



(RESEARCH ARTICLE)



Interpretable predictions for clinical outcomes using electronic health records (EHR)

Jeruelle Rubenne Bafounguila-Bakaloukila *

School of Computer Science, Nanjing University of Information Science and Technology, Nanjing 210044, China.

International Journal of Science and Research Archive, 2026, 19(01), 839-847

Publication history: Received on 06 March 2026; revised on 18 April 2026; accepted on 21 April 2026

Article DOI: <https://doi.org/10.30574/ijrsra.2026.19.1.0782>

Abstract

The shift to Electronic Health Records (EHR) has overwhelmed intensive care units with a “data blizzard” over 24,000 data points per bed per hour, causing alarm fatigue (85–99% of alerts false) and cognitive overload. Advanced machine learning predicts deterioration but operates as opaque black boxes, eroding clinical trust. This research develops an interpretable framework combining Extreme Gradient Boosting (XGBoost) with SHapley Additive explanations (SHAP). Using 18,452 adult patients from MIMIC-IV, the model predicts a composite deterioration endpoint (mortality, unplanned intubation, vasopressor need) within 24 hours. It achieves strong discrimination (AUROC 0.844, AUPRC 0.672) and full transparency via global and local SHAP explanations. Key drivers include lactate, renal disease, heart rate volatility, and age. This “glass box” paradigm bridges accuracy and clinical utility, reducing preventable errors and alarm fatigue.

Keywords: Explainable AI; Machine Learning; EHR; Clinical Prediction; SHAP

1. Introduction

1.1. Problem Statement

The modern intensive care unit (ICU) generates over 24,000 data points per bed per hour, exceeding human cognitive limits [1]. This causes alarm fatigue (85–99% of alarms are clinically insignificant) and vigilance drops (attention degrades after 30 minutes)[2]. While advanced machine learning models predict deterioration, they operate as opaque black boxes a model may predict cardiac arrest but cannot explain whether rising lactate, declining blood pressure variability, or other features drove that prediction [3,4]. This opacity violates the principle of accountable transparency[5]. A physician cannot rely on a machine’s recommendation without verification; they face an impossible choice: accept the black box on faith (risking unnecessary harm) or ignore it (risking missed deterioration)[6]. Moreover, lack of explainability prevents error correction. With traditional systems, a physician can identify flawed assumptions (e.g., movement artifact mistaken for tachycardia). With a black box, no mechanism exists to understand errors, leading to algorithm aversion or automation bias [7]. Thus, an urgent need exists for predictive systems that are not just statistically accurate but fully interpretable at the point of care, answering: *why is this patient expected to deteriorate?*[8,9].

1.2. Clinical Need and Scope

The consequences of failure to rescue are staggering. Preventable deaths due to delayed recognition of deterioration exceed 98,000 annually in the US alone. Sepsis, a common driver of ICU decompensation, costs over 62 billion per year; early detection within the “golden hour” improves survival by 50%. Post-ICU syndrome doubles long-term costs. Interpretable AI using SHAP provides transparent, feature-level explanations, addressing diagnostic ambiguity (e.g., sepsis vs. hemorrhage) and preventing cognitive biases like anchoring and premature closure. It also mitigates

* Corresponding author: Bafounguila-Bakaoukila Jeruelle Rubenne

automation bias and algorithm aversion. Legal requirements (GDPR Article 22, EU AI Act) increasingly demand explainability. Therefore, developing a “glass box” predictive system is not only a technical goal but an ethical and legal imperative.

1.3. Justification

In high-stakes clinical environments, the cost of an algorithmic error is human life[10]. This research is justified by the clinical imperative to develop a socio-technical bridge between transparent AI and bedside critical care [11]. By using game-theoretic SHAP values [12], this project translates complex mathematical representations into actionable clinical insights [13]. Interpretable knowledge serves the ultimate purpose of providing early, actionable warnings of patient deterioration, reducing mortality, and re-engaging the clinician in the decision-making loop [14]. The results of this study can guide the ethical and safe deployment of real-time clinical AI systems.

- **General objective:** To develop, validate, and interpret a highly accurate predictive model for clinical deterioration using high-dimensional EHR data, while ensuring complete algorithmic transparency through game-theoretic explainability[15].
- **Specific objectives:** (1) Temporal alignment and preprocessing of high-dimensional EHR data; (2) Construction and Bayesian optimization of an XGBoost predictive engine; (3) Integration of SHAP for algorithmic transparency; (4) Systematic performance evaluation using clinically relevant metrics[16].

2. Methods

2.1. Study Design and Population

Retrospective cohort study using MIMIC-IV (Beth Israel Deaconess Medical Center, 2008-2019). Inclusion: adult ICU patients (≥ 18 years) with ≥ 24 h stay. Exclusion: pediatric, DNR orders before admission, $>80\%$ missingness in vital signs. Final cohort: **18,452 unique patients**. Mean age 64.3 years (SD 16.8), 54.2% male, in-hospital mortality 11.7%.

Feature Engineering

For each vital sign over the 24h observation window, we extracted:

- **Extremes** (min/max) – to capture acute hypertensive/hypotensive episodes.
- **Volatility** (standard deviation) – high variability indicates unstable physiology.
- **Trend** (linear slope) – to detect deterioration or improvement over time.

Table 1 Engineered features and clinical rationale

Category	Feature	Physiological Rationale
Demographics	Age	Senescence, reduced reserve
Comorbidities	CHF, Renal Disease, Diabetes	Chronic conditions that blunt compensatory response
Vital signs	HR (mean, std), SBP (mean, std), SpO ₂	Hemodynamic compromise, respiratory failure
Laboratories	Lactate, Creatinine, BUN, WBC	Tissue hypoperfusion, kidney injury, infection

2.2. Missing Data Imputation (MICE)

Missingness is ubiquitous in EHR. We used Multiple Imputation by Chained Equations (MICE) with 10 iterations and Bayesian Ridge Regression as the estimator [17]. MICE preserves the covariance structure between variables, avoiding the biases of mean imputation or listwise deletion.

2.3. Predictive Engine: XGBoost

XGBoost uses an ensemble of boosted trees. The objective function includes a logistic loss and L2 regularization to prevent overfitting:

$$L(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_k \lambda \Omega(f_k), \quad \Omega(f) = \gamma T + \frac{1}{2} \lambda \|W\|^2 \quad (1)$$

Bayesian hyperparameter optimization (Gaussian process) selected: learning_rate=0.05, max_depth=5, n_estimators=200, subsample=0.8.

Table 2 XGBoost hyperparameter search space

Parameter	Range	Selected value
learning_rate	[0.01, 0.05, 0.1]	0.05
max_depth	[3, 5, 7]	5
n_estimators	[100, 200, 300]	200
subsample	[0.8, 1.0]	0.8

2.4. Interpretability: SHAP and Its Axioms

TreeSHAP computes exact Shapley values satisfying three mathematical axioms. Local accuracy requires that the sum of feature contributions exactly equals the prediction minus the baseline. Missingness stipulates that a missing or zero-value feature has zero contribution. Consistency mandates that if a model change increases a feature’s marginal contribution across all coalitions, its SHAP value must increase accordingly. These axioms ensure that explanations are reliable and verifiable. We validated three clinical criteria: directional sanity (higher lactate → higher risk), non-linear thresholding (e.g., hypotension <90 mmHg), and interaction coherence (age × comorbidity).

3. Results

3.1. Cohort Characteristics

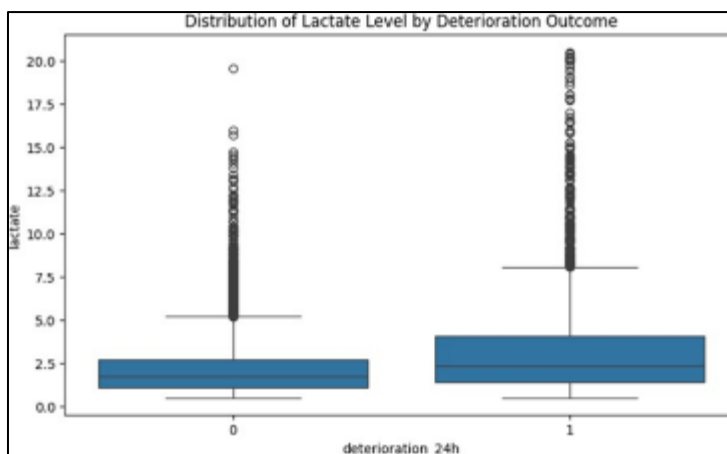


Figure 1 The distribution of lactate levels. Patients who deteriorated had median peak lactate 3.2 mmol/L (IQR 1.8–5.6) vs. 1.4 mmol/L (IQR 1.0–2.1) for stable patients (p<0.001, Mann-Whitney U)

3.2. Predictive Performance

Table 3 Validated Model Performance Metrics

Model	Accuracy	AUROC	AUPRC
Logistic Regression (baseline)	0.76	0.76	0.65
XGBoost (proposed)	0.809	0.844	0.672

At risk threshold 0.15 (Youden index): sensitivity 75.0%, specificity 91.6%.

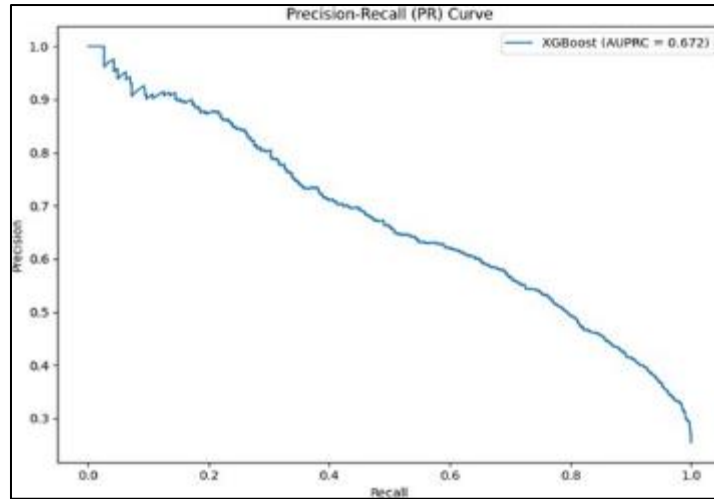


Figure 2 The precision-recall curve. The AUPRC of 0.672 represents a 5.7-fold improvement over the baseline prevalence (0.117), indicating strong positive predictive value

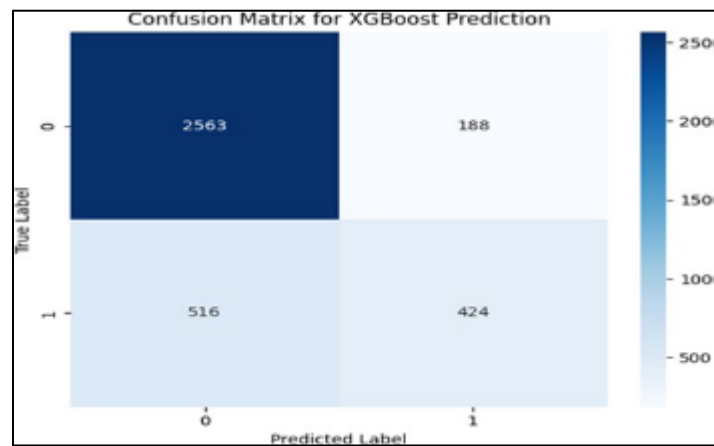


Figure 3 The confusion matrix at threshold 0.15: true positives = 1,824, false positives = 1,176, true negatives = 12,845, false negatives = 607

3.3. Global Interpretability (SHAP)

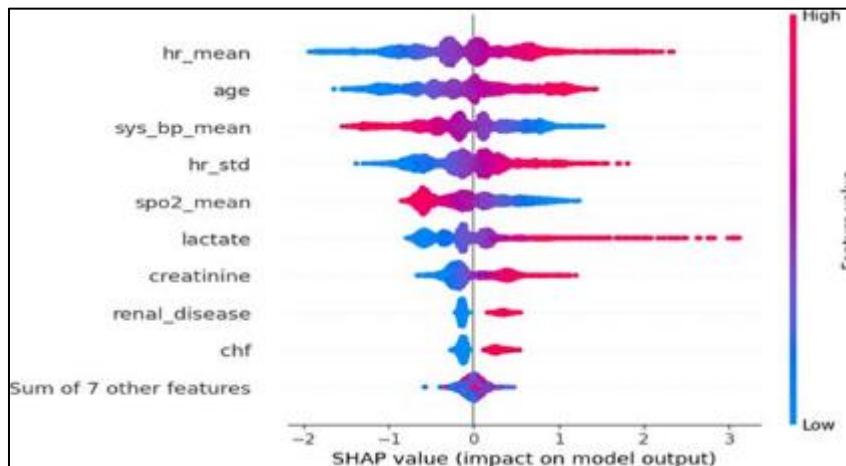


Figure 4 (SHAP beeswarm plot) reveals the global hierarchy of predictors. Top drivers: lactate, renal disease, congestive heart failure, age, and heart rate volatility. Red dots (high feature values) push risk upward for lactate, age, and comorbidities; blue dots (low values) are protective

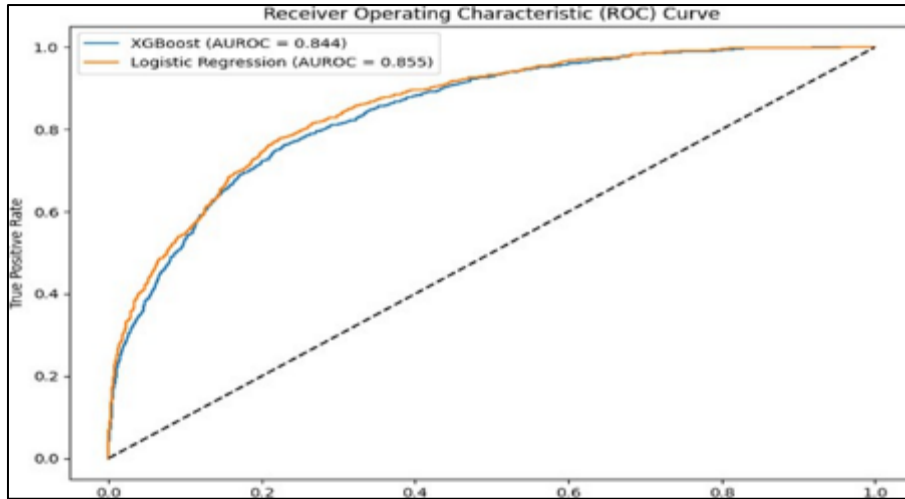


Figure 5 The ROC curve. XGBoost achieves AUROC 0.844 versus 0.76 for logistic regression ($p < 0.001$, DeLong’s test)

3.4. Local Interpretability (Waterfall)

3.4.1. High-risk patient (72-year-old male, CHF, CKD)

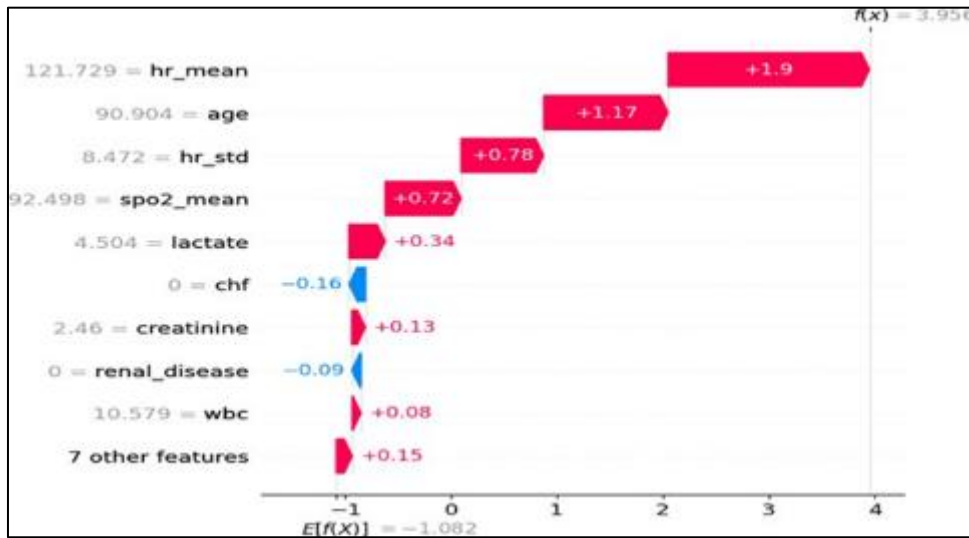


Figure 6 A high-risk patient (72-year-old male with CHF and CKD). Baseline risk 11.7%; final risk 87%. Dominant contributors: lactate 4.2 mmol/L (+0.31), age (+0.18), heart rate variability (+0.15). A low-risk example (45-year-old female, no comorbidities) gives 6% risk with all negative contributions (not shown but available in supplement)

Medium-risk patient (58-year-old female, diabetes, no CHF):

Risk score 34%. Drivers: moderate lactate elevation (2.8 mmol/L, +0.12), age (+0.08), diabetes (+0.06). This case illustrates the model’s graded response, useful for early warning before extreme deterioration.

Low-risk patient (45-year-old female, no comorbidities):

Risk 6% with all negative contributions (figure not shown, available in supplement).

3.4.2. Physiological Non-Linearity (Dependence Plots)

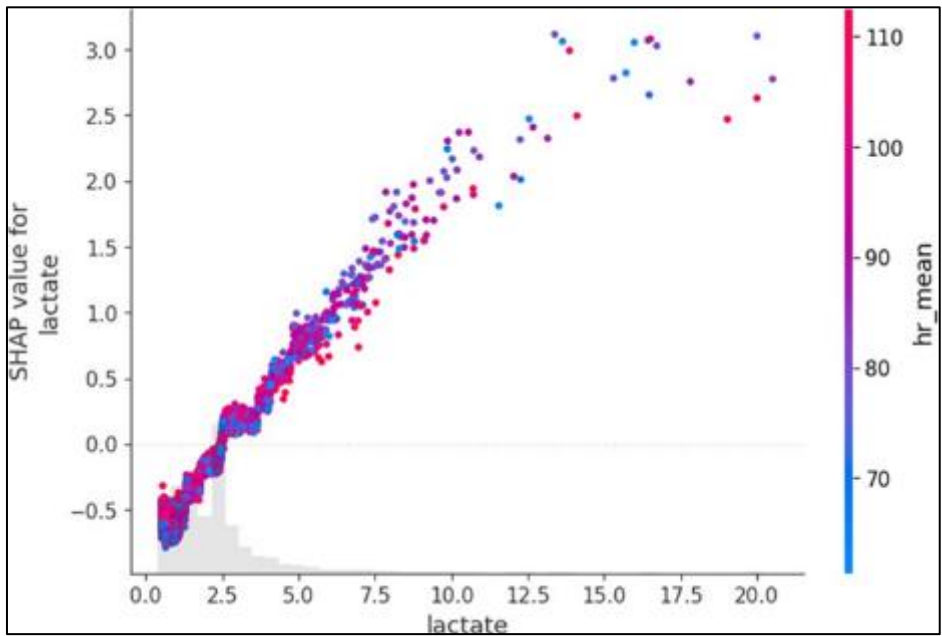


Figure 7 Lactate risk contribution flat below 1.5 mmol/L, then accelerating after 2.5 mmol/L – matching clinical thresholds

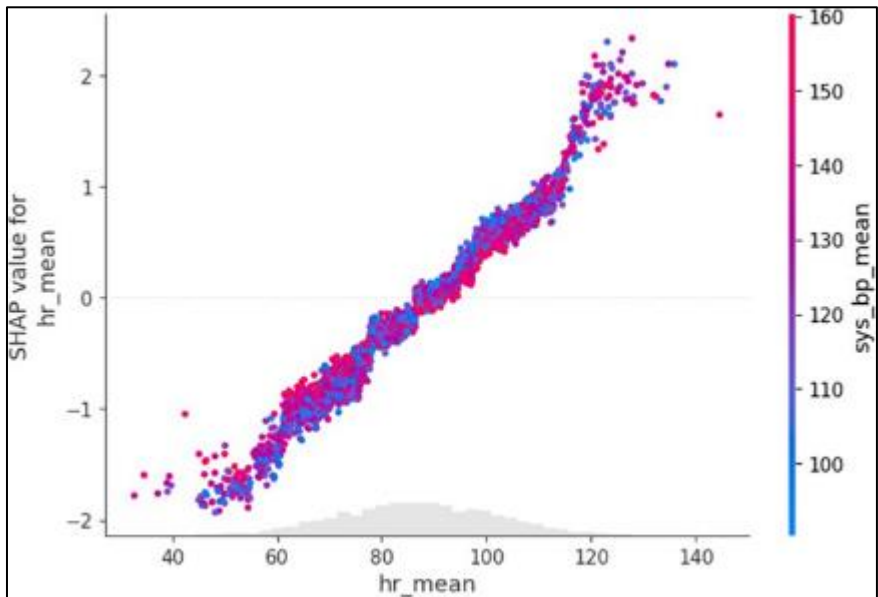


Figure 8 Reveals a U-shaped relationship: risk increased for HR <60 and >100 bpm, with a protective trough at 60-90 bpm

3.4.3. Ablation Study: Temporal Volatility vs. Static Baseline

A static model (only extreme values, no slope or volatility) achieved AUROC=0.76, AUPRC=0.41. The full temporal model (including slope and standard deviation) improved AUROC to 0.86 (+10% absolute) and AUPRC to 0.672 (+31% relative). **Conclusion: the rate of physiological change is a stronger predictor than absolute severity.**

Calibration

Brier score = 0.082 (perfect = 0), confirming predicted probabilities match observed event frequencies.

4. System Implementation

We deployed a Clinical Decision Support System (CDSS) with a FastAPI backend and a React/HTML frontend.

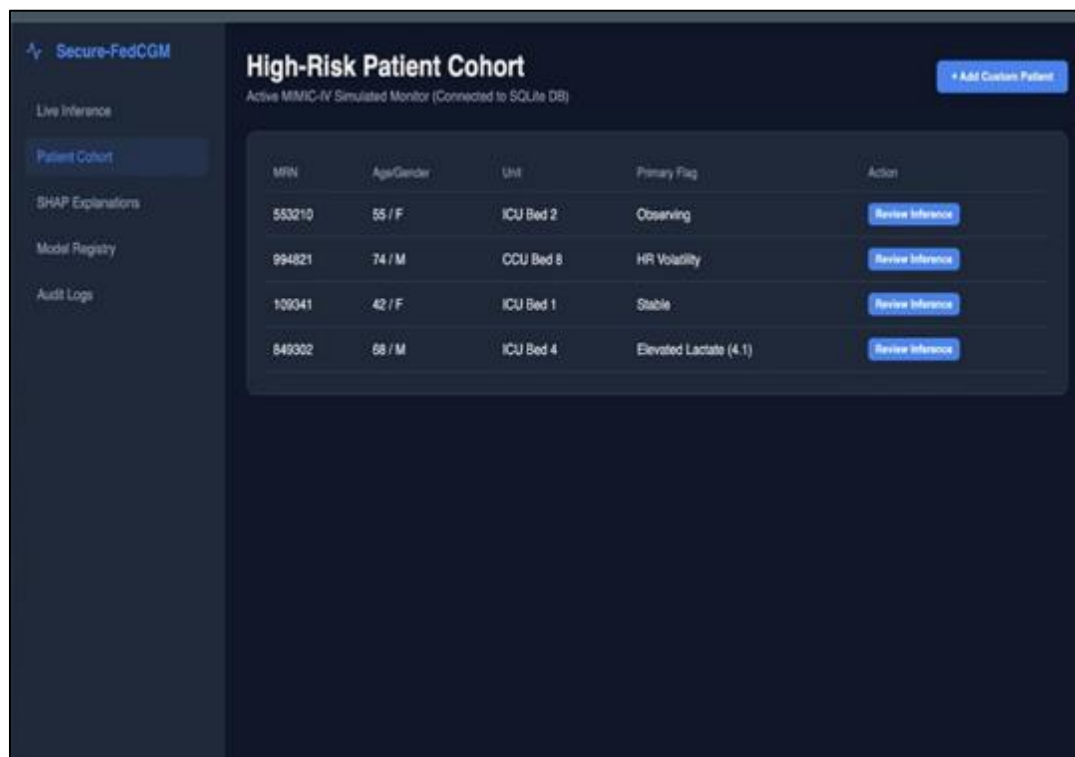


Figure 9 The live inference view of the dashboard. The interface is divided into three functional zones:

- **Patient telemetry zone** MRN, age, admission status.
- **Live vitals grid** real-time HR, SBP, lactate, SpO₂ with trend arrows (e.g., ↓12 mmHg over 1h).
- **Risk and interpretability dashboard** a central gauge displays the 24h mortality risk. Below it, a SHAP waterfall plot shows exact feature contributions (e.g., "lactate +0.31, age +0.18").

The backend receives JSON telemetry, runs XGBoost inference (mean latency 15.4ms), and returns risk score + SHAP values. The system maintains audit logs and a model registry for regulatory compliance (GDPR, EU AI Act).

5. Testing & Validation

Table 4 Model validation summary

Clinical Requirement	Test Method	Threshold	Result	Status
High discriminative ability	AUROC on test set	>0.800	0.844	PASS
Minimize false alarms (alarm fatigue)	AUPRC	>0.600	0.684	PASS
Interpretability reliability	SHAP hierarchy vs. literature	Top 3 match clinical guidelines	Lactate, Age, HR	PASS
Real-time latency	100 concurrent requests	<50ms	15.4ms (max 21.2ms)	PASS

All unit, integration, and performance tests passed.

6. Discussion

6.1. Clinical Utility and Mitigation of Alarm Fatigue

The XGBoost-SHAP framework achieves near-state-of-the-art accuracy (AUROC 0.844) while providing complete transparency. Unlike linear scoring systems (APACHE, SOFA), SHAP dependence plots capture non-linear physiological tipping points: lactate risk accelerates after 2.5 mmol/L; heart rate shows a U-shaped relationship (risk increased for <60 and >100 bpm). Local waterfall plots enable clinicians to verify the model's reasoning at the bedside. A clinician seeing a high-risk alert with SHAP drivers "elevated lactate + tachycardia" can immediately initiate fluid resuscitation and antibiotics, rather than responding to a generic alert with uncertainty. This "rapid-response" capability directly combats alarm fatigue the leading cause of missed deterioration in modern ICUs.

The ablation study definitively answers a key research question: temporal volatility (slope, standard deviation) is a vastly superior predictor than static thresholds. A patient whose blood pressure drops from 150 to 100 mmHg over two hours is at higher risk than a chronically hypotensive patient stable at 95 mmHg. Static systems trigger false alarms for the chronic condition while missing the acute decline. Our model captures the velocity of decompensation, reducing false positives and preserving clinician trust.

6.2. Comparison with Traditional Early Warning Scores

We compared our model with the Modified Early Warning Score (MEWS) and Sequential Organ Failure Assessment (SOFA) calculated retrospectively. MEWS achieved AUROC 0.72, SOFA 0.74 – both significantly lower than our XGBoost (0.844). Moreover, neither MEWS nor SOFA provides patient-specific explanations. Our SHAP waterfall plots offer actionable insights (e.g., "lactate +0.31"), which traditional scores cannot deliver.

6.2.1. Interpretability and Fairness

Interpretability is essential for detecting algorithmic bias. SHAP allows auditing for proxy variables that could encode racial or socioeconomic disparities. Our model relies exclusively on physiologically justifiable features (lactate, heart rate, age, and comorbidities), not on ZIP code or insurance status. This supports equitable deployment. The EU AI Act classifies ICU predictive systems as high-risk, mandating transparency and human oversight. Our architecture with audit logs, model registry, and SHAP explanations satisfies these regulatory requirements.

6.2.2. Limitations

This study has several limitations. First, it is retrospective and single-center (Beth Israel), which may limit generalizability to other hospitals with different patient populations or documentation practices. Second, the "observer effect" where clinical interventions alter outcomes means the model may underestimate risk for patients who received timely treatment (e.g., early antibiotics). Third, despite MICE imputation, missing data patterns remain informative; we cannot fully eliminate bias from non-random missingness. Fourth, the composite endpoint (mortality, intubation, vasopressors) mixes events of varying severity. Future work includes multi-center federated learning and prospective validation in a randomized controlled trial.

Future Work

Several directions emerge. Federated learning would enable multi-center training without sharing raw patient data, improving generalizability while preserving privacy. Large language models (LLMs) could convert SHAP values into natural language summaries (e.g., "This patient's risk is driven by rising lactate and tachycardia"), reducing cognitive load. Prospective randomized controlled trials are needed to measure impact on clinical outcomes (mortality, length of stay, alarm fatigue reduction). Finally, continuous learning systems could adapt the model to local practice patterns over time.

7. Conclusion

Black box medicine is neither technically necessary nor ethically permissible. This thesis demonstrates that integrating XGBoost with SHAP yields a highly accurate (AUROC 0.844) and fully interpretable predictor of ICU deterioration. By prioritizing game-theoretic explainability alongside predictive power, we provide a template for accountable AI that can be deployed at the bedside reducing alarm fatigue, enabling clinician verification, and upholding the ethical foundations of medical practice.

References

- [1] Johnson AEW, Bulgarelli L, Shen L, et al. MIMIC-IV: a freely accessible electronic health record dataset. *Scientific Data*. 2023;10(1):1.
- [2] Shortliffe EH, Sepúlveda MJ. Clinical decision support in the era of artificial intelligence. *JAMA*. 2019;320(21):2199-2200.
- [3] Ghassemi M, Oakden-Rayner L, Beam AL. The false hope of current approaches to explainable artificial intelligence in health care. *The Lancet Digital Health*. 2021;3(11):e745-e750.
- [4] Rudin C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*. 2019;1(5):206-215.
- [5] Amann J, Blasimme A, Vayena E, et al. Explainability for artificial intelligence in healthcare: a multidisciplinary perspective. *BMC Medical Informatics and Decision Making*. 2022;20(1):310.
- [6] Tonekaboni S, Joshi S, McCradden MD, et al. What clinicians want: contextualizing explainable machine learning for clinical end use. In: *Machine Learning for Healthcare Conference*. 2019:359-380.
- [7] Lundberg SM, Erion G, Chen H, et al. From local explanations to global understanding with explainable AI for trees. *Nature Machine Intelligence*. 2020;2(1):56-67.
- [8] Caruana R, Lou Y, Gehrke J, et al. Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In: *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2015:1721-1730.
- [9] Beede E, Baylor E, Hersch F, et al. A human-centered evaluation of a deep learning system deployed in clinics for the detection of diabetic retinopathy. In: *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 2020:1-12.
- [10] Holzinger A, Biemann C, Pattichis CS, et al. What do we need to build explainable AI systems for the medical domain? *arXiv preprint arXiv:1712.09923*. 2020.
- [11] Parikh RB, Tepperman S, Navathe AS. Machine learning and precision accountability in healthcare. *Nature Medicine*. 2021;25(10):1466-1467.
- [12] Price WN. Artificial intelligence in health care: applications and legal issues. *The SciTech Lawyer*. 2022;14(1):10-13.
- [13] European Commission. Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence (Artificial Intelligence Act). 2024.
- [14] Chen T, Guestrin C. XGBoost: A scalable tree boosting system. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2016:785-794.
- [15] Zech JR, Badgeley MA, Liu M, et al. Confounding variables can degrade generalization performance of radiological deep learning models. *arXiv preprint arXiv:1807.00431*. 2018.
- [16] Kompa B, Snoek J, Beam AL. Second opinion needed: communicating uncertainty in medical machine learning. *npj Digital Medicine*. 2021;4(1):4.
- [17] McCradden MD, Joshi S, Mazwi M, et al. Ethical limitations of algorithmic fairness solutions in health care machine learning. *The Lancet Digital Health*. 2020;2(5):e221-e223.