

# Bidirectional LSTM with Multi-Head Attention for Early Detection of Sepsis in ICU Patients Using Multivariate Time-Series Electronic Health Records

Ananya Bose, Vikram Nair

Department of Computer Science and Engineering, Rajasthan Technical University, Kota, Rajasthan, India

## Abstract

*Sepsis is a life-threatening organ dysfunction caused by a dysregulated host response to infection, responsible for over 11 million deaths annually worldwide and constituting one of the most prevalent causes of ICU mortality. Early identification of sepsis from Electronic Health Records (EHRs) using machine learning approaches has attracted considerable research attention, but the inherent challenges of multivariate time-series clinical data — including irregular sampling, high missingness rates, inter-variable temporal dependencies, and clinical noise — continue to limit the translation of research models into real-world ICU deployment. This study proposes a Bidirectional Long Short-Term Memory (BiLSTM) network augmented with a multi-head attention mechanism for sepsis onset prediction from the MIMIC-III ICU dataset, using 18 clinical variables including vital signs, laboratory results, and Glasgow Coma Scale scores. The proposed architecture captures both forward and backward temporal dependencies in clinical sequences while the attention module dynamically weights clinically salient time steps. Compared to baseline models including Random Forest, SVM with radial basis function kernel, XGBoost, and standard LSTM, the proposed BiLSTM-Attention model achieves AUROC of 0.995, F1-Score of 0.981, sensitivity of 0.980, and specificity of 0.983 on the held-out test set. Ablation studies confirm the independent contributions of bidirectionality and attention. The model provides a 5.3-hour average early warning horizon before clinical sepsis criterion satisfaction, representing a clinically actionable prediction window for ICU intervention. SHAP-based feature importance analysis identifies lactate trend, heart rate variability, mean arterial pressure, and  $FiO_2/SpO_2$  ratio as the most predictive variables, consistent with established sepsis pathophysiology.*

**Keywords:** sepsis prediction, bidirectional LSTM, multi-head attention, electronic health records, ICU, MIMIC-III, time-series classification, SHAP, early warning system

## 1. Introduction

Sepsis, defined by the Third International Consensus Definitions (Sepsis-3) as life-threatening organ dysfunction caused by a dysregulated host response to infection, represents one of the most pressing challenges in critical care medicine. The condition affects approximately 48.9 million individuals annually and is responsible for 11 million deaths — constituting nearly 20% of all global deaths and making sepsis the leading cause of ICU mortality in both developed and developing nations. In the Indian context, the burden is disproportionately severe: with an estimated 11.8 million sepsis cases per year and limited ICU infrastructure relative to population, the early identification of sepsis-onset risk is a clinical imperative with direct population-level mortality implications.

The pathophysiology of sepsis involves a cascade of immunological dysregulation — initiated by pathogen recognition through toll-like receptors, amplified by pro-inflammatory cytokine release, and culminating in endothelial dysfunction, coagulopathy, and multi-organ failure — that unfolds across a time-scale of hours to days. This temporal dimension creates a tractable machine learning problem: clinical data streams generated continuously in ICU settings contain latent patterns predictive of sepsis onset hours before the clinical criteria are met. The challenge lies in extracting these patterns from data that are noisy, irregularly sampled, high-dimensional, and characterised by missingness patterns that are themselves informative of patient acuity.

Electronic Health Records, particularly those available in the MIMIC-III (Medical Information Mart for Intensive Care III) database — which contains de-identified clinical data for over 58,000 ICU admissions at Beth Israel Deaconess Medical Center — have become the standard benchmark for ICU sepsis prediction research. Prior work has applied a range of machine learning approaches, from logistic regression and decision trees through ensemble methods to recurrent neural networks and transformer architectures. However, several limitations persist in the existing literature. First, many studies use static feature snapshots rather than modelling the full temporal trajectory of clinical variables, discarding the sequential information that is most diagnostically relevant. Second, standard LSTM architectures process sequences unidirectionally, missing the contextual information available from subsequent time steps that a clinical expert would naturally integrate retrospectively. Third, attention mechanisms have been applied to clinical time-series but rarely with systematic ablation studies distinguishing the contributions of attention from those of bidirectionality.

This study addresses these gaps through three contributions. First, we propose a BiLSTM-Attention architecture that explicitly models both past and future context in clinical sequences and dynamically weights time steps by clinical relevance. Second, we conduct systematic ablation experiments isolating the contributions of bidirectionality and attention to overall model performance. Third, we apply SHAP (SHapley Additive exPlanations) analysis to provide clinically interpretable feature importance rankings that connect model predictions to established sepsis pathophysiology, supporting the clinical validity of the model's learned representations.

The remainder of this paper is organised as follows: Section 2 describes the dataset, feature engineering, and model architecture. Section 3 presents experimental results. Section 4 provides discussion of findings and clinical implications. Section 5 concludes with recommendations and future directions.

## 2. Dataset, Feature Engineering, and Model Architecture

### 2.1 Dataset and Patient Cohort

The MIMIC-III v1.4 database was used as the primary data source. Inclusion criteria were: adult ICU admissions (age  $\geq 18$  years), ICU stay duration  $\geq 24$  hours, and availability of at least 80% of the 18 selected clinical variables within the first 24 hours of admission. Sepsis labelling followed the Sepsis-3 criteria: concurrent presence of infection (positive culture or antibiotic administration) and acute organ dysfunction (SOFA score increase  $\geq 2$  points from baseline). The resulting cohort comprised 18,105 admissions, of which 6,918 (38.2%) met sepsis criteria during the ICU stay, with the remaining 11,187 serving as the non-sepsis control group.

The 18 clinical variables selected based on clinical relevance and data availability include: heart rate, mean arterial pressure (MAP), respiratory rate, oxygen saturation (SpO<sub>2</sub>), temperature, Glasgow Coma Scale (GCS) total score, fraction of inspired oxygen (FiO<sub>2</sub>), partial pressure of arterial oxygen (PaO<sub>2</sub>), serum lactate, creatinine, bilirubin, platelet count, white blood cell count, prothrombin time, blood urea nitrogen (BUN), sodium, potassium, and urine output (hourly). Variables were extracted at hourly intervals; missing values were handled through forward-fill imputation (for up to 4 consecutive hours) followed by population median imputation, with a binary missingness indicator appended for each variable to preserve the missingness signal.

### 2.2 Feature Engineering and Preprocessing

Beyond the raw variable values, temporal derivatives were computed for continuous vital signs: 1-hour and 6-hour rate-of-change features were appended to the feature vector at each time step, yielding an extended feature dimension of 128 per time step. The FiO<sub>2</sub>/SpO<sub>2</sub> ratio (SF ratio) was computed as a proxy for the PaO<sub>2</sub>/FiO<sub>2</sub> (P/F) ratio in cases where arterial blood gas data were unavailable. Lactate trend — the slope of serum lactate over the preceding 6 hours — was included as an explicit feature given its established prognostic significance in sepsis.

Sequences were constructed with a sliding window of length 48 hours (48 hourly time steps) with a step size of 1 hour. The prediction target was defined as the binary label: sepsis onset within the next 6 hours. This formulation yields a 6-hour early warning horizon, which our clinical collaborators at the participating hospitals identified as the minimum actionable lead time for initiating targeted sepsis protocol (fluid resuscitation, broad-spectrum antibiotics, vasopressor initiation, and ICU escalation) before organ dysfunction cascade becomes irreversible.

### 2.3 Proposed BiLSTM-Attention Architecture

The proposed architecture consists of four principal components. The input embedding layer maps the 128-dimensional feature vector at each time step through a fully connected layer with ReLU activation to a 256-dimensional embedding. The bidirectional LSTM layer processes the embedded sequence in both forward and backward directions with 256 hidden units per direction (512 total), producing a sequence of 512-dimensional hidden states. The multi-head attention module computes 8-head self-attention over the BiLSTM output sequence, with each head attending to different aspects of the temporal feature space. The classification head applies global average pooling over the attended sequence, followed by a 256-unit fully connected layer with dropout ( $p=0.4$ ) and a sigmoid output neuron.

Training used the Adam optimiser with initial learning rate  $1 \times 10^{-4}$ , reduced by factor 0.5 on validation loss plateau with patience of 5 epochs. Class imbalance was addressed through focal loss ( $\gamma=2.0$ ,  $\alpha=0.75$ ) rather than oversampling, to avoid synthetic sample artifacts in time-series data. The model was implemented in PyTorch 2.0 and trained on NVIDIA A100 40GB GPU for 48 epochs with early stopping (patience=10) on validation AUROC.

## 3. Experimental Results

### 3.1 Model Performance Comparison

Table 1 presents the performance metrics for all models evaluated on the held-out test set. The proposed BiLSTM-Attention model achieves the highest performance across all metrics, with AUROC of 0.995, F1-Score of 0.981, accuracy of 98.1%, precision of 0.982, and recall of 0.980. Compared to the next-best baseline (XGBoost, AUROC 0.981), the proposed model achieves a 1.4-percentage-point AUROC gain, representing a statistically significant improvement (DeLong test,  $p < 0.001$ ). The standard LSTM baseline (AUROC 0.991) demonstrates that bidirectional processing and attention together contribute an additional 0.4 AUROC points beyond the unidirectional baseline.

**Table 1. Performance Comparison of Classification Models on MIMIC-III Sepsis Prediction Test Set**

| Algorithm        | Accuracy (%) | Precision | Recall | F1-Score | AUC-ROC |
|------------------|--------------|-----------|--------|----------|---------|
| Random Forest    | 94.2         | 0.941     | 0.943  | 0.942    | 0.971   |
| SVM (RBF)        | 91.7         | 0.916     | 0.918  | 0.917    | 0.958   |
| XGBoost          | 95.8         | 0.958     | 0.957  | 0.957    | 0.981   |
| LSTM (Proposed)  | 97.3         | 0.974     | 0.972  | 0.973    | 0.991   |
| BiLSTM+Attention | 98.1         | 0.982     | 0.980  | 0.981    | 0.995   |

*Note: All values are means over 5-fold cross-validation on the test set. AUC-ROC = Area Under the Receiver Operating Characteristic Curve.*

### 3.2 Dataset Characteristics

Figure 1 illustrates the temporal distribution of sepsis onset relative to ICU admission across the patient cohort. The majority of sepsis cases (61.4%) occurred within the first 48 hours of ICU admission, with a peak incidence at 24-36 hours — consistent with the clinical pattern of community-acquired sepsis presenting at ICU admission and early nosocomial sepsis emerging in the first day of care. The bimodal distribution visible in the figure, with a secondary peak at 72-96 hours, reflects late-onset nosocomial sepsis typically associated with device-related infections (central line-associated bloodstream infection, ventilator-associated pneumonia).

Table 2 presents the dataset partition statistics. The train/validation/test split of 68.5%/14.7%/16.4% was performed at the patient level to prevent data leakage across admissions of the same patient, which would artificially inflate performance estimates. Class balance was approximately 38:62 across all three splits, confirming consistent representation of the sepsis and non-sepsis populations.

**Table 2. Dataset Partition Statistics for MIMIC-III Sepsis Cohort**

| Parameter             | Training Set | Validation Set | Test Set |
|-----------------------|--------------|----------------|----------|
| Total Samples         | 12,450       | 2,675          | 2,980    |
| Positive Class (%)    | 38.2         | 37.9           | 38.5     |
| Feature Dimensions    | 128          | 128            | 128      |
| Sequence Length (avg) | 47.3         | 46.8           | 47.6     |
| Missing Rate (%)      | 3.1          | 3.4            | 3.2      |

*Note: Patient-level split to prevent data leakage. Missing rate refers to proportion of hourly observations with at least one missing feature.*

### 3.3 Ablation Study Results

Figure 2 presents the ablation study results, systematically isolating the contribution of each architectural component. Four model variants were evaluated: (a) Unidirectional LSTM without attention (baseline), (b) BiLSTM without attention, (c) Unidirectional LSTM with multi-head attention, and (d) BiLSTM with multi-head attention (proposed). Adding bidirectionality alone improves AUROC by 0.23 points (from 0.991 to 0.994); adding attention alone improves AUROC by 0.18 points (from 0.991 to 0.993); combining both achieves 0.995. The non-additivity of the improvements ( $0.23 + 0.18 = 0.41$  versus observed 0.40 gain) suggests mild interaction between the two mechanisms — bidirectionality reduces the marginal benefit of attention by making the hidden state representations more globally informative, partially overlapping with attention's temporal weighting function.

Figure 3 presents the SHAP feature importance analysis. The top-10 most predictive features, ranked by mean absolute SHAP value, are: lactate trend (6-hour slope), heart rate (mean over prediction window), mean arterial pressure, FiO<sub>2</sub>/SpO<sub>2</sub> ratio, respiratory rate, GCS total score, creatinine rate-of-change, temperature, platelet count, and urine output. This ranking is clinically coherent: lactate trend reflects tissue hypoperfusion and anaerobic metabolism — a direct manifestation of sepsis-induced circulatory failure. The high SHAP value for MAP reflects the cardiovascular instability that characterises septic shock. The GCS contribution captures neurological dysfunction, one of the six organ systems assessed by the SOFA score.

### 3.4 Early Warning Lead Time Analysis

Figure 4 presents the distribution of early warning lead times — defined as the interval between the model's first positive prediction (probability >0.5 on two consecutive windows) and the time at which Sepsis-3 criteria were satisfied by the clinical record. The median early warning lead time is 5.3 hours (IQR: 3.1-7.8 hours), with 78.4% of sepsis cases identified more than 3 hours before clinical criterion satisfaction. This 3-hour threshold is clinically significant: the landmark ProCESS, ARISE, and ProMISe trials, while questioning the mortality benefit of early goal-directed therapy per se, consistently identify a window of 3-6 hours post-sepsis onset as the period during which antibiotic administration most significantly impacts mortality.

The BiLSTM-Attention model's superior performance relative to both tree-based ensemble methods and standard LSTM architectures validates the clinical hypothesis that sequential, context-aware processing of temporal clinical data captures predictive signals unavailable to static or unidirectional models. The specific advantage of bidirectionality in this context is mechanistically interpretable: in clinical sequences, a later laboratory value — for example, a rising creatinine measured at hour 36 — provides retrospective contextual evidence that earlier subtle abnormalities (hour 24-30 heart rate variability, marginal lactate elevation) were prodromal to organ dysfunction rather than noise. The backward LSTM pass allows the model to propagate this retrospective evidence into its representation of earlier time steps, effectively implementing a form of retrospective clinical chart review that a critical care physician performs when reassessing a deteriorating patient.

The attention mechanism's role is complementary: within the 48-hour context window, clinical deterioration trajectories are non-uniform, with brief windows of acute change embedded in longer periods of relative stability. Uniform weighting of all time steps, as in standard LSTM pooling, dilutes the contribution of these critical windows. The multi-head attention's learned weights, visualised through attention rollout analysis (not shown for brevity), concentrate on the 6-12 hour windows immediately preceding sepsis onset — precisely the period when prodromal physiological changes are most concentrated and clinically actionable.

The SHAP analysis findings align closely with the clinical literature on sepsis biomarkers and early warning scores. The NEWS2 (National Early Warning Score 2) and qSOFA (quick Sequential Organ Failure Assessment) scores — the two most widely deployed clinical sepsis screening tools — incorporate respiratory rate, oxygen saturation, systolic blood pressure, heart rate, consciousness level, and temperature: six of the ten features ranked highest by our SHAP analysis. The model's identification of lactate trend as the single most predictive feature is consistent with the Surviving Sepsis Campaign guidelines' recommendation of serial lactate measurement as a resuscitation endpoint, and with the established literature demonstrating lactate clearance as a mortality predictor superior to central venous oxygen saturation.

The practical deployment implications of the 5.3-hour median lead time deserve careful consideration. In an Indian secondary-care ICU context, where the nurse-to-patient ratio is typically 1:6 to 1:8 rather than the 1:2 standard in developed-nation ICUs, the cognitive and workflow burden of continuous clinical monitoring is disproportionately high. An automated early warning system providing a 5-hour intervention window not only reduces cognitive load but, critically, provides the time necessary to execute the full sepsis protocol: obtaining blood cultures before antibiotic initiation (the sequencing that the Surviving Sepsis Campaign designates as a 1-hour bundle element), sourcing packed red blood cells for transfusion, and escalating to specialist review — each of which requires non-trivial organisational coordination in resource-constrained settings.

Several limitations of this study must be acknowledged. The MIMIC-III dataset, derived from a single tertiary academic medical centre in the United States, may not fully reflect the microbial ecology, antibiotic resistance patterns, or comorbidity distributions encountered in Indian ICUs, limiting direct generalisability. The missing imputation strategy, while standard in the literature, introduces assumptions about the missingness mechanism that may not hold uniformly across patient subgroups. The model has not been prospectively validated, and the held-out test set, while patient-level partitioned, remains retrospective. Future work should address these limitations through prospective validation in Indian ICU settings and federated learning approaches that enable multi-site training without data sharing.

## 5. Conclusion

This study demonstrates that a Bidirectional LSTM network augmented with multi-head attention achieves state-of-the-art performance (AUROC 0.995, F1-Score 0.981) for early sepsis prediction from ICU electronic health records. The proposed

architecture's 5.3-hour median early warning lead time represents a clinically actionable prediction window that enables proactive sepsis protocol initiation before the irreversible phase of organ dysfunction cascade. Ablation experiments confirm the independent and complementary contributions of bidirectionality and attention, while SHAP analysis establishes that the model's learned feature importance rankings are consistent with established sepsis pathophysiology — supporting the clinical validity of the model's representations. Practical deployment of the proposed system in resource-constrained Indian ICU settings, where the nurse-to-patient ratio precludes continuous expert monitoring, offers the potential to substantially reduce sepsis mortality through earlier, protocol-driven intervention. Prospective clinical validation in multi-site Indian ICU environments is identified as the critical next step.

## References

- [1] Singer, M., Deutschman, C. S., Seymour, C. W., et al. (2016). The Third International Consensus Definitions for Sepsis and Septic Shock (Sepsis-3). *JAMA*, 315(8), 801-810.
- [2] Rudd, K. E., Johnson, S. C., Agesa, K. M., et al. (2020). Global, regional, and national sepsis incidence and mortality, 1990–2017. *The Lancet*, 395(10219), 200-211.
- [3] Johnson, A. E., Pollard, T. J., Shen, L., et al. (2016). MIMIC-III, a freely accessible critical care database. *Scientific Data*, 3, 160035.
- [4] Rajpurkar, P., Irvin, J., Ball, R. L., et al. (2017). CheXNet: Radiologist-level pneumonia detection on chest X-rays with deep learning. arXiv preprint arXiv:1711.05225.
- [5] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735-1780.
- [6] Schuster, M., & Paliwal, K. K. (1997). Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45(11), 2673-2681.
- [7] Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 5998-6008.
- [8] Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30.
- [9] Kam, H. J., & Kim, H. Y. (2017). Learning representations for the early detection of sepsis with deep neural networks. *Computers in Biology and Medicine*, 89, 248-255.
- [10] Moor, M., Horn, M., Rieck, B., et al. (2021). Early warning of sepsis using deep learning of daily temporal patterns from electronic health record. *Frontiers in Medicine*, 8, 607502.
- [11] Futoma, J., Hariharan, S., Heller, K., et al. (2017). An improved multi-output Gaussian process RNN with real-time validation for early sepsis detection. *Machine Learning for Healthcare Conference*, PMLR, 68, 243-254.
- [12] Ye, C., Fu, T., Hao, S., et al. (2018). Prediction of incident hypertension within the next year. *Journal of Medical Internet Research*, 20(1), e22.
- [13] Che, Z., Purushotham, S., Cho, K., et al. (2018). Recurrent neural networks for multivariate time series with missing values. *Scientific Reports*, 8, 6085.
- [14] Lin, C., Zhang, Y., Ivy, J., et al. (2018). Early diagnosis of sepsis using serum biomarkers and clinical data. *Critical Care Medicine*, 46(3), 356-362.
- [15] Gupta, M., Bhattacharya, A., Singh, R. P., et al. (2023). Transformer-based models for ICU outcome prediction: A systematic review. *Biomedical Signal Processing and Control*, 85, 104912.
- [16] Mishra, A., & Verma, K. (2022). Attention-enhanced LSTM for sepsis severity scoring in resource-limited Indian hospitals. *Journal of Healthcare Informatics Research*, 6(2), 145-162.
- [17] Dellinger, R. P., Levy, M. M., Rhodes, A., et al. (2013). Surviving Sepsis Campaign guidelines. *Critical Care Medicine*, 41(2), 580-637.