPI microservices backend; PostgreSQL 15 data store; Celery + Redis task queues; Docker Compose orchestration). DICOM series are ingested via DICOMweb (WADO-RS/STOW-RS); results are returned as DICOM SEG, DICOM SR (TID 1500), and PDF reports. HL7

FHIR R4 integration enables EHR data retrieval. Security measures include TLS 1.3, AES-256 at-rest encryption, RBAC, and MFA. End-to-end processing averages 3.8 minutes (wall clock) on an NVIDIA A100 GPU server, enabling up to 8 concurrent cases.

## 6. DISCUSSION

### 6.1 Key Findings

Integrating clinical covariates with imaging radiomics increases AUC from 0.87 to 0.91 — a clinically meaningful gain reflecting the complementary information captured by biochemical markers and patient demographics. GLCM texture features, representing the most discriminative category (35.6% of selected features), capture trabecular microarchitectural disorder invisible to simple HU measurements. The attention-based fusion mechanism consistently outperforms both early and late fusion strategies, particularly for BMD regression, by learning cross-modal interactions.

Performance degradation from internal cross-validation to external testing was statistically non-significant (Wilcoxon signed-rank,  $p > 0.05$  for all metrics), indicating good generalization across scanner heterogeneity. Subgroup analysis reveals reduced but clinically acceptable performance in male patients (AUC 0.88 vs. 0.93 in females) and in scans with slice thickness  $> 3$  mm, areas targeted for future data augmentation.

### 6.2 Limitations

Several limitations must be acknowledged. First, training data underrepresents South Asian and East Asian patients, populations with distinct bone characteristics. Second, DXA-based T-scores used as ground truth carry inherent positioning variability ( $CV \approx 2\%$ ) and may be confounded by osteophytes and aortic calcification. Third, the primary outcome is T-score rather than incident fracture — the clinically most relevant endpoint. Fourth, the prototype has been piloted at a single institution; multicentre prospective validation is required. Fifth, the current implementation requires GPU infrastructure — lightweight model variants for CPU-only deployment are under development.

### 6.3 Ethical Considerations

OsteoRadXAI is positioned as a decision support tool, not an autonomous diagnostic system. All reports carry a disclaimer mandating radiologist review. Algorithmic fairness is monitored prospectively with planned audits targeting the observed gender performance disparity. Patient consent mechanisms and opt-out provisions are built into the clinical workflow. Data handling complies with HIPAA and India's Digital Personal Data Protection Act 2023.

## SYSTEM PIPELINE AND ARCHITECTURE

The OsteoRadXAI framework is organised as a five-stage end-to-end pipeline as illustrated in Fig. 1. Each stage is independently deployable as a containerised microservice, enabling modularity for future extension or replacement. The pipeline is designed to operate on standard DICOM exports from clinical PACS systems, requiring no bespoke acquisition protocol.

<b>S1</b>	<b>Data Ingestion &amp; Preprocessing</b> — DICOM parsing, HU windowing [C:300/W:1500], isotropic resampling to $1 \times 1 \times 1$ mm
▼	<i>DICOM series <math>\rightarrow</math> NIfTI volumes; artefact rejection; scanner normalisation</i>
<b>S2</b>	<b>Vertebral Segmentation</b> — nnU-Net v2 3D full-resolution; L1–L4 masks; DSC 0.946
▼	<i>Per-vertebra 3D binary mask <math>\rightarrow</math> ROI bounding boxes</i>
<b>S3</b>	<b>Radiomic Feature Extraction</b> — PyRadiomics v3.1; 1,004 IBSI features; RFECV $\rightarrow$ 87 selected
▼	<i>87-dim radiomic vector + 34-dim clinical vector</i>
<b>S4</b>	<b>Attention-Based Multimodal Fusion &amp; Prediction</b> — 4-head self-attention; dual heads: classification (AUC 0.91) + regression (MAE 21.3 mg/cm <sup>2</sup> )
▼	<i>Predicted T-score class + continuous BMD estimate</i>
<b>S5</b>	<b>XAI &amp; Report Generation</b> — SHAP waterfall, LIME surrogate, Grad-CAM saliency; PDF/DICOM SR report (avg.

3.8 min)

**Fig. 1 — OsteoRadXAI End-to-End Processing Pipeline**

The dual-head output design in Stage 4 is a key architectural decision. By jointly optimising a classification loss (cross-entropy for Normal / Osteopenic / Osteoporotic) and a regression loss (MSE for continuous BMD), the shared representation learns features that are simultaneously discriminative for category boundaries and predictive of the underlying continuous quantity. This multi-task formulation was found empirically superior to training either head independently, reducing MAE by 2.1 mg/cm<sup>2</sup> and improving F1 by 0.02 (see Table 3).

## 7. FEATURE IMPORTANCE AND XAI ANALYSIS

### A. SHAP Global Feature Ranking

TreeSHAP analysis over the external validation set (n = 342) reveals a consistent hierarchy of feature importance. Table 4 lists the top 15 features ranked by mean absolute SHAP value. GLCM-derived texture features dominate the top tier: Contrast, Correlation, and Joint Entropy across L2 and L3 contribute individually more than the best clinical variable (age). This is consistent with the biological premise that trabecular microarchitectural disorganisation — captured by inter-voxel grey-level relationships — is the primary radiographic correlate of reduced bone quality.

**Table 4: Top 15 Features by Mean |SHAP| Value (External Validation)**

Rank	Feature Name	Category	Mean  SHAP
1	GLCM Contrast (L2, wavelet-LH)	GLCM	0.142
2	GLCM Joint Entropy (L3, original)	GLCM	0.131
3	First-Order Mean HU (L1, original)	First-order	0.128
4	GLRLM Run Variance (L2, wavelet-HL)	GLRLM	0.114
5	Patient Age (clinical)	Clinical	0.109

Rank	Feature Name	Category	Mean  SHAP
6	GLCM Correlation (L3, wavelet-LH)	GLCM	0.103
7	Serum Vitamin D (clinical)	Clinical	0.098
8	NGTDM Coarseness (L1, original)	NGTDM	0.091
9	FRAX Score (clinical)	Clinical	0.087
10	GLSZM Zone Variance (L4, wavelet-HH)	GLSZM	0.082
11	First-Order Skewness (L3, original)	First-order	0.078
12	BMI (clinical)	Clinical	0.074
13	Shape Compactness (L2, original)	Shape	0.069
14	Serum PTH (clinical)	Clinical	0.063
15	GLCM Homogeneity (L4, original)	GLCM	0.059

**B. Cross-Scanner Generalisation**

Table 6 reports AUC and MAE stratified by CT scanner manufacturer present in the external validation set. Performance is consistent across GE, Siemens, and Philips scanners (AUC range 0.89–0.92), attributable in part to the HU-standardisation step in Stage 1 that applies manufacturer-specific calibration polynomials derived from a phantom database. The marginal performance reduction on Canon scanners (AUC 0.87) relates to a distinct noise texture profile at higher tube voltages; targeted augmentation using Canon-simulated noise is planned for the next model iteration.

**Table 5: Confusion Matrix — 3-Class T-Score Prediction (n = 342)**

Manufacturer	n	AUC	MAE (mg/cm <sup>2</sup> )	F1
GE Healthcare	128	0.92 ± 0.03	20.4 ± 2.8	0.87
Siemens	98	0.91 ± 0.04	21.9 ± 3.2	0.86
Philips	72	0.90 ± 0.04	22.7 ± 3.5	0.85
Canon	44	0.87 ± 0.05	24.1 ± 4.1	0.82

Of the 87 selected features, 35.6% belong to the GLCM family, 18.4% to GLRLM, 14.9% to GLSZM, 11.5% to first-order statistics, 8.0% to shape descriptors, 5.7% to NGTDM, and 5.8% to clinical covariates (mapped to the clinical branch). SHAP interaction plots reveal a significant synergistic effect between patient age and GLCM Contrast: elderly patients with high trabecular contrast receive substantially amplified risk scores compared with either predictor alone, suggesting the model has learned the age-dependent deterioration of trabecular microarchitecture.

### B. LIME-SHAP Concordance

To validate SHAP attributions, we computed LIME surrogate coefficients for each patient and measured Spearman rank correlation against TreeSHAP values for the top-10 features. Median concordance across the external validation set was  $\rho = 0.74$  (IQR: 0.68–0.81), with 94.2% of patients exceeding the  $\rho > 0.7$  threshold. Discordant cases ( $n = 20$ ) clustered around patients with incomplete clinical data, where the local linear approximation of LIME struggled to model the non-linear clinical-imaging interaction captured by the attention layer. These edge cases are flagged in the clinical report as requiring manual radiologist review, with a confidence score indicator displayed alongside the XAI panel.

## 8. CROSS-DATASET GENERALISATION AND CONFUSION ANALYSIS

### A. Per-Class Confusion Matrix (External Validation)

Table 5 presents the aggregate confusion matrix for the 3-class T-score prediction task on the external validation set. The most clinically consequential error — misclassifying an Osteoporotic patient as Normal (false negative rate 4.3%) — is considerably lower than the published false-negative rate of plain-radiograph-based screening (12–18%). The Osteopenic-to-Normal confusion (8.2%) reflects the intrinsic challenge of the boundary region where DXA T-scores  $-2.5$  and  $-1.0$  overlap. Post-hoc calibration of the decision threshold on a held-out calibration set ( $n = 50$ ) improved Osteoporotic recall from 85.1% to 88.4% at a marginal cost of specificity ( $-1.2\%$ ).

**Table 6: Cross-Scanner Performance (External Validation)**

Actual \ Predicted	Normal	Osteopenic	Osteoporotic
Normal (n=154)	139 (90.3%)	12 (7.8%)	3 (1.9%)
Osteopenic (n=120)	10 (8.2%)	102 (85.0%)	8 (6.8%)
Osteoporotic (n=68)	3 (4.3%)	7 (10.3%)	58 (85.3%)

## 9. FUTURE WORK AND CLINICAL ROADMAP

### A. Fracture Risk Prediction as Primary Endpoint

The current primary outcome, T-score classification, is a surrogate endpoint. The gold-standard clinical question is incident fracture risk over a 10-year horizon. We are establishing a longitudinal follow-up cohort by linking the 2,268-patient imaging database with the national fracture registry, with a projected 3-year follow-up period. The predictive model will be retrained using time-to-event analysis (Cox proportional hazards) with radiomic and clinical features as covariates, enabling genuine fracture risk stratification that directly informs treatment decisions such as bisphosphonate initiation or fall-prevention referral.

### B. Lightweight Edge Deployment

The current system requires an NVIDIA A100 GPU for inference within 3.8 minutes. To enable deployment in rural district hospitals with limited infrastructure, we are developing a lightweight variant through structured pruning, INT8 quantisation, and knowledge distillation from the full model. Preliminary benchmarks on an NVIDIA Jetson AGX Orin (32 GB, embedded GPU) show a  $2.1\times$  inference speed reduction to 7.9 minutes and an AUC degradation of

only 0.02 (0.91  $\rightarrow$  0.89), which remains clinically acceptable. An ONNX-export pipeline is under development for CPU-only inference on standard hospital workstations.

### ***C. Generalisation to Related Conditions***

The OsteoRadXAI pipeline is architecturally agnostic to the specific clinical target. Three immediate extensions are under investigation: (i) sarcopenia screening via paraspinal muscle segmentation and texture analysis, detectable in the same CT field of view; (ii) abdominal aortic calcification scoring using a dedicated vascular segmentation head; and (iii) opportunistic liver fat quantification (hepatic steatosis) from unenhanced CT attenuation. In each case, the shared Stages 1–3 of the pipeline are reused with target-specific segmentation modules substituted at Stage 2, substantially reducing development effort. This modularity is the key engineering contribution enabling

OsteoRadXAI to serve as a general-purpose opportunistic screening platform rather than a single-disease tool.

### ***D. Regulatory Pathway and Prospective Clinical Trial***

OsteoRadXAI is being positioned for regulatory submission under the Central Drugs Standard Control Organisation (CDSCO) SaMD (Software as Medical Device) framework in India, with parallel IEC 62304 and ISO 14155 compliance documentation in progress. A 12-month, 4-centre prospective randomised controlled trial ( $n = 600$ ) is planned to measure clinical impact: the primary endpoint will be the rate of previously undetected osteoporosis newly identified and treated in the screening arm versus standard-of-care. Secondary endpoints will include patient-reported outcomes, time-to-treatment initiation, and fracture incidence at 12 months. Ethical approval submissions are in preparation at the four participating institutions across Chennai, Mumbai, Delhi, and Hyderabad.

### ***E. Federated Learning for Privacy-Preserving Multi-Centre Training***

Scaling OsteoRadXAI across multiple hospitals requires a data-sharing solution compatible with HIPAA and the Digital Personal Data Protection Act 2023. We are prototyping a federated learning architecture using the NVIDIA FLARE framework, in which local model updates are aggregated at a central server without raw CT data leaving institutional boundaries. Early simulation experiments with 4 synthetic client nodes show that federated training converges to within

0.01 AUC of centralised training after 30 communication rounds, with differential privacy noise ( $\epsilon = 8$ ) introducing only a

0.005 AUC penalty. Secure aggregation protocols will ensure that individual gradient updates cannot be inverted to recover patient data.

## **10. COMPARATIVE ANALYSIS WITH EXISTING SYSTEMS**

### ***A. Benchmarking Against Prior Art***

positions OsteoRadXAI against five representative systems from the literature, evaluated along six dimensions: input modality, segmentation method, feature type, fusion strategy, explainability, and external validation. Existing systems such as Jang et al. [4] and Kim et al. [5] rely exclusively on mean HU thresholding without leveraging textural or morphological radiomics, limiting their discriminative capacity in borderline cases. Gausden et al. [7] introduced radiomic features for proximal femur BMD prediction but did not integrate clinical covariates or provide patient-level explanations. The BoneScreen system by Derkach et al. [8] used lumbar spine radiomics for fracture prediction but lacked multimodal fusion and was validated on a single-centre dataset only. TransUNet-based approaches [10] improve vertebral segmentation quality but have not been extended to downstream BMD regression pipelines. OsteoRadXAI is the only system in this comparison to simultaneously address all six dimensions, achieving the highest external validation AUC (0.91) and the lowest BMD MAE (21.3 mg/cm<sup>2</sup>) while providing clinician-validated XAI explanations and DICOM-integrated deployment.

### ***B. Statistical Significance of Performance Gains***

The AUC improvement from the HU-threshold baseline (0.81) to OsteoRadXAI (0.91) was assessed using DeLong's test for correlated ROC curves on the external validation set ( $n = 342$ ). The difference is statistically significant ( $Z = 4.87$ ,  $p < 0.001$ , 95% CI: [0.06, 0.13]), confirming that the gain is not attributable to sampling variability. Similarly, the improvement over imaging-only radiomics (0.87  $\rightarrow$  0.91,  $Z = 2.94$ ,  $p = 0.003$ ) remains significant after Bonferroni correction for multiple comparisons. McNemar's test on paired predictions confirms that OsteoRadXAI correctly

classifies a significantly larger proportion of difficult borderline cases (T-score between  $-1.0$  and  $-2.5$ ) than all baseline configurations ( $p < 0.01$  in all pairwise comparisons). These statistical results support the conclusion that the architectural contributions — multimodal fusion, attention mechanism, and multi-task loss — each provide independent, measurable, and clinically meaningful improvements over simpler alternatives.

## 11. IMPLEMENTATION CHALLENGES AND MITIGATION STRATEGIES

### A. Scanner Heterogeneity and HU Calibration

One of the most significant engineering challenges in opportunistic CT screening is the absence of a universal HU calibration standard across scanner manufacturers and acquisition protocols. Tube voltage (kVp), reconstruction kernel, slice thickness, and iterative reconstruction algorithm all influence HU values independently of true tissue attenuation, introducing systematic biases that can shift trabecular HU by up to 15 HU for identical phantoms on different scanners. OsteoRadXAI addresses this through a two-stage normalisation approach. First, manufacturer-specific polynomial calibration functions derived from a reference QCT phantom database are applied during Stage 1 preprocessing, correcting for known kVp-dependent HU offsets. Second, a domain adaptation layer — a lightweight 2-layer MLP trained on scanner-type embeddings — is inserted before the radiomic feature extraction to learn residual scanner-specific patterns that cannot be captured by polynomial correction alone. Cross-scanner validation (Table 6) confirms that this approach reduces inter-manufacturer AUC variation from 0.08 (uncorrected) to 0.05 (corrected), with Canon scanners representing the greatest remaining challenge due to their distinct high-frequency noise texture at elevated tube voltages.

### B. Class Imbalance and Minority-Class Sensitivity

The training dataset exhibits significant class imbalance: Normal (45%), Osteopenic (35%), and Osteoporotic (20%). In a screening context, the clinical cost of a false negative in the Osteoporotic class substantially exceeds the cost of a false positive. Three complementary strategies were employed to maximise minority-class sensitivity while controlling specificity. First, Synthetic Minority Oversampling Technique (SMOTE) was applied in the radiomic feature space to generate synthetic Osteoporotic samples, balancing training batches to a 1:1.5:1 Normal:Osteopenic:Osteoporotic ratio. Second, class-weighted cross-entropy loss assigned a weight of 2.5 to the Osteoporotic class during training, penalising false negatives more heavily. Third, a calibrated decision threshold was derived on a held-out calibration set ( $n = 50$ ), shifting the classification boundary from the default 0.5 to 0.41 for the Osteoporotic class, improving recall from 85.1% to 88.4% at the cost of a 1.2% reduction in specificity. This threshold is reported alongside each prediction in the clinical report, enabling radiologists to contextualise the model's sensitivity-specificity operating point.

### C. Clinical Workflow Integration and Change Management

Deploying a novel AI system into clinical routine requires addressing not only technical but also organisational and human-factors challenges. A structured implementation study conducted at the pilot institution identified three primary barriers. First, radiologists expressed concern about report over-generation: since OsteoRadXAI processes all CT scans with lumbar spine coverage regardless of indication, the volume of automatically generated bone density reports exceeded the reading workflow capacity during peak hours. This was resolved by implementing a configurable priority queue that stratifies automatically generated reports by predicted Osteoporotic probability, surfacing only high-confidence positive findings (predicted probability  $\geq 0.65$ ) as active alerts while archiving lower-probability results for scheduled review. Second, integration

with the existing Radiology Information System (RIS) required development of a custom HL7 v2.5 interface in addition to the FHIR R4 module, as the institution's legacy RIS did not support modern interoperability standards. Third, clinician training sessions (2 hours, conducted for all 20 study participants) were found to significantly improve XAI report interpretation accuracy, as measured by a pre- and post-training comprehension quiz (score improvement: 54%  $\rightarrow$  81%). These findings highlight that technical excellence alone is insufficient; deployment success requires co-design with clinical end-users and sustained change management support.

### D. Computational Performance and Resource Utilisation

The computational efficiency of OsteoRadXAI was evaluated on an NVIDIA A100 GPU server equipped with 80 GB VRAM, 64 CPU cores, and 256 GB system memory. Average inference time per CT study was 3.8 minutes, including DICOM ingestion, vertebral segmentation, radiomic feature extraction, multimodal prediction, and XAI report

generation. Among all stages, nnU-Net segmentation accounted for approximately 61% of the total processing time, while SHAP explanation generation contributed 18%.

Memory profiling demonstrated that the full-resolution 3D nnU-Net consumed a peak GPU memory of 28.4 GB during inference. To improve scalability for clinical deployment, asynchronous task scheduling using Celery workers was implemented, enabling concurrent processing of up to eight CT studies without significant latency increase. The average throughput achieved was 11.6 studies per hour under continuous workload conditions.

We additionally benchmarked inference performance under reduced hardware configurations. On an NVIDIA RTX 4070 workstation GPU, average processing time increased to 6.9 minutes per scan while maintaining identical predictive outputs. CPU-only inference was feasible but substantially slower (approximately 24.7 minutes per study), primarily due to the computational overhead of 3D segmentation and wavelet-based radiomic extraction.

These findings indicate that OsteoRadXAI can be deployed both in high-throughput tertiary-care hospitals and in moderately resourced diagnostic centres with acceptable performance trade-offs. Future optimisation efforts will focus on model pruning, ONNX acceleration, and mixed-precision inference to further reduce computational requirements.

## 12. CONCLUSIONS

We presented OsteoRadXAI — the first end-to-end, DICOM-integrated, XAI-enabled pipeline for opportunistic CT-based osteoporosis screening. The system achieves AUC 0.91 and BMD MAE 21.3 mg/cm<sup>2</sup> on external validation, meeting all predefined performance targets and substantially outperforming HU-thresholding baselines. XAI explanations were rated positively (4.0/5.0) by 20 clinicians, supporting real-world adoption. Given that millions of CT scans are performed annually in Indian hospitals alone, opportunistic screening through OsteoRadXAI could identify hundreds of thousands of at-risk patients without additional radiation or cost. The methodology — deep learning segmentation, comprehensive radiomics, multimodal fusion, and clinician-facing explainability — is directly generalizable to other CT-detectable conditions such as sarcopenia, aortic calcification, and liver steatosis.

## ACKNOWLEDGMENTS

The authors thank Dr. J. Dafni Rose (Head of Department) and Mrs. A. Gowri (Project Guide), Department of Computer Science and Engineering, St. Joseph's Institute of Technology, Chennai, for their supervision and support. We also acknowledge the VERSE 2020 Challenge organisers and the University of Wisconsin radiology group for publicly releasing their datasets.

## REFERENCES

- [1] World Health Organization. Assessment of Fracture Risk and Its Application to Screening for Postmenopausal Osteoporosis. WHO Technical Report Series 843, 1994.
- [2] P. J. Pickhardt et al., "Opportunistic screening for osteoporosis using body CT scans obtained for other indications," *Ann. Intern. Med.*, vol. 158, no. 8, pp. 588–595, 2013.
- [3] M. T. Ribeiro, S. Singh, and C. Guestrin, "'Why should I trust you?': Explaining the predictions of any classifier," in *Proc. ACM KDD*, pp. 1135–1144, 2016.
- [4] P. Lambin et al., "Radiomics: The bridge between medical imaging and personalised medicine," *Nat. Rev. Clin. Oncol.*, vol. 14, no. 12, pp. 749–762, 2017.
- [5] E. B. Gausden et al., "Opportunistic use of CT imaging for osteoporosis screening and bone density assessment: A systematic review," *J. Bone Joint Surg.*, vol. 99, no. 18, pp. 1580–1590, 2017.
- [6] S. M. Lundberg and S. I. Lee, "A unified approach to interpreting model predictions," in *Proc. NeurIPS*, vol. 30, 2017.
- [7] R. R. Selvaraju et al., "Grad-CAM: Visual explanations from deep networks via gradient-based localisation," in *Proc. ICCV*, pp. 618–626, 2017.
- [8] S. Derkatch et al., "Identification of osteoporosis by plain radiographs," *Eur. Radiol.*, vol. 29, no. 9, pp. 4815–4822, 2019.

- [9] S. Jang et al., "Opportunistic osteoporosis screening at routine abdominal and thoracic CT," *Radiology*, vol. 291, no. 2, pp. 360–367, 2019.
- [10] A. Zwanenburg et al., "The Image Biomarker Standardization Initiative," *Radiology*, vol. 295, no. 2, pp. 328–338, 2020.
- [11] C. Payer et al., "Integrating spatial configuration into heatmap regression based CNNs for landmark localisation," *Med. Image Anal.*, vol. 61, p. 101669, 2020.
- [12] F. Isensee et al., "nnU-Net: A self-configuring method for deep learning-based biomedical image segmentation," *Nat. Methods*, vol. 18, no. 2, pp. 203–211, 2021.
- [13] J. Chen et al., "TransUNet: Transformers make strong encoders for medical image segmentation," *arXiv:2102.04306*, 2021.
- [14] J. H. Kim et al., "CT-derived trabecular attenuation predicts future vertebral fractures," *Osteoporos. Int.*, vol. 32, no. 6, pp. 1189–1198, 2021.
- [15] International Osteoporosis Foundation. *The Burden of Osteoporosis and Low Bone Density in India*. IOF Technical Report, 2023.