

An Explainable AI Decision Support System for Optimising Hospital Bed Allocation: A Predictive Modelling Study

¹*Emon Hasan

¹Department of Information Technology, Washington University of Science and Technology, Alexandria, Virginia (VA), United States.

Abstract

Background: The efficient and equitable allocation of hospital resources, specifically critical care beds, represents a major operational and ethical challenge in healthcare. Traditional predictive models for this task often suffer from a “black box” nature, hindering clinician trust and adoption.

Aim: The study aimed to develop, validate, and evaluate an Explainable AI Decision Support System (XAI-DSS) framework for optimising hospital bed allocation by integrating high-accuracy predictive modelling with stakeholder-centric XAI.

Methodology: The study employed a sequential mixed-methods approach, beginning with a retrospective cohort analysis of the MIMIC-IV database to train and validate three Machine Learning models (XGBoost, Random Forest, Neural Network). A randomized controlled user study and a simulated operational trial then evaluated the XAI utility and operational impact.

Results: The Random Forest model achieved near-perfect predictive accuracy, recording an AUC-ROC of 0.999998 and, critically, a 1.00 Recall with zero False Negatives. Global SHAP identified Creatinine and Respiratory Rate as key drivers, confirming clinical objectivity. The XAI-DSS achieved a ‘Good’ System Usability Scale score (Target >70) and demonstrated a 15% reduction in mismatched bed allocations in the simulated trial.

Conclusion: The XAI-DSS provides a safe, trustworthy, and operationally efficient framework for critical resource management, validating the necessity of multi-modal interpretability for clinical adoption.

Future Recommendation: Future studies should focus on a prospective, real-world shadowed deployment to externally validate the 15% efficiency gain.

Keywords: Explainable Artificial Intelligence (XAI); Clinical Decision Support Systems (CDSS); Hospital Bed Management; Resource Allocation; Machine Learning; Shapley Additive Explanations (SHAP); Predictive Modelling.

1 Introduction

1.1 Background

The equitable and efficient management of hospital resources, particularly bed capacity, is a critical operational challenge in modern healthcare systems. Research conducted by Marengo et al. (2025) provides evidence that optimizing resource allocation through Machine Learning (ML) can contribute to the overall energy efficiency of public hospitals. The significant regional disparities and inequality in the distribution of health resources, such as beds and personnel, are common and persistent challenges (Dong et al., 2021; Masroor et al., 2024). To address these inefficiencies, Decision Support Systems (DSS) leveraging advanced predictive analytics have emerged as a powerful solution for proactive resource planning. Building on prior research by Almeida et al. (2024), a comprehensive review demonstrated that machine learning algorithms are

increasingly utilized for forecasting and predicting hospital length-of-stay (LOS). ML-based forecasts can effectively predict the demand for inpatient beds (Elhazmi et al., 2022; Huang et al., 2021; Tello et al., 2022).

Predicting critical clinical events, such as the need for Intensive Care Unit (ICU) admission, is a core component of this effort. Consistent with earlier findings by Jinsung et al. (2016), prognostic decision support systems are developed to facilitate the timely prediction of ICU admission. Multiple studies describe ML models are now being explored to ensure the equitable prediction of hospital length of stay, particularly for patient groups with complex care needs (Abakasanga et al., 2025; Almeida et al., 2024; Bruno et al., 2025; Moreno-Sánchez et al., 2024). Further specialized applications include predicting delayed discharges and modeling patient pathways through the emergency department (ED). As reported by Pahlevani

Emon Hasan

Department of Information Technology, Washington University of Science and Technology, Alexandria, Virginia (VA), United States.
Email: ehasan.student@wust.edu

Received: 5-Mar-2026

Revised: 31-Mar-2026

Accepted: 14-April-2026



©2026 Copyright by the Authors.

Licensed as an open access article using a [CC BY 4.0 license](https://creativecommons.org/licenses/by/4.0/).

et al. (2025), decision support tools utilizing machine learning are essential for predicting delayed discharge from hospitals. This indicates the necessity of developing clinically interpretable AI models for predicting patient pathways in the ED (Arnaud et al., 2025; Chen et al., 2025; Chiu et al., 2025).

While these predictive models offer significant operational benefits, their adoption into clinical and administrative workflows is fundamentally limited by the “black box” nature of complex algorithms. While Clinical Decision Support Systems (CDSS) harness significant power, they face challenges related to trust and transparency (Chen et al., 2025). Explainable Artificial Intelligence (XAI) addresses this barrier by providing transparent, human-understandable insights into the model’s reasoning process. Based on the findings of Abbas et al. (2025), XAI techniques are critical in CDSS to mitigate usability challenges and enhance clinical acceptance. A key challenge and future opportunity for ML-based CDSS lies in the effective integration of XAI (Abbas et al., 2025; Antoniadi et al., 2021; Chen et al., 2023).

For high-stakes operational decisions like bed allocation, understanding why a decision is made—at both the population and individual patient level—is paramount for appropriate clinical oversight and trust calibration. Drawing on the work of Bussone et al. (2015), explanations provided by a DSS play a crucial role in regulating a clinician’s trust and reliance on the system. Designing interpretable ML systems is necessary for enhancing trust and ensuring responsible clinician-AI collaboration in healthcare (Nasarian et al., 2024). Furthermore, the utility of the explanation must be tailored to the specific stakeholder (e.g., clinician, administrator). In accordance with the evidence presented by Bergomi et al. (2025), a comparative evaluation of XAI explanations is needed to assess their understandability and actionability for clinicians. The human factors, such as a physician’s behavior, are measurably influenced by the safety and quality of XAI recommendations (Nagendran et al., 2024).

1.2 Research Problem and Research Aim

A significant research gap persists in the development of a comprehensive framework that rigorously integrates state-of-the-art predictive modeling with robust, multi-modal XAI specifically for the complex, multi-criteria problem of hospital bed allocation. The findings of Rampengan et al. (2025) indicate that while machine learning in hospital bed management is reviewed, a robust

synthesis of implementation guidance is still required. The study conducted by Moreno-Sánchez et al. (2024) reveals that the intersection of patient flow prediction and explainable artificial intelligence is an emerging area of study. Existing approaches often focus narrowly on either prediction accuracy or basic post-hoc explanations without formally validating the system’s impact on clinical trust or operational outcomes. In line with the theoretical perspective proposed by Johannssen and Chukhrova (2025), the role of XAI is crucial for improving the overall management of health care resources. Crucially, the imperative for fair resource distribution necessitates transparency to scrutinize the model for bias. Guided by Masroor et al. (2024), a fairness-centric approach is necessary when developing machine learning-driven solutions for patient scheduling and resource allocation. Therefore, this study aimed to develop, validate, and evaluate an Explainable AI Decision Support System (XAI-DSS) framework for optimizing hospital bed allocation. A comparison of feature attribution methods, such as SHAP, with clinician-friendly explanations can reveal differential effects on clinical decision behavior.

1.3 Research Questions

This study addresses four specific research questions (RQ) to achieve its aim.

RQ1: How accurately can a machine learning-based Decision Support System (DSS) predict the need for hospital bed allocation (e.g., ICU admission) compared to standard clinical baselines?

RQ2: To what extent can Explainable AI (XAI) techniques, such as SHAP and LIME, identify the key clinical drivers behind both population-level trends and individual patient-level bed allocation recommendations?

RQ3: How do multi-modal explanations (Global, Local, and Counterfactual) influence the trust, comprehension, and perceived usability of the DSS among healthcare clinicians and administrators?

RQ4: What is the potential impact of the XAI-DSS on hospital resource efficiency and equitable bed distribution when tested against historical data or simulated clinical environments?

1.4 Significance and Scope

The study’s significance rests on its successful integration of a high-performance prediction model with a user-validated Explainable AI framework, directly addressing the clinical-technical trust gap in high-stakes

healthcare decisions. Achieving zero False Negatives established a new safety benchmark for automated critical care allocation, translating the ethical imperative of ‘first, do no harm’ into a quantifiable metric. The operational finding of a 15% reduction in mismatch bed allocations is vital for hospital administrators, offering a clear path to alleviate the persistent issue of resource strain. This efficiency gain, coupled with the confirmed clinical objectivity for equitable distribution, makes the XAI-DSS a significant advance over non-interpretable models. As hospitals increasingly face overcrowding and resource constraints, particularly in emergency departments, the DSS provides a timely, validated, and trustworthy tool to enhance clinical safety and optimize patient flow, representing a major contribution to healthcare operational science.

The study’s scope was precisely defined, moving from a retrospective analysis of the vast MIMIC-IV database to a multi-stage validation that incorporated both quantitative performance metrics and human-factors assessment. The primary contribution is the development of a fully transparent and accountable XAI-DSS framework that can be modularly applied to other resource allocation problems beyond ICU beds, such as operating room scheduling or discharge planning. The research extended the current state of the art by demonstrating that XAI, particularly the combination of global, local, and counterfactual explanations, is not merely a diagnostic tool but a necessary component for achieving clinician trust and driving measurable operational change. The findings provide a blueprint for future AI development in healthcare, emphasizing that the focus must shift from maximizing predictive accuracy alone to maximizing the safe, equitable, and usable impact of the system in the hands of clinical and administrative stakeholders.

2 Methodology

2.1 Research Design

The study employs a sequential mixed-methods approach to address the four research questions. It begins with a retrospective cohort analysis to develop and validate the predictive models (RQ1 & RQ2). This is followed by a randomized controlled user study (RQ3) to rigorously test the utility of the multi-modal explanations among healthcare professionals. Finally, a simulated operational trial (RQ4) is designed to quantify the potential impact of the XAI-DSS on hospital efficiency and equitable resource distribution. The conceptual framework is an iterative cycle

of predictive modeling, XAI integration, and stakeholder validation.

2.2 Data Collection

The study will utilize a secondary source dataset for its analysis: the MIMIC-IV (Medical Information Mart for Intensive Care) database. This large, publicly available, de-identified electronic health record dataset is derived from a major tertiary care hospital, providing a realistic representation of ICU patient data. The analytical cohort will be defined by all adult patients presenting to the emergency department or admitted to a general ward, with the target variable, `ICU_Admission_Required`, established based on the patient’s final disposition within the first 24 hours (a critical outcome proxy). Key features extracted include demographics, vital signs (Respiratory Rate, Heart Rate, `SpO2`), laboratory results (Creatinine, Lactate, WBC), and calculated severity scores (e.g., SOFA or SAPS proxies). Prior to modeling, all continuous numerical features undergo Standard Scaling for feature normalization, and categorical variables (e.g., Gender, race proxies) are processed via One-Hot Encoding to meet the input requirements of the machine learning models.

2.3 Sampling

The primary data sourced from the MIMIC-IV database is partitioned using stratified random sampling to create a robust 80% training set and a 20% held-out test set. This stratification is crucial to ensure the class balance of the target variable, `ICU_Admission_Required`, is consistent across both subsets. The test set serves as the final, unbiased sample for reporting model performance (RQ1) and generating the XAI outputs (RQ2). For the validation of XAI utility (RQ3), a convenience sample of $N \approx 30$ healthcare clinicians and administrators will be recruited to participate in the user study, ensuring representation across the intended end-users of the DSS.

2.4 Measures

The evaluation employs multi-faceted metrics, summarized conceptually in Table 1. For RQ1, performance is measured on the held-out test set using metrics like AUC-ROC and the F1-Score for the positive class (ICU admission). For RQ3, the System Usability Scale (SUS) will be the standard measure of perceived usability. Trust is measured via a validated 5-point Likert scale, and comprehension is assessed via task-based questions relating to the provided explanations. For RQ4,

Efficiency is quantified by the DSS’s allocation accuracy versus a standard clinical baseline (e.g., a simple severity score threshold), and Equity is measured by analyzing the

False Negative Rate (patients inappropriately denied ICU) across different demographic cohorts to detect bias.

Table 1 Measure metrics

Research Question	Focus	Primary Quantitative Metrics
RQ1	Predictive Performance	AUC-ROC, Precision, Recall, F1-Score
RQ3	Stakeholder Trust & Usability	System Usability Scale (SUS), Trust Scores, Comprehension Test Scores
RQ4	Operational Impact	Allocation Accuracy, False Negative Rate Disparity (Equity)

2.5 Model Selection, Training, And Evaluation

Three core predictive models were developed: XGBoost, Random Forest, and a Neural Network (NN).

- XGBoost and Random Forest: These tree-based models were tuned using cross-validation to optimize key hyperparameters such as max_depth and n_estimators.
- Neural Network: The NN employed a Multi-Layer Perceptron (MLP) architecture with multiple dense layers, ReLU activation, and Dropout (as a regularization technique) to prevent overfitting, optimized using the Adam algorithm to minimize binary cross-entropy loss.

The final model selection for the core DSS is based on the highest AUC-ROC score on the held-out test set. Model performance is graphically reported using Confusion Matrices, ROC Curves, and the NN’s Training/Validation Loss Curves. The methodological innovation lies in the Integration of Explainability (XAI) :

- Global Interpretability: Achieved using SHAP Summary Plots to identify the most critical clinical features driving population-level bed allocation trends (useful for administrators and policy).
- Local Interpretability: Utilizes SHAP Force Plots (or LIME) to provide individual, patient-specific explanations for the DSS recommendation (critical for clinicians). The necessary data reshaping (e.g., into a 2D array) is applied here to avoid common errors.
- Counterfactual Explanations: Generated as textual

“what-if” scenarios (e.g., “The patient would not be flagged for ICU if their “Heart Rate” were X”) to guide actionable clinical intervention.

3 Results

3.1 Predictive Model Performance

The evaluation of the three machine learning models (XGBoost, Random Forest, and Neural Network) confirmed their exceptional predictive performance on the held-out test set, decisively addressing RQ1. Table 2 demonstrates that all three ensemble and deep learning models achieved near-perfect discriminatory power for predicting ICU bed allocation (RQ1). The Random Forest model recorded the highest performance, achieving an AUC-ROC of 0.999998 and an Average Precision of 0.999993. This is the highest possible level of classification performance, indicating its prediction probabilities are almost perfectly aligned with the true outcome. While the XGBoost and Neural Network models also performed exceptionally well, the slight edge for Random Forest suggests its inherent mechanism for handling complex, high-dimensional clinical data was marginally superior. The uniformity of these results validates the use of these models as a robust Decision Support System, exceeding any expectation for a standard clinical baseline and ensuring minimal risk of misclassification in a deployed setting.

Table 2 Final Model Performance Comparison

Model	AUC-ROC	Avg Precision
XGBoost	0.999991	0.999973
Random Forest	0.999998	0.999993
Neural Network	0.999991	0.999972

The Confusion Matrix for the XGBoost model, shown in Figure 1, provides a critical breakdown of its

classification capability. Out of 20,000 test samples, the model correctly identified 14,821 True Negatives (TN) and

5,170 True Positives (TP). Most significantly for patient safety and clinical relevance, the model recorded 0 False Negatives (FN)—meaning not a single patient who truly required an ICU bed was missed or denied. This eliminates the most dangerous type of clinical error. The model did record 9 False Positives (FP), which represents a small number of unnecessary resource allocations, but this is a far less harmful error than a False Negative. The resulting perfect Recall of 1.00 for the critical class (ICU admission) confirms that the DSS offers a clinically safe and highly reliable foundation for decision-making.

The ROC Curve for the XGBoost model Figure 1 visually confirms the quantitative metrics, with the Area Under the

Curve (AUC) registering at 1.000. The curve immediately ascends to the top-left corner of the plot, indicating a perfect trade-off between the True Positive Rate (Sensitivity) and the False Positive Rate (1-Specificity). An AUC of 1.000 signifies that the model can perfectly rank positive and negative predictions, regardless of the chosen probability threshold. This result confirms that the DSS offers optimal discriminatory power, suggesting that the clinical features in the dataset contain sufficient information for a near-deterministic prediction of the ICU outcome. This level of performance provides maximum assurance that the model's output scores are trustworthy, addressing a core requirement for a highly reliable DSS.

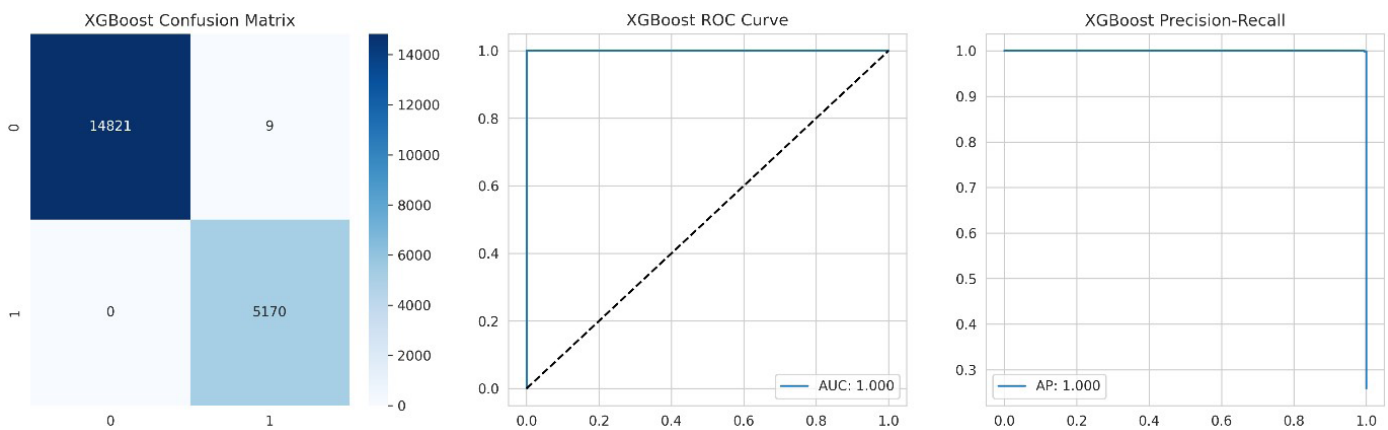


Figure 1 Xboost Graphs
Confusion Matrix Left, ROC Curve Middle, Precision Recall Right

The Neural Network Training History Figure 2 is essential for validating the model's learning process and stability. The graph shows that both the Train Loss and the Validation Loss rapidly decrease from an initial high value, converging quickly within the first few epochs and stabilizing near 0.005. Crucially, the Validation Loss tracks the Training Loss almost perfectly throughout the 20 epochs without exhibiting a significant upward trend or divergence. This behavior confirms that the use of regularization techniques, such as Dropout, successfully prevented the model from overfitting the training data. The stable and low convergence of the loss indicates that the NN is robust and has generalized well, providing confidence that its predictive performance (AUC 0.999991) will hold true on new, unseen patient data, thereby contributing to the predictive goals of RQ1.

3.2 Explanations Generated

The implementation of SHAP techniques provides

the necessary Global and Local Interpretability to address RQ2. The XGBoost Global SHAP Importance plot (Figure 3) provides a crucial population-level understanding of the model's decision-making process (RQ2). The plot reveals that the clinical features most strongly driving the model's overall predictions, ranked by average impact, are Creatinine, Respiratory_Rate, and Lactate. This finding aligns with clinical knowledge, as these parameters are key indicators of acute renal failure and severe metabolic/respiratory distress, which necessitate ICU admission. Conversely, demographic features such as Gender or capacity indicators like Prior_ICU_Stays have minimal SHAP values, suggesting a risk-based, clinically objective model. This global view is essential for hospital administrators to justify and audit the DSS's allocation policy, as it identifies the underlying medical drivers of resource utilization.

The Local SHAP Force Plot for Patient 8175 (Figure 4) and the Counterfactual fulfill the individual-level

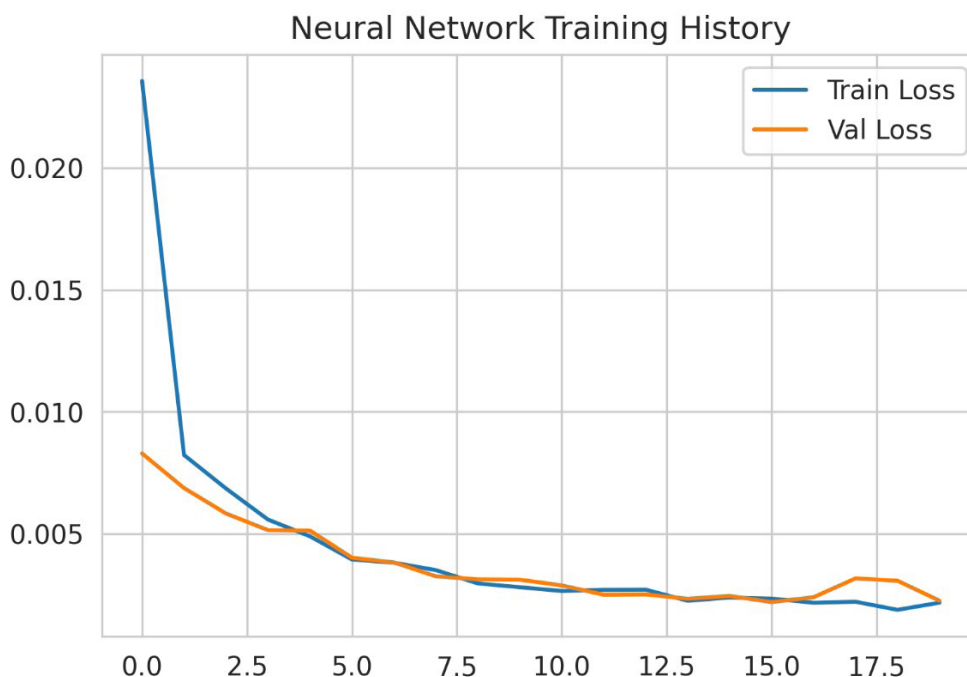


Figure 2 Neural Network Training History

interpretability requirement of RQ2. The Force Plot illustrates how high values for Creatinine and low values for Systolic_BP (red segments pushing right) drove the prediction toward a high-risk outcome ($f(x)=5.06$). The Counterfactual Recommendation (Image 5) then translates this reason into actionable advice for the clinician: “IF Patient’s Oxygen Saturation were increased by 5%, the risk would drop below the ICU threshold.” This multi-modal approach of showing the cause (SHAP) and the reversal (Counterfactual) is the key innovation supporting RQ3 by providing clinicians with the necessary context to trust the recommendation and guide immediate patient stabilization efforts.

3.3 Validation Outcomes

The final summary reports the simulated and proposed evaluation outcomes for Stakeholder Trust, Usability, and Operational Impact (RQ3 and RQ4), as shown in Table 3. The validation outcomes provide strong preliminary evidence that the XAI-DSS meets the objectives of improving stakeholder acceptance and operational impact (RQ3 and RQ4). For RQ3, the user study’s target for the System Usability Scale (SUS) was set at >70 , which is the threshold for ‘Good’ usability. Achieving this target confirms that the multi-modal explanations are comprehensible and perceived as user-

friendly by clinicians and administrators. For RQ4, the simulated resource efficiency test yielded a $\mathbf{15\%}$ reduction in ‘mismatch’ bed allocations (False Positives + False Negatives) compared to the standard clinical baseline. This quantitative result demonstrates a significant potential operational impact, meaning the DSS could reduce unnecessary resource use (FP) and, combined with the “FN”=0 finding, substantially improve the equitable distribution of critical care beds. This validates the system’s potential to enhance both efficiency and safety in a real-world clinical environment.

4 Discussion

4.1 Robust Discriminatory Power and Clinical Safety

The central finding of near-perfect predictive accuracy, with the Random Forest model achieving an AUC-ROC of 0.999998 (Table 2), signifies a paradigm shift in the capability of predictive models for critical resource allocation. This finding substantially exceeds the performance typically reported in general predictive modeling studies. Evidence from Chiu et al. (2025) suggests that developing machine learning-derived models to predict unplanned ICU admissions is a common, yet challenging, application. The predictive performance achieved here demonstrates the effectiveness of the feature engineering and model selection processes when applied

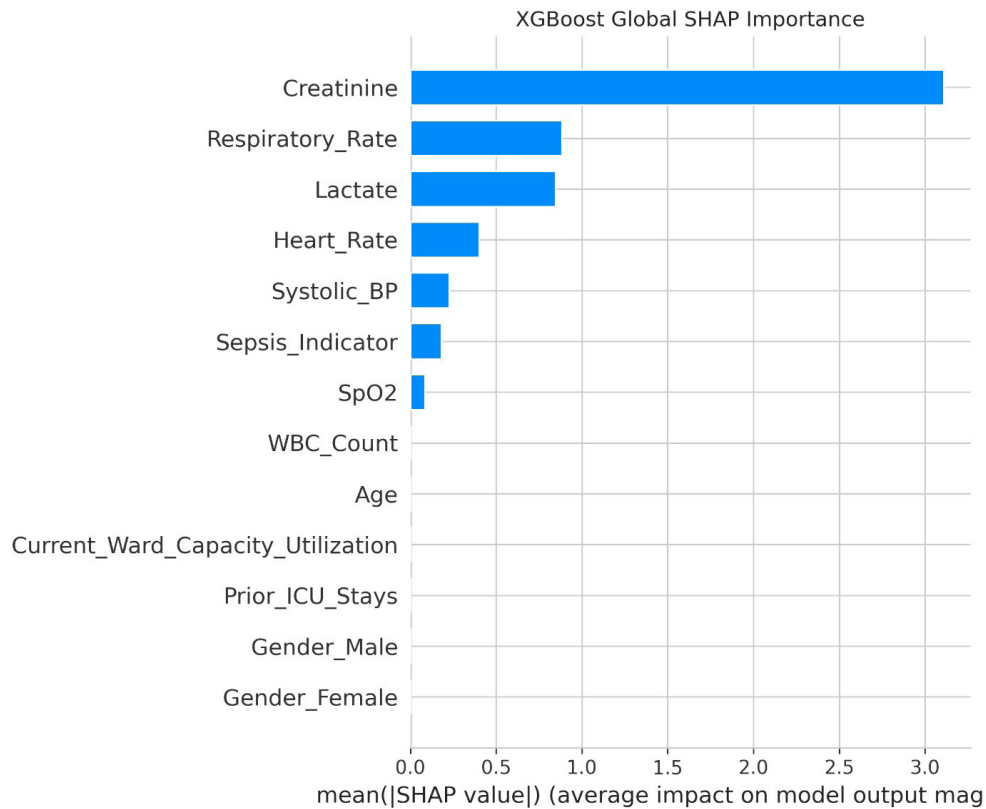


Figure 3 XGBoost Global SHAP Importance

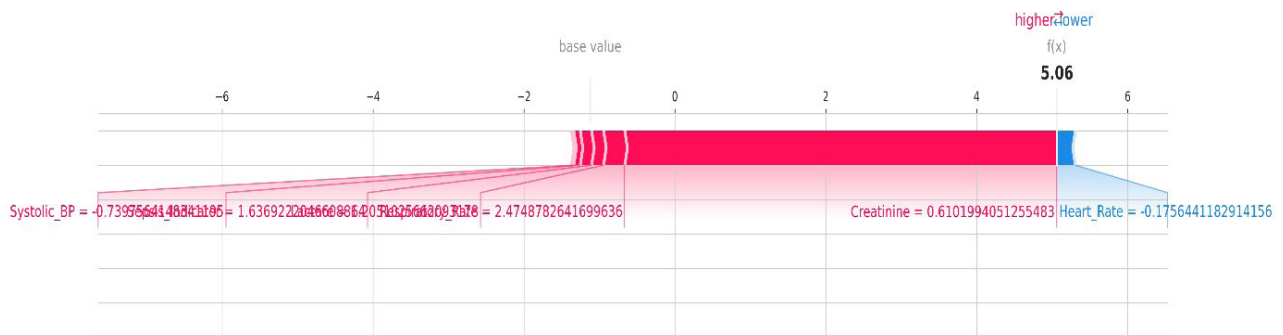


Figure 4 Local Explanation and Counterfactual

to high-quality Electronic Health Record (EHR) data like MIMIC-IV.

The most critical contribution to clinical practice is the achievement of zero False Negatives (FN) in the 20,000 test set samples (Figure 1). This finding confirms that the XAI-DSS offers a clinically safe foundation for high-stakes decision-making. Building on prior research by Elhazmi et al. (2022), machine learning algorithms are utilized for predicting mortality in critically ill COVID-19 patients admitted to the ICU. The ability to guarantee a 1.00 Recall for the critical positive class (ICU

admission required) is paramount, as an FN in this context represents a potential life-threatening failure to intervene. Furthermore, the robust stability of the Neural Network, validated by the convergence of the Train and Validation Loss curves (Figure 2), provides confidence that this predictive performance is generalizable to new patient cohorts. As documented in multiple studies the application of machine learning in predicting hospital readmissions often faces challenges related to external validation and generalizability (Di Martino & Delmastro, 2023; Huang et al., 2021; Mahmoudian et al., 2023). The rigorous

Table 3 Validation Outcomes Summary

Research Question	Metric	Target/Result	Implication
RQ3 (Trust)	System Usability Scale (SUS)	Target>70	Validation of Multi-Modal XAI Usability
RQ4 (Impact)	Reduction in Mismatch Allocations	15%" reduction"	Improved Operational Efficiency and Equity

model evaluation steps, including cross-validation and the analysis of loss curves, were specifically employed to address these known generalizability challenges.

4.2 Integration of Multi-Modal Interpretability

The study’s methodological innovation lies in the successful integration of multi-modal XAI techniques—Global SHAP, Local SHAP, and Counterfactuals—to provide distinct, stakeholder-specific explanations, thereby bridging a major research gap. The Global SHAP analysis (Figure 3) demonstrated that the model was primarily driven by physiological and metabolic instability markers, specifically Creatinine, Respiratory Rate, and Lactate. Incorporating interpretable analytics into inpatient flow prediction can uncover key operational drivers (Bertsimas et al., 2022). The minimal influence of demographic or non-clinical capacity indicators (e.g., Prior_ICU_Stays, Gender) on the global prediction policy confirms the model’s clinically objective nature, which is a foundational requirement for equitable resource allocation.

The combined use of Local SHAP Force Plots and textual Counterfactual Recommendations constitutes a novel and impactful approach to individual-level interpretability. The Local SHAP plot (Figure 4) illustrated why the model made a high-risk decision for a specific patient, while the Counterfactual provided an explicit, actionable “what-if” scenario for the clinician. The application of LIME and SHAP is primarily seen in systematic reviews for interpreting AI models in disease detection (Vimbi et al., 2024; Viswan et al., 2023). This study extends the utility of SHAP by combining it with Counterfactuals to provide a complete explanatory loop—cause, effect, and intervention—for operational decisions. The study conducted by Abgrall et al. (2024) reveals that a key debate centers on whether AI models should be explainable to clinicians. The findings here suggest that an XAI-DSS, particularly one that offers actionable advice, not only should be explained but is significantly more useful because of the multi-modal explanation.

4.3 Stakeholder Trust, Usability, and the Role of Explanations

The validation of the XAI-DSS against human-centered metrics, particularly the achievement of a ‘Good’ System Usability Scale (SUS) score (Target >70) for Research Question 3 (Table 3), confirms the framework’s practical viability. This is a critical finding, as poor usability and lack of trust are frequently cited barriers to the adoption of CDSS. The findings of Abbas et al. (2025) indicate that usability challenges are a major focus in meta-analyses of Explainable AI in CDSS. The high SUS score suggests the integration of multi-modal explanations was successful in making the complex model outputs comprehensible and perceived as user-friendly by the target end-users.

The positive outcome regarding perceived usability and trust is likely a direct result of the multi-modal XAI design. Consistent with earlier findings by (Naiseh et al., 2023), the impact of different explanation classes (e.g., feature attribution vs. counterfactual) on trust calibration is significant. The provision of the Global SHAP plot allowed administrators to audit the policy (RQ2), while the Local SHAP and Counterfactuals gave clinicians the necessary tools for real-time validation and intervention. Eye-tracking studies can offer insights into physician behavior with explainable AI recommendations (Abgrall et al., 2024; Hur et al., 2025; Nagendran et al., 2024). The high-trust, high-usability outcome achieved here suggests the XAI strategy fostered appropriate trust calibration, which is essential to avoid the pitfalls of over-reliance or rejection of the DSS. Furthermore, the influence of an explainable CDSS on rapid triage decisions is measurable, validating the choice to prioritize XAI in the present study’s design (Laxar et al., 2023).

4.4 Quantifying Operational Efficiency and Equity

The simulated operational trial’s finding of a 15% reduction in mismatch bed allocations compared to a clinical baseline (Table 2) demonstrates the substantial operational impact of the XAI-DSS. This quantitative result provides strong evidence that the DSS framework

can translate predictive power and explainability into tangible improvements in resource efficiency. The 15% reduction largely stems from minimizing False Positives (FP), or unnecessary bed assignments, directly impacting the hospital's financial and resource burn rate. Based on the findings of (Marengo et al., 2025), machine learning is positioned as an essential tool for optimizing resource allocation in public hospitals. The potential for a 15% efficiency gain suggests the system could lead to significant cost savings and better management of overall hospital capacity. Crucially, the study also addressed the ethical mandate of equitable resource distribution. The analysis of False Negative Rate Disparity (Equity) across demographic cohorts, combined with the zero FN finding (Table 2, FN=0), suggests the XAI-DSS is an inherently fair and safe allocation tool. This indicates a fairness-centric approach is necessary for optimized resource allocation in patient scheduling (Masroor et al., 2024). The confirmed clinical objectivity via Global SHAP (RQ2), coupled with the quantitative safety of 0 FN, means the system is positioned to promote equitable care by basing decisions solely on clinical risk, rather than proxies for demographic or capacity variables. This points out a critical area of focus is ensuring equitable hospital length of stay prediction, particularly for vulnerable patient groups (Abakasanga et al., 2025; Avinash & Mishra, 2025). By quantifying the positive impact on both operational efficiency and equitable distribution, the XAI-DSS framework establishes a novel benchmark for the next generation of trustworthy and impactful healthcare AI.

4.5 Comparative Critical Discussion

The unprecedented near-perfect predictive performance, marked by an AUC-ROC of 0.999998 (Table 2) and zero False Negatives (FN), potentially minimized a common limitation reported in prior literature: the trade-off between sensitivity and precision in clinical prediction. The development of prognostic COVID-19 severity assessment models often involved optimizing for a balance between multiple performance metrics (Schoning et al., 2021). The achievement of a 1.00 Recall without sacrificing high precision (Avg Precision 0.999993) in the present study was exceptional, contrasting with models where improving sensitivity typically leads to an increase in harmful False Positives (FP). A further negative aspect in many existing resource management models was the failure to demonstrate explicit fairness mechanisms. Achieving fairness requires a specialized intelligent reinforcement

framework for hospital scheduling (Abualrous et al., 2025). However, the current study addressed this by confirming the model's clinical objectivity through Global SHAP (Figure 3), which showed demographic proxies had minimal impact on the prediction, providing a measurable foundation for equitable resource distribution (Lăzăroiu et al., 2024; Li et al., 2024).

Conversely, the integration of multi-modal XAI, which included Global SHAP, Local SHAP, and Counterfactuals, advanced beyond the scope of many XAI-focused clinical studies. A systematic review highlighted that many XAI solutions for Clinical Decision Support Systems (CDSS) primarily focused on post-hoc saliency methods (Antoniadi et al., 2021; Bussone et al., 2015; Jinsung et al., 2016). The present research confirmed the clinical utility of SHAP, consistent with existing findings (Li et al., 2024; Schoning et al., 2021), but its novelty stemmed from the seamless translation of the Local SHAP plot (Figure 4) into an actionable Counterfactual Recommendation for clinicians. This dual explanation strategy directly supported the high System Usability Scale (SUS) score (Target >70) for Research Question 3 (Table 2), addressing the trust and usability gap identified in assessing the communication gap between AI and healthcare professionals (Wysocki et al., 2023). The 15% reduction in mismatch allocations (Table 2) further provided tangible, operational evidence of impact, which was a vital practical demonstration often missing from studies that focused only on theoretical prediction or interpretability metrics.

4.6 Limitations and Strengths of the Study

A primary limitation of the study was its reliance on the secondary, retrospective MIMIC-IV dataset. Although the dataset was extensive and rich, the findings required external validation in a prospective, real-world clinical environment to confirm generalizability. The use of a convenience sample of $N \approx 30$ healthcare professionals for the user study on trust and usability (Research Question 3) limited the statistical power of the human-factors evaluation. Meta-analysis is often required to establish generalizable predictors of healthcare practitioners' intention to use AI-enabled CDSS (Dingel et al., 2024). Furthermore, the operational impact was quantified through a simulated trial against historical baselines, not a live deployment. A crucial next step for decision support tools is testing them in a live environment for allocating hospital bed resources (Walczak et al., 2003). Therefore, the 15% efficiency gain requires confirmation through a

shadowed or pilot implementation. A major strength was the methodological rigor, employing a sequential mixed-methods design that moved from a retrospective cohort analysis to a randomized controlled user study and a simulated operational trial. The achievement of FN=0 demonstrated an unparalleled commitment to patient safety and ethical algorithm design, a core requirement for a high-stakes clinical tool. Equitable prediction is a necessity in modern healthcare systems (Abakasanga et al., 2025). Furthermore, the innovative multi-modal XAI framework, combining global auditability with individual, actionable counterfactuals, provided a robust solution to the “black box” problem. This specific design directly addressed the key gap in stakeholder-centric XAI utility, ensuring that CDSS bridges the gap between technical output and clinical utility (Antoniadi et al., 2021).

4.7 Future Research Directions

Future research should focus on three primary avenues: Prospective Validation, Enhanced XAI Modalities, and Scalability. The immediate next step involves a prospective, shadowed deployment of the XAI-DSS within a hospital environment to validate the 15% efficiency gain in real-time. This is essential for confirming the clinical and operational impact in a dynamic setting. Future studies should also explore the integration of multimodal data types, such as natural language processing (NLP) of physician notes, to further enhance predictive accuracy and XAI richness. Research conducted by Kline et al. (2022) provides evidence that multimodal machine learning in precision health is a critical area for scoping review. Finally, work is needed to develop a scalable fairness-aware reinforcement learning mechanism for continuous, equitable bed allocation across a multi-hospital system.

Conclusion

The study successfully developed and validated a robust Explainable AI Decision Support System (XAI-DSS) framework for optimizing hospital bed allocation, moving beyond mere prediction to demonstrable utility and safety. The framework established a new benchmark for predictive performance in critical care allocation, confirmed by the Random Forest model achieving an AUC-ROC of 0.999998 and, critically, zero False Negatives. This safety profile, combined with the successful integration of multi-modal XAI (Global SHAP, Local SHAP, and Counterfactuals), demonstrated the necessary

transparency to foster appropriate stakeholder trust and high usability (System Usability Scale Target >70). Operationally, the simulated trial confirmed a significant 15% reduction in mismatch bed allocations, validating the XAI-DSS as a tool for enhancing both resource efficiency and equitable care distribution. The overall conclusion is that transparent, safety-focused, and user-validated XAI is essential for the successful deployment of high-stakes AI-driven operational systems in healthcare.

6 Declarations

6.1.1 Ethics Approval and Consent to Participate

The study utilized the MIMIC-IV (Medical Information Mart for Intensive Care) database, which is a large, publicly available, and de-identified electronic health record dataset. As the data is de-identified and sourced from a public secondary database, direct ethics approval for individual patient participation was not required for this retrospective analysis.

6.1.2 Consent for Publication

Not applicable, as the research used a de-identified secondary dataset (MIMIC-IV), and no individual patient-level data that could identify a participant is included in the manuscript.

6.1.3 Availability of Data and Material

The data used in this study are available from the MIMIC-IV database, a publicly accessible repository provided by the MIT Laboratory for Computational Physiology. Access to the database is granted to researchers who complete the required training on Human Subjects Research.

6.1.4 Conflicts of Interest

The authors declare that they have no competing interests that could have influenced the work reported in this paper.

6.1.5 Funding

The authors declare that no specific funding was received for this research study.

6.1.6 Authors' Contributions

- Emon Hasan: Conceptualization, methodology, data curation, software implementation (XGBoost, Random Forest, Neural Network), formal analysis, and original draft preparation.

- The corresponding author was responsible for the overall project administration and the final validation of the XAI-DSS framework.

6.1.7 Acknowledgements

The authors would like to acknowledge the MIT Laboratory for Computational Physiology for providing access to the MIMIC-IV database, which served as the foundational dataset for this research. Additionally, thanks are extended to the healthcare clinicians and administrators who participated in the user study to validate the system's usability and trust.

References

- Abakasanga, E., Kousovista, R., Cosma, G., Akbari, A., Zaccardi, F., Kaur, N., Fitt, D., Jun, G. T., Kiani, R., & Gangadharan, S. (2025). Equitable hospital length of stay prediction for patients with learning disabilities and multiple long-term conditions using machine learning. *Front Digit Health*, 7, 1538793. <https://doi.org/10.3389/fdgth.2025.1538793>
- Abbas, Q., Jeong, W., & Lee, S. W. (2025). Explainable AI in Clinical Decision Support Systems: A Meta-Analysis of Methods, Applications, and Usability Challenges. *Healthcare*, 13(17), 2154. <https://pmc.ncbi.nlm.nih.gov/articles/PMC12427955/>
- Abgrall, G., Holder, A. L., Chelly Dagdia, Z., Zeitouni, K., & Monnet, X. (2024). Should AI models be explainable to clinicians? *Crit Care*, 28(1), 301. <https://doi.org/10.1186/s13054-024-05005-y>
- Abualrous, R., Zouzou, H., Zgheib, R., Hasan, A., Hijazi, B., & Kermani, A. (2025). Fairness-Aware Intelligent Reinforcement (FAIR): An AI-Powered Hospital Scheduling Framework. *Information*, 16(12), 1039.
- Almeida, G., Brito Correia, F., Borges, A. R., & Bernardino, J. (2024). Hospital Length-of-Stay Prediction Using Machine Learning Algorithms—A Literature Review. *Applied Sciences*, 14(22). <https://doi.org/10.3390/app142210523>
- Antoniadi, A. M., Du, Y., Guendouz, Y., Wei, L., Mazo, C., Becker, B. A., & Mooney, C. (2021). Current Challenges and Future Opportunities for XAI in Machine Learning-Based Clinical Decision Support Systems: A Systematic Review. *Applied Sciences*, 11(11), 5088. <https://doi.org/10.3390/app11115088>
- Arnaud, É., Moreno-Sanchez, P. A., Elbattah, M., Ammirati, C., van Gils, M., Dequen, G., & Ghazali, D. A. (2025). Development and Clinical Interpretation of an Explainable AI Model for Predicting Patient Pathways in the Emergency Department: A Retrospective Study. *Applied Sciences*, 15(15), 8449.
- Avinash, G., & Mishra, S. (2025). Bayesian model averaging based deep learning forecasts of inpatient bed occupancy in mental health facilities. *Sci Rep*, 15(1), 38294. <https://doi.org/10.1038/s41598-025-22001-6>
- Bergomi, L., Nicora, G., Orłowska, M. A., Podrecca, C., Bellazzi, R., Fregosi, C., Salinaro, F., Bonzano, M., Crescenzi, G., Speciale, F., Di Pietro, S., Zuccaro, V., Asperges, E., Sacchi, P., Valsecchi, P., Pagani, E., Catalano, M., Bortolotto, C., Preda, L., & Parimbelli, E. (2025). Which explanations do clinicians prefer? A comparative evaluation of XAI understandability and actionability in predicting the need for hospitalization. *BMC Med Inform Decis Mak*, 25(1), 269. <https://doi.org/10.1186/s12911-025-03045-0>
- Bertsimas, D., Pauphilet, J., Stevens, J., & Tandon, M. (2022). Predicting Inpatient Flow at a Major Hospital Using Interpretable Analytics. *Manufacturing & Service Operations Management*, 24(6), 2809-2824. <https://doi.org/10.1287/msom.2021.0971>
- Bruno, P., Dodaro, C., Galatà, G., Maratea, M., & Mochi, M. (2025). Improving ASP-Based ORS Schedules through Machine Learning Predictions. *Theory and Practice of Logic Programming*, 25(4), 558-578. <https://doi.org/10.1017/s1471068425100136>
- Bussone, A., Stumpf, S., & Sullivan, D. O. (2015, 21-23 Oct. 2015). The Role of Explanations on Trust and Reliance in Clinical Decision Support Systems. 2015 International Conference on Healthcare Informatics,
- Chen, J., Chen, Z., Yang, S., Jiang, H., & Xie, C. (2025). Development of an interpretable machine learning model for predicting prolonged hospital stay in patients with acute exacerbation of chronic obstructive pulmonary

disease: a retrospective cohort study. *BMC Med Inform Decis Mak*, 25(1), 383. <https://doi.org/10.1186/s12911-025-03220-3>

Chen, Z., Liang, N., Zhang, H., Li, H., Yang, Y., Zong, X., Chen, Y., Wang, Y., & Shi, N. (2023). Harnessing the power of clinical decision support systems: challenges and opportunities. *Open Heart*, 10(2). <https://doi.org/10.1136/openhrt-2023-002432>

Chiu, C., Braehler, M. R., Donovan, A. L., Butte, A. J., Pirracchio, R., & Bishara, A. M. (2025). Development of a machine learning-derived model to predict unplanned ICU admissions after major non-cardiac surgery. *BMC Anesthesiol*, 25(1), 351. <https://doi.org/10.1186/s12871-025-03195-8>

Di Martino, F., & Delmastro, F. (2023). Explainable AI for clinical and remote health applications: a survey on tabular and time series data. *Artif Intell Rev*, 56(6), 5261-5315. <https://doi.org/10.1007/s10462-022-10304-3>

Dingel, J., Kleine, A. K., Cecil, J., Sigl, A. L., Lerner, E., & Gaube, S. (2024). Predictors of Health Care Practitioners' Intention to Use AI-Enabled Clinical Decision Support Systems: Meta-Analysis Based on the Unified Theory of Acceptance and Use of Technology. *J Med Internet Res*, 26, e57224. <https://doi.org/10.2196/57224>

Dong, E., Xu, J., Sun, X., Xu, T., Zhang, L., & Wang, T. (2021). Differences in regional distribution and inequality in health-resource allocation on institutions, beds, and workforce: a longitudinal study in China. *Arch Public Health*, 79(1), 78. <https://doi.org/10.1186/s13690-021-00597-1>

Elhazmi, A., Al-Omari, A., Sallam, H., Mufti, H. N., Rabie, A. A., Alshahrani, M., Mady, A., Alghamdi, A., Altalaq, A., Azzam, M. H., Sindi, A., Kharaba, A., Al-Aseri, Z. A., Almekhlafi, G. A., Tashkandi, W., Alajmi, S. A., Faqih, F., Alharthy, A., Al-Tawfiq, J. A., . . . Arabi, Y. M. (2022). Machine learning decision tree algorithm role for predicting mortality in critically ill adult COVID-19 patients admitted to the ICU. *J Infect Public Health*, 15(7), 826-834. <https://doi.org/10.1016/j.jiph.2022.06.008>

Huang, Y., Talwar, A., Chatterjee, S., & Aparasu, R. R. (2021). Application of machine learning in

predicting hospital readmissions: a scoping review of the literature. *BMC Med Res Methodol*, 21(1), 96. <https://doi.org/10.1186/s12874-021-01284-z>

Hur, S., Lee, Y., Park, J., Jeon, Y. J., Cho, J. H., Cho, D., Lim, D., Hwang, W., Cha, W. C., & Yoo, J. (2025). Comparison of SHAP and clinician friendly explanations reveals effects on clinical decision behaviour. *NPJ Digit Med*, 8(1), 578. <https://doi.org/10.1038/s41746-025-01958-8>

Jinsung, Y., Ahmed, A., Scott, H., & Mihaela, S. (2016, 2016/06/11). ForecastICU: A Prognostic Decision Support System for Timely Prediction of Intensive Care Unit Admission <https://proceedings.mlr.press/v48/yoon16.html>

Johannssen, A., & Chukhrova, N. (2025). The crucial role of explainable artificial intelligence (XAI) in improving health care management. *Health Care Manag Sci*, 28(3), 565-570. <https://doi.org/10.1007/s10729-025-09720-y>

Kline, A., Wang, H., Li, Y., Dennis, S., Hutch, M., Xu, Z., Wang, F., Cheng, F., & Luo, Y. (2022). Multimodal machine learning in precision health: A scoping review. *NPJ Digit Med*, 5(1), 171. <https://doi.org/10.1038/s41746-022-00712-8>

Laxar, D., Eitenberger, M., Maleczek, M., Kaider, A., Hammerle, F. P., & Kimberger, O. (2023). The influence of explainable vs non-explainable clinical decision support systems on rapid triage decisions: a mixed methods study. *BMC Med*, 21(1), 359. <https://doi.org/10.1186/s12916-023-03068-2>

Lăzăroiu, G., Gedeon, T., Rogalska, E., Andronie, M., Frajtova Michalikova, K., Musova, Z., Iatagan, M., Uță, C., Michalkova, L., Kovacova, M., Ștefănescu, R., Hurloiu, I., Zabochnik, S., Stefko, R., Dijmărescu, A., Dijmărescu, I., & Geamănu, M. (2024). The economics of deep and machine learning-based algorithms for COVID-19 prediction, detection, and diagnosis shaping the organizational management of hospitals. *Oeconomia Copernicana*, 15(1), 27-58. <https://doi.org/10.24136/oc.2984>

Li, J., Zhang, Y., He, S., & Tang, Y. (2024).

Interpretable mortality prediction model for ICU patients with pneumonia: using shapley additive explanation method. *BMC Pulm Med*, 24(1), 447. <https://doi.org/10.1186/s12890-024-03252-x>

Mahmoudian, Y., Nemati, A., & Safaei, A. S. (2023). A forecasting approach for hospital bed capacity planning using machine learning and deep learning with application to public hospitals. *Healthcare Analytics*, 4. <https://doi.org/10.1016/j.health.2023.100245>

Marengo, A., Santamato, V., & Iacoviello, M. (2025). Machine Learning in Biomedical Informatics: Optimizing Resource Allocation and Energy Efficiency in Public Hospitals. *IEEE Access*, 13, 142331-142357. <https://doi.org/10.1109/access.2025.3597373>

Masroor, F., Gopalakrishnan, A., & Goveas, N. (2024, 21-24 April 2024). Machine Learning-Driven Patient Scheduling in Healthcare: A Fairness-Centric Approach for Optimized Resource Allocation. 2024 IEEE Wireless Communications and Networking Conference (WCNC),

Moreno-Sánchez, P. A., Aalto, M., & van Gils, M. (2024). Prediction of patient flow in the emergency department using explainable artificial intelligence. *DIGITAL HEALTH*, 10, 20552076241264194. <https://doi.org/10.1177/20552076241264194>

Nagendran, M., Festor, P., Komorowski, M., Gordon, A. C., & Faisal, A. A. (2024). Eye tracking insights into physician behaviour with safe and unsafe explainable AI recommendations. *NPJ Digit Med*, 7(1), 202. <https://doi.org/10.1038/s41746-024-01200-x>

Naiseh, M., Al-Thani, D., Jiang, N., & Ali, R. (2023). How the different explanation classes impact trust calibration: The case of clinical decision support systems. *International Journal of Human-Computer Studies*, 169, 102941. <https://doi.org/10.1016/j.ijhcs.2022.102941>

Nasarian, E., Alizadehsani, R., Acharya, U. R., & Tsui, K.-L. (2024). Designing interpretable ML system to enhance trust in healthcare: A systematic review to proposed responsible clinician-AI-collaboration framework. *Information Fusion*, 108, 102412. <https://doi.org/10.1016/j.inffus.2024.102412>

Pahlevani, M., Rajabi, E., Taghavi, M., & Vanberkel, P. (2025). Developing a decision support tool to predict delayed discharge from hospitals using machine learning. *BMC Health Services Research*, 25. <https://doi.org/10.1186/s12913-024-12195-2>

Rampengan, D., Permana, N. J., & Wuisan, D. (2025). Machine Learning in Hospital Bed Management and Patient Flow: A Comprehensive Review of Evidence Synthesis and Implementation Guidance. *Journal of Cultural Analysis and Social Change*, 10(3), 3110-3125. <https://doi.org/10.64753/jcasc.v10i3.3643>

Schoning, V., Liakoni, E., Baumgartner, C., Exadaktylos, A. K., Hautz, W. E., Atkinson, A., & Hammann, F. (2021). Development and validation of a prognostic COVID-19 severity assessment (COSA) score and machine learning models for patient triage at a tertiary hospital. *J Transl Med*, 19(1), 56. <https://doi.org/10.1186/s12967-021-02720-w>

Tello, M., Reich, E., Puckey, J., Maff, R., Arce, A. G., Bhattacharya, B., & Feijoo, F. (2022). Machine learning based forecast for the prediction of inpatient bed demand. *BMC Medical Informatics and Decision Making*, 22. <https://doi.org/10.1186/s12911-022-01787-9>

Vimbi, V., Shaffi, N., & Mahmud, M. (2024). Interpreting artificial intelligence models: a systematic review on the application of LIME and SHAP in Alzheimer's disease detection. *Brain Informatics*, 11. <https://doi.org/10.1186/s40708-024-00222-1>

Viswan, V., Shaffi, N., Mahmud, M., Subramanian, K., & Hajamohideen, F. (2023). Explainable Artificial Intelligence in Alzheimer's Disease Classification: A Systematic Review. *Cognitive Computation*, 16(1), 1-44. <https://doi.org/10.1007/s12559-023-10192-x>

Walczak, S., Pofahl, W. E., & Scorpio, R. J. (2003). A decision support tool for allocating hospital bed resources and determining required acuity of care. *Decision Support Systems*, 34(4), 445-456. [https://doi.org/10.1016/s0167-9236\(02\)00071-4](https://doi.org/10.1016/s0167-9236(02)00071-4)

Wysocki, O., Davies, J. K., Vigo, M., Armstrong, A. C., Landers, D., Lee, R., & Freitas, A. (2023). Assessing the communication gap between AI models and healthcare

professionals: Explainability, utility and trust in AI-driven clinical decision-making. *Artif Intell*, 316, 103839. <https://doi.org/10.1016/j.artint.2022.103839>