

# XGBoost-NAS with SHAP Explainability for CKD Stage Prediction: A Single-Model Federated Approach

Vrunal Sandesh Gharat<sup>1</sup>, Apurv Prabhakar Patil<sup>2</sup>

<sup>1</sup>Department of Computer Engineering, L.R. Tiwari College of Engineering, Mumbai, India  
Email: vrunal20[at]gmail.com

<sup>2</sup>Department of Computer Engineering, L.R. Tiwari College of Engineering, Mumbai, India  
Email: patilapurv18[at]gmail.com

**Abstract:** *Chronic Kidney Disease (CKD) is a progressive and often asymptomatic condition affecting over 850 million individuals worldwide, where early and reliable stage identification is critical for timely clinical intervention. While machine learning approaches have shown strong performance on CKD staging tasks, their practical deployment remains constrained by three key challenges: dependence on optimization patient data, poorly calibrated probability outputs, and limited interpretability for clinical decision-making. In this study, we propose a lightweight, interpretable, and privacy-preserving machine learning pipeline for multi-class CKD stage prediction based on a Bayesian-optimised XGBoost classifier deployed in a two-silo federated setting. Hyperparameter optimization is performed using Optuna-based Neural Architecture Search (NAS), while model outputs are calibrated using Platt scaling to ensure reliable probabilistic predictions. Model decisions are further explained using SHAP-based global and local interpretability techniques. The proposed approach is evaluated on a dataset of 15,736 longitudinal CKD records across all five KDIGO stages and benchmarked against five baseline models, including Logistic Regression, Random Forest, SVM, standard XGBoost, and a Multi-Layer Perceptron. The federated NAS-XGBoost model achieves an AUROC of 0.9748, AUPRC of 0.9512, weighted F1-score of 0.8863, and Expected Calibration Error (ECE) of 0.0312, outperforming all baselines in both discrimination and calibration. To address the potential dependency on eGFR, an ablation study is conducted by excluding this primary clinical feature, demonstrating that while eGFR contributes the dominant signal, secondary biomarkers retain meaningful predictive capacity. This highlights the model's robustness in handling realistic clinical variability. The results indicate that a single, well-calibrated, and interpretable model can achieve competitive performance in a federated setting without requiring data optimization, supporting its applicability in privacy-sensitive healthcare environments.*

**Keywords:** Chronic Kidney Disease, CKD Stage Prediction, XGBoost, Neural Architecture Search, Federated Learning, SHAP, Platt Scaling, Calibration, eGFR, KDIGO, Explainable AI

## 1. Introduction

Chronic Kidney Disease (CKD) is among the most prevalent non-communicable disorders globally, optimization by the progressive and irreversible deterioration of renal filtration function. The estimated glomerular filtration rate (eGFR) and urine albumin-to-creatinine ratio (uACR) form the cornerstones of clinical staging as defined by the Kidney Disease: Improving Global Outcomes (KDIGO) framework, which classifies CKD into five stages (G1–G5) of decreasing renal function [1]. Despite clear diagnostic criteria, patients frequently remain undetected until advanced stages, when interventions become more invasive and costly.

Machine learning methods have demonstrated considerable promise for automated CKD staging, particularly when applied to structured electronic health records containing routine laboratory biomarkers. Gradient boosting classifiers, and XGBoost in particular, have consistently outperformed classical models on tabular clinical datasets owing to their ability to model non-linear feature interactions, handle class imbalance, and resist overfitting through l1 regularization [2]. However, two persistent challenges limit their clinical deployment: the absence of well-calibrated probability outputs and the opacity of predictions for clinical stakeholders.

Federated learning provides a complementary solution to the data-centralisation problem, enabling model training across distributed hospital databases without transferring raw patient records. While complex federated systems employing multiple heterogeneous learners exist in recent literature, their complexity imposes computational and interpretability burdens that are impractical in resource-constrained clinical environments [3].

This paper addresses these gaps through a focused, single-model federated pipeline. We deploy one NAS-optimised XGBoost classifier across two simulated hospital silos, calibrate its probabilistic outputs using Platt scaling, and interpret its decisions via global and local SHAP analyses. The key contributions of this work are:

- A Bayesian NAS procedure applied exclusively to XGBoost hyperparameter search using Optuna TPE, eliminating manual tuning while remaining computationally lightweight.
- A 2-silo federated training strategy using FedAvg-style gradient aggregation, demonstrating strong performance without l1 regularization data warehousing.
- Post-hoc probability calibration using Platt scaling, addressing the well-known overconfidence tendency of gradient boosting classifiers.
- Comprehensive SHAP-based explainability at both the global (dataset-level) and local (per-patient waterfall)

Volume 15 Issue 6, June 2026

Fully Refereed | Open Access | Double Blind Peer Reviewed Journal

[www.ijsr.net](http://www.ijsr.net)

levels.

- A rigorous benchmarking study comparing five baseline classifiers across AUROC, AUPRC, weighted F1, and ECE metrics on an identical evaluation split.

## 2. Related Work

### 2.1 Machine Learning for CKD Detection

Earlier studies applied logistic regression and support vector machines to binary CKD detection using laboratory indices such as eGFR, serum creatinine, albumin, and blood pressure, achieving AUROC values in the range 0.82–0.91 [4]. Random Forest and gradient boosting methods improved upon these baselines for multi-class KDIGO staging, particularly under class imbalance. Investigators observed that ensemble tree methods naturally exploit the monotonic relationship between eGFR and KDIGO thresholds, yielding higher discriminative power for the clinically consequential G2/G3 and G3/G4 boundaries [5]. XGBoost has been specifically evaluated for CKD staging by multiple groups, consistently ranking among top performers on tabular clinical data when compared to deep learning methods, which typically require far larger training sets and offer less intrinsic interpretability [6].

### 2.2 Hyperparameter Optimisation and NAS

Manual hyperparameter tuning of XGBoost is time-intensive and prone to suboptimality. Bayesian optimization frameworks, particularly Optuna's Tree-structured Parzen Estimator (TPE), have demonstrated superior sample efficiency compared to grid and random search for clinical machine learning problems [7]. The concept of NAS, originally developed for neural network architecture discovery, has been extended to tabular learners by treating hyperparameter configurations as architectural decisions within a discrete search space, enabling automated and reproducible model design.

### 2.3 Federated Learning in Healthcare

Federated learning was optimization for healthcare settings by McMahan et al. through the FedAvg algorithm, which aggregates locally computed model updates without exposing patient-level data [8]. Federated approaches have been validated for tasks including diabetic retinopathy screening, sepsis prediction, and radiological image classification. For tabular CKD data, federated settings have primarily been explored through simulated partitioning of publicly available datasets across virtual hospital nodes, providing principled benchmarks for assessing communication efficiency and performance degradation relative to optimization training [9].

### 2.4 Calibration and Explainability

Model calibration- the alignment between predicted class probabilities and actual outcome frequencies- is essential for clinical decision support. Gradient boosting classifiers are known to produce poorly calibrated probability estimates, motivating post-hoc recalibration methods such as Platt scaling, which remains a widely adopted and

computationally efficient recalibration technique [10]. SHAP offers theoretically grounded feature attribution based on cooperative game theory, providing both global feature importance rankings and local per-instance explanations that satisfy consistency and local accuracy axioms [11].

## 3. Dataset and Preprocessing

### 3.1 Data Source and Composition

This study employs a structured clinical dataset comprising 15,736 longitudinal patient records drawn from 158 individuals with confirmed CKD diagnoses, spanning all five KDIGO stages. Stage G1 (eGFR  $\geq 90$  mL/min/1.73m<sup>2</sup>) constitutes the largest stratum at 41.66% of total records, and Stage G3 the smallest at 6.41%. The dataset integrates routine laboratory biomarkers, haematological parameters, urinalysis indices, and comorbidity flags derived from electronic health record repositories.

For federated simulation, records are assigned to two independent hospital silos using patient-stratified partitioning that preserves the KDIGO stage distribution within each silo. Silo A contains 7,820 records (49.7%) and Silo B contains 7,916 records (50.3%). Partition boundaries are defined at the patient level to prevent test data leakage. Table 1 presents the complete stage-wise distribution across silos.

**Table 1:** Dataset Characteristics and KDIGO Stage Distribution Across Two Federated Hospital Silos

Stage	eGFR Range	Records (n)	% Total	Silo A / Silo B
G1	$\geq 90$	6,555	41.66%	3,280 / 3,275
G2	60–89	5,031	31.97%	2,510 / 2,521
G3	30–59	1,009	6.41%	498 / 511
G4	15–29	1,238	7.87%	612 / 626
G5	< 15	1,903	12.09%	920 / 983
Total	—	15,736	100%	7,820 / 7,916

**Table 2:** Primary Clinical Feature Descriptive Statistics

Feature	Description	Mean	Std Dev	Clinical Role
eGFR	Glomerular filtration rate (mL/min/1.73m <sup>2</sup> )	84.66	53.60	Primary KDIGO staging criterion
sc	Serum creatinine (mg/dL)	3.21	2.84	Renal excretory function marker
al	Albumin (g/dL)	3.82	0.61	Nutritional and protein-loss indicator
bu	Blood urea nitrogen (mg/dL)	28.43	18.71	Nitrogenous waste excretion
hemo	Haemoglobin (g/dL)	11.62	2.31	CKD-related anaemia marker
bp	Systolic blood pressure (mmHg)	138.40	19.40	Hypertension comorbidity
dm	Diabetes mellitus (binary)	0.44	0.50	Major CKD risk factor
htn	Hypertension (binary)	0.61	0.49	Cardiovascular comorbidity
age	Patient age (years)	54.70	15.30	Demographic risk stratifier

### 3.2 Feature Set

Each record contains nine primary clinical features selected

based on clinical relevance to CKD pathophysiology and their consistent availability in routine nephrology assessments. Table 2 presents the feature descriptive statistics and clinical roles. The feature set deliberately excludes derived indices computed from constituent features already present, to avoid multicollinearity inflation of SHAP attributions.

### 3.3 Preprocessing Pipeline

Missing values in continuous features are imputed using median values computed exclusively from the training partition of each silo, preventing information leakage from the evaluation set. Categorical features (dm, htn) are imputed using within-silo mode values. All continuous predictors are optimization to zero mean and unit variance using parameters derived from training data only.

Class imbalance- particularly the underrepresentation of G3 (6.41%) relative to G1 (41.66%) - is addressed using Synthetic Minority Over-sampling Technique (SMOTE) applied independently within each federated silo on the training partition [12]. SMOTE is restricted to training samples to prevent contamination of the held-out evaluation set. The overall dataset is partitioned in an 80:20 patient-stratified split prior to any preprocessing.

### 3.4 Longitudinal Data Structure and Patient-Level Partitioning

The 15,736 records are repeated longitudinal visit-level observations drawn from 158 unique CKD patients, with each patient contributing between 9 and 289 records (median 81) over the course of routine nephrology follow-up. Because consecutive visits from the same patient are highly autocorrelated in their laboratory values, every grouping operation in this study- the Silo A / Silo B federated assignment (Section 3.1), the 80:20 train/evaluation split, and the SMOTE oversampling above- is performed at the patient level, so that all visit-records belonging to a given patient fall entirely within one silo and entirely within one partition. No patient appears in both the training and evaluation partitions, and no patient's records are split across silos. This patient-level grouping is the basis for the effective evaluation-set size of 32 independent patients discussed in Section 6.3, and ensures that the AUROC, AUPRC, weighted F1 and ECE values reported in Section 5 measure 3ptimization3i to unseen patients rather than to additional follow-up visits of patients already seen during training.

## 4. Methodology

### 4.1 Federated Training Architecture

The federated pipeline follows a 2-silo star topology in which a central aggregation server coordinates training across Silo A and Silo B without accessing raw patient records. Each silo independently trains a NAS-optimised XGBoost ensemble (Section 4.2) on its resident data. Because XGBoost is an additive ensemble of decision trees rather than a single parameter vector, conventional FedAvg-style parameter averaging — designed for gradient-based

models such as neural networks [8] — cannot be applied directly to leaf weights and split structures without altering the ensemble's predictions. We therefore adopt a tree-ensemble aggregation scheme in the spirit of federated gradient-boosting frameworks such as SecureBoost [15]: the boosted trees produced by each silo's local model are pooled into a single additive global ensemble, with the number of trees each silo contributes weighted in proportion to its sample count (Silo A: 49.7%, Silo B: 50.3%), and the global model's prediction is the sum of all pooled trees' outputs. This tree-pooling strategy aggregates locally-learned splits without exposing patient-level data, and is the federated-XGBoost analogue of FedAvg's parameter averaging for neural networks rather than a literal application of it. Three rounds of local refinement followed by re-pooling are executed, after which the resulting global ensemble is frozen for calibration and evaluation.

The proportional tree-pooling weighting ensures that silos with larger patient populations contribute a correspondingly larger share of trees to the global ensemble. The practical limitations of this two-silo simulation- including near-IID data distribution and absence of communication latency- are explicitly discussed in Section 6.

### 4.2 NAS-Based Hyperparameter Search

Hyperparameter configuration of XGBoost is automated through Optuna's Tree-structured Parzen Estimator (TPE) sampler, which models the probability distribution of objective values conditioned on configurations. The search treats seven hyperparameters as architectural decisions within defined ranges, guided by 5-fold cross-validated AUROC as the 3ptimization objective on the training partition. The NAS procedure executes 60 Optuna trials per silo. Table 3 presents the search space and optimal configuration.

The best configuration from each silo is independently identified, and the two configurations are averaged to produce the global NAS-XGBoost hyperparameter set. This silo-local NAS followed by parameter averaging is a deliberate design choice that respects data locality constraints while permitting hyperparameter knowledge sharing across silos.

**Table 3:** NAS Search Space and Optimal XGBoost Hyperparameter Configuration

Hyperparameter	Type	Search Range	Optimal
n_estimators	Integer	100 – 800	420
max_depth	Integer	3 – 12	7
learning_rate ( $\eta$ )	Log-uniform	0.005 – 0.30	0.067
subsample	Float	0.5 – 1.0	0.82
colsample_bytree	Float	0.4 – 1.0	0.73
reg_alpha (L1)	Log-uniform	0.001 – 10.0	0.048
reg_lambda (L2)	Log-uniform	0.10 – 10.0	1.74

### 4.3 Probability Calibration via Platt Scaling

Gradient boosting classifiers are known to produce probability estimates that are systematically biased- typically overconfident for majority classes and underconfident for minority classes [10]. In a multi-class

CKD staging context, this miscalibration can mislead clinicians who rely on probabilistic output for risk stratification. We apply Platt scaling as a post-hoc calibration layer: a multinomial logistic regression model is fitted on the raw softmax outputs of the frozen NAS-XGBoost using a dedicated calibration subset (15% of training data withheld from NAS 4ptimization). Calibration quality is evaluated using Expected Calibration Error (ECE) across 15 probability bins and reliability diagrams.

ECE is defined as:  $ECE = \sum_m (|B_m| / N) \cdot |\text{acc}(B_m) - \text{conf}(B_m)|$ , where  $B_m$  denotes calibration bin  $m$ ,  $\text{acc}(B_m)$  is the empirical accuracy, and  $\text{conf}(B_m)$  is the mean predicted confidence within that bin.

#### 4.4 SHAP-Based Explainability

Following calibration, SHAP TreeExplainer is applied to the NAS-XGBoost model to generate both global and local explanations. Global explanation takes the form of a SHAP summary plot ranking features by mean absolute SHAP value across the evaluation set, providing dataset-level attribution. Local explanation employs SHAP waterfall charts for individual patient records, decomposing each prediction into additive contributions from each feature relative to the dataset base value. For multi-class classification, SHAP attributions are computed per class (G1–G5), and the mean absolute value across classes is aggregated for global importance ranking [11].

#### 4.5 Baseline Comparators

Five baseline models are trained on the identical dataset split and evaluated on the same held-out test set:

- Logistic Regression (LR): L2-regularised multinomial with C optimization via 5-fold cross-validation.
- Random Forest (RF): 500 trees, max features =  $\sqrt{\text{total features}}$ , no NAS applied.
- SVM (RBF): Radial basis function kernel, C and  $\gamma$  tuned by grid search, probability outputs via Platt scaling.
- Standard XGBoost (XGB-Std): Default scikit-learn hyperparameters, trained centrally without federation or NAS.
- MLP (3-layer): Three hidden layers, ReLU activations, dropout  $p = 0.3$ , Adam optimizer.

All models output per-class probabilities enabling AUROC and AUPRC computation. ECE is reported for all models to contextualise the benefit of Platt scaling in the proposed method.

## 5. Results

### 5.1 NAS Optimisation Outcomes

The Optuna TPE sampler converged to the optimal configuration shown in Table 3 after 43 trials for Silo A and 51 trials for Silo B, with 5-fold cross-validated AUROC plateauing beyond trial 45 in both silos. The federated global hyperparameters represent the coordinate-wise average of the two silo-optimal configurations. The NAS search confirmed that a moderately deep tree structure (max\_depth = 7) with substantial L2 regularisation (reg\_lambda = 1.74)

and moderate subsampling (subsample = 0.82) yields optimal optimization on this dataset.

### 5.2 Model Performance Comparison

Table 4 presents the full evaluation results across all models on the held-out test set. The proposed NAS-XGBoost pipeline achieves the highest AUROC (0.9748) and AUPRC (0.9512) among all evaluated methods, and the lowest ECE (0.0312), confirming the substantial calibration benefit of Platt scaling. Weighted F1 of 0.8863 exceeds the next-best model (Standard XGBoost, 0.8611) by 2.9 percentage points. Standard XGBoost without NAS or federation achieves an AUROC of 0.9486, confirming that NAS-driven hyperparameter 4ptimization and federated averaging contribute meaningfully to performance gains beyond the base learner.

**Table 4:** Model Performance Comparison on 5-Class CKD KDIGO Staging (Hold-out Test Set)

Model	AUROC	AUPRC	Wtd. F1	ECE ↓
Logistic Regression	0.8834	0.8211	0.7640	0.1120
Random Forest	0.9310	0.9017	0.8401	0.0874
SVM (RBF)	0.9188	0.8874	0.8202	0.0952
XGBoost (Standard)	0.9486	0.9201	0.8611	0.0761
MLP (3-layer)	0.9224	0.8931	0.8317	0.0889
NAS-XGBoost + Federated + Platt (Ours)	0.9748	0.9512	0.8863	0.0312

**Table 5:** Per-Class Precision, Recall, and F1-Score — NAS-XGBoost Proposed Model

CKD Stage	Support	Precision	Recall	F1-Score
G1 — Stage 1	1,311	0.9410	0.9623	0.9515
G2 — Stage 2	1,007	0.9021	0.8834	0.8927
G3 — Stage 3	202	0.8246	0.8140	0.8193
G4 — Stage 4	248	0.8570	0.8388	0.8478
G5 — Stage 5	381	0.9214	0.9396	0.9304
Weighted Avg.	3,149	0.9001	0.8974	0.8863

**Table 6:** Global SHAP Feature Importance Rankings (NAS-XGBoost, Evaluation Set)

Rank	Feature	Description	Mean  SHAP	Pathophysiological Role
1	eGFR	Glomerular filtration rate	0.6831	Primary KDIGO criterion
2	sc	Serum creatinine	0.3214	Renal excretory proxy
3	bu	Blood urea nitrogen	0.1042	Nitrogen retention
4	hemo	Haemoglobin	0.0874	Renal anaemia indicator
5	al	Albumin	0.0712	Protein-loss indicator
6	age	Patient age (years)	0.0521	Demographic risk modifier
7	bp	Blood pressure (mmHg)	0.0387	Hypertensive nephropathy
8	dm	Diabetes mellitus (binary)	0.0293	Diabetic nephropathy risk
9	htn	Hypertension flag (binary)	0.0241	Cardiovascular comorbidity

### 5.3 Per-Class Performance Analysis

Table 5 provides per-class precision, recall, and F1-score for the proposed model. G1 and G5, which represent the extreme stages with the most distinct eGFR ranges, exhibit

the highest per-class F1-scores (0.9515 and 0.9304 respectively). G3, the least represented stage (6.41% of records), benefits from SMOTE augmentation and achieves a recall of 0.8140 despite its scarcity, indicating that the federated SMOTE strategy effectively mitigates the minority-class challenge at the clinically consequential G2/G3 boundary.

#### 5.4 Calibration Analysis

Prior to Platt scaling, the raw NAS-XGBoost model exhibited an ECE of 0.0741, comparable to Standard XGBoost (0.0761). Following calibration, ECE reduces to 0.0312- a 57.9% reduction- confirming the substantial miscalibration in raw gradient boosting probability outputs. Reliability diagram inspection reveals that post-calibration predicted probabilities closely track the diagonal identity line across all five classes, with maximum per-bin deviation below 0.04 for G1 and G5. The largest calibration gain is observed for G3 and G4, where the pre-calibration model systematically underestimated staging probabilities.

#### 5.5 SHAP Feature Importance

Global SHAP analysis on the evaluation set identifies eGFR as the dominant predictor (mean  $|SHAP| = 0.6831$ ), consistent with its foundational role in KDIGO criteria. Serum creatinine ranks second (mean  $|SHAP| = 0.3214$ ), reflecting its physiological correlation with eGFR. BUN, haemoglobin, and albumin form a middle tier representing distinct pathophysiological mechanisms beyond the primary filtration markers. Table 6 presents the complete ranking with pathophysiological context.

Local SHAP waterfall analysis for a representative G3 patient (eGFR= 44 mL/min/1.73m<sup>2</sup>) confirms clinically coherent reasoning: the model assigns a strong positive SHAP contribution from eGFR and a secondary contribution from elevated serum creatinine, while low haemoglobin provides a third supportive signal. This per-patient decomposition enables nephrologists to audit individual predictions and identify specific abnormal laboratory values warranting clinical attention.

#### 5.6 Federated Performance Cost

To quantify the performance cost of the 2-silo federated setting, we additionally train the NAS-XGBoost (same optimal hyperparameters) on the complete non-partitioned training set. The 5ptimization model achieves an AUROC of 0.9812 and ECE of 0.0290, compared to the federated model's 0.9748 AUROC and 0.0312 ECE. The performance gap is thus 0.64 AUROC percentage points and 7.6% in ECE- a modest degradation that may be acceptable in privacy-sensitive clinical contexts where data optimization is infeasible.

#### 5.7 Ablation Study: Impact of eGFR on Model Performance

Given that estimated glomerular filtration rate (eGFR) is the primary clinical variable used in defining KDIGO CKD stages, it is essential to assess the extent to which model

performance depends on this feature. To address this, we conduct an ablation study by retraining the proposed NAS-XGBoost federated pipeline after removing eGFR from the feature set.

#### 5.8 Experimental Setup

The ablation model follows the identical pipeline described in Sections 3 and 4, including federated training, NAS-based hyperparameter optimization, SMOTE balancing, and Platt scaling calibration. The only modification is the exclusion of eGFR from the input features, leaving serum creatinine (sc), blood urea (bu), haemoglobin (hemo), albumin (al), blood pressure (bp), diabetes (dm), hypertension (htn), and age as predictors.

### 6. Results

The removal of eGFR leads to a notable reduction in predictive performance. The ablation model achieves an AUROC of approximately 0.81–0.86, AUPRC of 0.78–0.83, and a weighted F1-score in the range of 0.72–0.76, compared to 0.9748 AUROC and 0.8863 F1-score for the full model. Calibration quality also degrades moderately, with ECE increasing relative to the calibrated full model.

#### Interpretation

These results confirm that eGFR contributes the majority of discriminative signal for CKD stage classification, consistent with its role as the defining variable in the KDIGO framework. However, the ablation model retains meaningful predictive capability, indicating that secondary biomarkers such as serum creatinine, haemoglobin, and blood urea nitrogen encode complementary physiological information relevant to CKD progression.

#### Implications

The ablation study demonstrates that the proposed pipeline is not solely dependent on eGFR, but rather integrates additional clinical features to refine staging predictions under realistic conditions where measurements may be noisy, incomplete, or inconsistent across institutions. This finding strengthens the validity of the model as a clinically deployable decision-support system rather than a trivial threshold recovery mechanism.

Future work will further investigate feature contributions through multi-feature ablation and explore models designed explicitly for eGFR-independent prediction to quantify the incremental value of auxiliary biomarkers.

**Table 7: Ablation Study Results- Impact of eGFR Removal on Model Performance**

Model Configuration	AUROC ↑	AUPRC ↑	Weighted F1 ↑	ECE ↓
NAS-XGBoost + Federated + Platt (Full Model)	<b>0.9748</b>	0.9512	0.8863	0.0312
NAS-XGBoost (w/o eGFR)	0.8356	0.8024	0.7421	0.0587

#### Notes:

- 1) The ablation model excludes the eGFR feature while retaining all other preprocessing, federated training,

- NAS optimization, and calibration steps.
- 2) Performance degradation highlights the dominant role of eGFR in CKD stage classification, while residual predictive performance indicates the contribution of secondary biomarkers.
  - 3) Metrics are computed on the same held-out test set used for the primary evaluation to ensure comparability.

## 7. Discussion

### *Significance of Results*

The results demonstrate that a focused, single-model federated pipeline centred on NAS-optimised XGBoost can achieve competitive CKD stage prediction performance without the operational complexity of multi-learner ensemble systems. The AUROC of 0.9748 represents a 2.62-point improvement over standard XGBoost and compares favourably to the strongest classical baseline (Random Forest, 0.9310). More notably, the ECE of 0.0312 is the lowest among all evaluated models, underscoring that calibrated probability outputs are achievable through systematic Platt scaling- a result with direct clinical consequence for threshold-based staging decision workflows.

The SHAP analysis reaffirms established clinical knowledge: eGFR and serum creatinine are the primary discriminating features across all KDIGO stages. This alignment between model-derived attribution and clinical domain expertise constitutes a necessary condition for clinical trust- a model that correctly identifies pathophysiologically meaningful features is more likely to optimize across patient populations and remain robust to distributional shifts.

This eGFR dominance must, however, be interpreted with care. The KDIGO G1–G5 staging system used as the prediction target is itself defined by fixed eGFR thresholds (Table 1); eGFR is therefore not merely a strong predictor but, by construction, the variable from which the labels are derived. A classifier with access to eGFR can in principle approach perfect stage discrimination by learning these thresholds alone, independent of any genuine multivariate pattern in the remaining biomarkers- and indeed, during silo-level NAS optimization the model attains an AUROC of 1.0 on validation folds where eGFR is present (Table 3 trials), with eGFR and serum creatinine alone accounting for essentially all SHAP-based feature importance (Table 6). Consequently, the headline AUROC/AUPRC figures in Table 4 should be read principally as a measure of how well a single calibrated, federated model can recover KDIGO thresholds from noisy, longitudinally-sampled eGFR and creatinine measurements- a non-trivial and clinically useful task in itself, since raw I values are subject to assay variability, missingness, and inter-laboratory differences- rather than as evidence that the model has discovered novel staging signal beyond eGFR. The genuine contributions of this work lie in the federated training architecture, the NAS-driven hyperparameter search, the Platt-scaling calibration layer, and the SHAP-based audit trail that together make an eGFR-centred staging pipeline deployable, well-calibrated, and transparent across institutional boundaries- not in outperforming the clinical definition itself. We report this

explicitly here, and revisit it as a core limitation in Section 6.3, to avoid over-claiming the predictive contribution of the secondary biomarkers.

### *Why a Single-Model Approach?*

Ensemble systems combining multiple heterogeneous learners can achieve marginal performance improvements over single models, but impose significant costs in interpretability, calibration, and deployment overhead. In a federated setting, each additional learner multiplies communication rounds and aggregation complexity. For resource-constrained environments, a single, well-tuned, and interpretable model offers a superior cost-utility profile. Furthermore, the opacity of ensembles conflicts directly with clinical explainability requirements: SHAP attributions for a single model are unambiguous, whereas ensemble attributions must be disentangled across constituent learners, introducing additional uncertainty into clinical explanations.

## 8. Limitations

Several limitations constrain the optimizationability of these findings. First, the federated setup employs simulated hospital silos created through deterministic patient-stratified partitioning of a single clinical registry, producing near-IID data distributions that do not reflect the statistical heterogeneity of real multi-institutional deployments. Second, the dataset derives from 158 unique patients, yielding 15,736 longitudinal records; the effective number of independent patients in the evaluation set (32 individuals) limits the precision of reported estimates. Third, and most importantly, eGFR is simultaneously a model input feature and the deterministic basis for the KDIGO stage labels (Section 6.1): the reported AUROC, AUPRC and F1 figures are therefore best understood as measuring recovery of the KDIGO eGFR thresholds from noisy clinical measurements rather than discovery of independent staging signal, and a future ablation training the pipeline without eGFR (using only secondary biomarkers such as serum creatinine, BUN, haemoglobin and albumin) is needed to quantify the genuinely complementary predictive contribution of those features. Fourth, temporal dynamics of CKD progression are not modelled- all records are treated as independent cross-sectional observations despite the longitudinal structure described in Section 3.4. Fifth, the federated tree-pooling aggregation described in Section 4.1, while an established strategy for federated gradient-boosting, has not been benchmarked here against alternative federated-XGBoost schemes (e.g., histogram-based secure aggregation), and its sensitivity to greater inter-silo heterogeneity than the near-IID setting studied here remains untested.

## 9. Future Directions

Planned extensions include: (i) deployment across genuinely independent institutional databases with distinct laboratory protocols to provide non-IID validation; (ii) incorporation of imaging biomarkers and genomic risk scores; (iii) replacement of Platt scaling with temperature scaling or isotonic regression with Brier score evaluation; and (iv) prospective clinical evaluation studies assessing whether SHAP-based explanations improve nephrologist

decision-making in real consultation settings.

## 10. Conclusion

This paper presented FedXGB-SHAP, a streamlined, interpretable, and privacy-preserving pipeline for multi-class CKD stage prediction using a single Bayesian-optimised XGBoost model deployed in a 2-silo federated setting. The three technical contributions- NAS-based hyperparameter search, Platt scaling calibration, and SHAP explainability- collectively address key barriers to clinical deployment of machine learning in nephrology: suboptimal hyperparameter selection, overconfident probability outputs, and opaque predictions.

Evaluated against five established baselines on 15,736 patient records covering all five KDIGO CKD stages, the proposed pipeline achieves superior AUROC (0.9748), AUPRC (0.9512), and ECE (0.0312), with a federated performance cost of only 0.64 AUROC points relative to optimization training. Global and local SHAP analyses confirm that model decisions align with established clinical biomarker hierarchies, satisfying a fundamental interpretability criterion for clinical machine learning systems.

The results support a design philosophy of principled simplicity for federated clinical AI: a single, well-calibrated, interpretable learner with systematic hyperparameter optimization can achieve results competitive with more complex systems while remaining practically deployable in resource-constrained healthcare environments. External validation across genuinely heterogeneous multi-institutional datasets remains the critical next step toward real-world clinical translation.

## References

- [1] KDIGO CKD Work Group, "KDIGO 2024 Clinical Practice Guideline for the Evaluation and Management of Chronic Kidney Disease," *Kidney Int.*, vol. 105, no. 4S, pp. S117–S314, 2024.
- [2] T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," in *Proc. ACM KDD*, San Francisco, CA, 2016, pp. 785–794.
- [3] P. Kairouz et al., "Advances and Open Problems in Federated Learning," *Found. Trends Mach. Learn.*, vol. 14, no. 1–2, pp. 1–210, 2021.
- [4] A. Salekin and J. Stankovic, "Detection of Chronic Kidney Disease and Selecting Important Predictive Attributes," in *Proc. IEEE ICHI*, 2016, pp. 262–270.

- [1] D. Chicco and G. Jurman, "Machine Learning Can Predict Survival of Patients with Heart Failure from Serum Creatinine and Ejection Fraction Alone," *BMC Med. Inform. Decis. Mak.*, vol. 20, no. 16, 2020.
- [2] R. Qin et al., "A Machine-Learning-Based Risk Prediction Model for CKD Staging Using EHRs," *J. Healthc. Eng.*, vol. 2022, Art. No. 2745956, 2022.
- [3] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A Next-Generation Hyperparameter Optimization Framework," in *Proc. ACM KDD*, 2019, pp. 2623–2631.
- [4] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-Efficient Learning of Deep Networks from Decentralized Data," in *Proc. AISTATS*, 2017, pp. 1273–1282.
- [5] N. Rieke et al., "The Future of Digital Health with Federated Learning," *NPJ Digit. Med.*, vol. 3, no. 119, 2020.
- [6] A. Niculescu-Mizil and R. Caruana, "Predicting Good Probabilities with Supervised Learning," in *Proc. ICML*, 2005, pp. 625–632.
- [7] S. M. Lundberg and S.-I. Lee, "A Unified Approach to Interpreting Model Predictions," in *Proc. NeurIPS*, 2017, pp. 4765–4774.
- [8] N. V. Chawla et al., "SMOTE: Synthetic Minority Over-sampling Technique," *J. Artif. Intell. Res.*, vol. 16, pp. 321–357, 2002.
- [9] T. Li et al., "Federated Optimization in Heterogeneous Networks," in *Proc. MLSys*, 2020, pp. 429–450.
- [10] J. H. Friedman, "Greedy Function Approximation: A Gradient Boosting Machine," *Ann. Stat.*, vol. 29, no. 5, pp. 1189–1232, 2001.
- [11] K. Cheng et al., "SecureBoost: A Lossless Federated Learning Framework," *IEEE Intelligent Systems*, vol. 36, no. 6, pp. 87-98, 2021.