

SepsisCalc: Integrating Clinical Calculators into Early Sepsis Prediction via Dynamic Temporal Graph Construction

Changchang Yin
yin.731@osu.edu
The Ohio State University
Columbus, Ohio, USA

Shihan Fu
sh.fu@northeastern.edu
Northeastern University
Boston, Massachusetts, USA

Bingsheng Yao
b.yao@northeastern.edu
Northeastern University
Boston, Massachusetts, USA

Thai-Hoang Pham
pham.375@osu.edu
The Ohio State University
Columbus, Ohio, USA

Weidan Cao
weidan.cao@osu.edu
The Ohio State University
Columbus, Ohio, USA

Dakuo Wang
d.wang@northeastern.edu
Northeastern University
Boston, Massachusetts, USA

Jeffrey Caterino
jeffrey.caterino@osumc.edu
The Ohio State University Wexner
Medical Center
Columbus, Ohio, USA

Ping Zhang
zhang.10631@osu.edu
The Ohio State University
Columbus, Ohio, USA

ABSTRACT

Sepsis is an organ dysfunction caused by a deregulated immune response to an infection. Early sepsis prediction and identification allow for timely intervention, leading to improved clinical outcomes. Clinical calculators (e.g., the six-organ dysfunction assessment of SOFA in Figure 1) play a vital role in sepsis identification within clinicians' workflow, providing evidence-based risk assessments essential for sepsis diagnosis. However, artificial intelligence (AI) sepsis prediction models typically generate a single sepsis risk score without incorporating clinical calculators for assessing organ dysfunctions, making the models less convincing and transparent to clinicians. To bridge the gap, we propose to mimic clinicians' workflow with a novel framework SepsisCalc to integrate clinical calculators into the predictive model, yielding a clinically transparent and precise model for utilization in clinical settings. Practically, clinical calculators usually combine information from multiple component variables in Electronic Health Records (EHR), and might not be applicable when the variables are (partially) missing. We mitigate this issue by representing EHRs as temporal graphs and integrating a learning module to dynamically add the accurately estimated calculator to the graphs. Experimental results on real-world datasets show that the proposed model outperforms state-of-the-art methods on sepsis prediction tasks. Moreover, we developed a system to identify organ dysfunctions and potential sepsis risks, providing a human-AI interaction tool for deployment, which can help

clinicians understand the prediction outputs and prepare timely interventions for the corresponding dysfunctions, paving the way for actionable clinical decision-making support for early intervention.

CCS CONCEPTS

• **Information systems** → **Data mining**; • **Applied computing** → **Health informatics**.

KEYWORDS

Early sepsis prediction, Electronic health record, Deep learning

ACM Reference Format:

Changchang Yin, Shihan Fu, Bingsheng Yao, Thai-Hoang Pham, Weidan Cao, Dakuo Wang, Jeffrey Caterino, and Ping Zhang. 2018. SepsisCalc: Integrating Clinical Calculators into Early Sepsis Prediction via Dynamic Temporal Graph Construction. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 14 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

1 INTRODUCTION

Sepsis, defined as life-threatening organ dysfunction in response to infection, contributes to up to half of all hospital deaths and is associated with more than \$24 billion in annual costs in the United States [24]. Existing studies [25] have shown that sepsis patients may benefit from a 4% higher chance of survival if diagnosed 1 hour earlier, so developing early sepsis prediction systems can significantly improve clinical outcomes. Over the past few decades, the rapid growth in volume and diversity of electronic health records (EHR) has made it possible to adopt state-of-the-art artificial intelligence (AI) methods to early identify patients with sepsis risk.

Deep learning (DL) methods have been widely used for early sepsis risk prediction tasks and have achieved superior performance. However, most existing DL or AI methods are data-driven and fail to incorporate the important and widely used clinical calculators. For example, Sequential Organ Failure Assessment (SOFA) [43] serves as an important screening tool for the identification of organ dysfunctions for sepsis in clinicians' workflow, as shown in Figure 1.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://www.acm.org).
Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/18/06
<https://doi.org/XXXXXXXX.XXXXXXX>

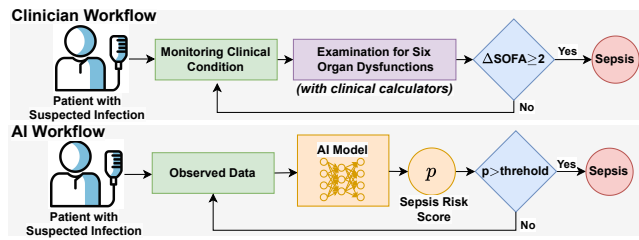


Figure 1: Workflows of clinicians and AI for sepsis identification. Clinicians examine sepsis by assessing organ dysfunctions with multiple clinical calculators as evidence, while AI workflow only gives an overall sepsis risk score.

Such a gap in the workflows makes the AI models less convincing and hinders their application to real-world clinical scenarios.

To bridge the gap, we propose to mimic clinicians' workflow by incorporating clinical calculators into automatic sepsis prediction tasks. However, directly feeding calculators into DL models presents a significant challenge: the calculators usually combine information from multiple component variables (e.g., SOFA has 6 component scores with 12 clinical variables, as shown in Table 6), while the variables might have high missing rates (e.g., the missing rates > 60% for most variables in Table 8), making calculators sometimes not applicable. An intuitive method to handle missing values is imputation. However, when the missing rate is high, the imputation models [27, 45, 47] become less accurate and introduce more bias, which could be harmful for downstream sepsis prediction tasks.

To address the challenges, we propose a novel early **Sepsis** prediction model with clinical **Calculators (SepsisCalc)**. For each patient, we first construct a temporal graph containing all the observed variables (including demographics, vital signs, lab tests, procedures, and medications). Then we use the graph to estimate the clinical calculators that can summarize six organ function statuses, which are dynamically added to the patient temporal graph, as shown in Figure 2(C). We only include the accurately estimated calculators, filtering out those with all components unobserved due to low confidence. Finally, we use a graph neural network (GNN) to extract the features of the dynamic temporal heterogeneous graph and make predictions for both organ dysfunctions and sepsis risks.

To demonstrate the effectiveness of our SepsisCalc, we conducted extensive experiments on real-world datasets MIMIC-III [17], AmsterdamUMCdb [40], and one proprietary dataset extracted from Ohio State University Wexner Medical Center (OSUWMC). Experimental results show that the proposed model outperforms state-of-the-art methods on the early sepsis prediction tasks. Moreover, we developed a system integrated into OSUWMC EHR system to identify organ dysfunctions and potential sepsis risks, providing a human-AI interaction tool for deployment and paving the way for actionable clinical decision-making support for early intervention.

We summarize our contributions as follows:

- We propose a novel sepsis prediction model SepsisCalc, which can represent patients' EHR data as dynamic temporal graphs, and effectively extract temporal information, clinical event

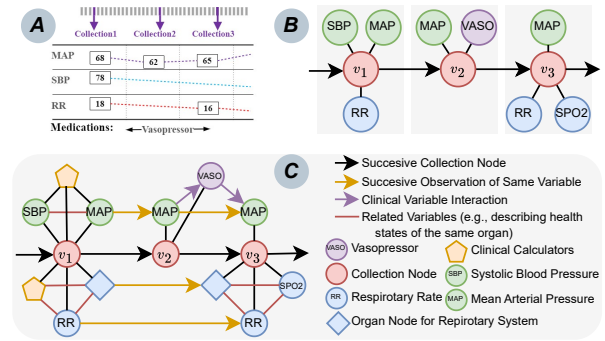


Figure 2: Different EHR representation methods. (A) An example of sequential EHR representation. (B) Example of graph representation with temporal information of clinical observations. (C) The proposed dynamic temporal graph representation with clinical event interaction and clinical calculators. Note that only partial calculator and organ nodes and edges are plotted for graph illustration in subfigure C.

interaction, and organ dysfunction information from the EHR data.

- We incorporate the widely-used and well-validated clinical calculators by dynamically generating the calculator nodes, which can significantly improve prediction performance and make the model more stable and convincing to clinicians.
- We conducted extensive experiments on various real-world datasets and the experimental results show that the proposed models outperform state-of-the-art methods on sepsis prediction tasks, demonstrating the effectiveness of the proposed SepsisCalc.
- We developed a system integrated into EHR system, allowing clinicians to easily use and effectively interact with the models.

2 RELATED WORK

In this section, we briefly review the existing work related to sepsis prediction systems.

Sepsis Screening Tools. Sepsis is a heterogeneous clinical syndrome that is the leading cause of mortality in hospital intensive care units (ICUs) [36, 47]. Early prediction and diagnosis may allow for timely treatment and lead to more targeted clinical interventions. Screening tools have been used clinically to recognize sepsis, including qSOFA [36], MEWS [38], NEWS [37], and SIRS [2]. However, those tools were designed to screen existing symptoms as opposed to explicitly early predicting sepsis before its onset.

Sequence-Based Models. With recent advances, deep learning methods have shown great potential for accurate sepsis prediction [15, 18, 35, 49]. Most of the studies represent patients' EHRs as the sequence of observations (see Figure 2(A)). Although the methods achieved superior performance, they face a critical challenge due to the data representation. The models need to completely observe a list of variables (including vital signs and lab tests), while many variables are missing in real-world data. Existing studies [15, 18, 49] usually impute the missing values before the

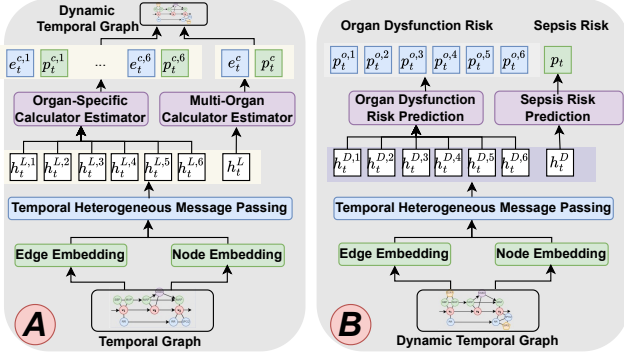


Figure 3: Framework of SepsisCalc. (A) Dynamic temporal graph construction. (B) Sepsis prediction framework.

prediction, raising a new problem that the sepsis prediction models will heavily rely on the imputation methods. The imputation bias would also be propagated to downstream prediction models.

Graph-Based Models. Graph representation can naturally handle the missing variables without imputation. However, most existing graph-based models [4, 26, 29, 46, 48] are designed to model longitudinal sparse EHRs (e.g., diagnosis codes and procedures with binary values) for chronic disease prediction (e.g., heart failure and COPD), the studies on dense float-value variables (e.g., vital signs and lab tests with multiple observations in the same visit) for acute diseases (e.g., sepsis) are still limited. Existing works [6, 23] construct a temporal graph with the observed variables (see Figure 2(B)), eliminating the need for additional imputation methods and avoiding potential imputation bias. However, they still suffer from two limitations: (i) Lack of consideration of clinical event interaction (e.g., vasopressor is used due to the extremely low MAP in Figure 2(B)); (ii) Lack of consideration of clinical calculators that provide clinicians with evidence-based risk assessments essential for accurate diagnosis and prognostic evaluation [16]. These limitations make these models unconvincing to clinicians, hindering their application in real-world clinical scenarios.

To address the challenges, we propose to dynamically construct temporal heterogeneous graphs (see Figure 2(C)) that (i) contain temporal relations between observations, (ii) include clinical event interaction, and (iii) estimate and integrate clinical calculators into graphs, to mimic clinicians’ workflow. Based on the dynamic temporal graph, we adopt graph neural networks to predict sepsis risk scores with potential organ dysfunction, paving the way for actionable clinical decision-making support for early intervention.

3 METHOD

In this section, we introduce the proposed SepsisCalc that dynamically constructs temporal heterogeneous graphs with clinical calculators and adopts novel GNNs to predict sepsis risks.

3.1 Notation and Problem Statement

We aim to predict the risk of sepsis with the observed clinical variables. We consider the following setup. A patient has a sequence of clinical variables (e.g., lab test and vital sign data) with timestamps.

Algorithm 1 Temporal Heterogeneous Message Passing

Input: temporal heterogeneous graph G ;

Output: collection node feature \mathbf{h}_t^L and organ node feature $\mathbf{h}_t^{L,i}$;

- 1: Obtain node embedding \mathbf{h}_v and edge embedding \mathbf{h}_r for each node v and edge r with Equation 1 and Equation 2;
- 2: **for** $l \leftarrow 1$ **to** L **do**
- 3: Calculate attention result $\text{Att}^i(v_e, v_t)$ with Equation 3;
- 4: Concatenate multiple head attention \mathbf{h}_v^l with Equation 5;
- 5: **end for**
- 6: Obtain collection node feature \mathbf{h}_t^L and organ node feature $\mathbf{h}_t^{L,i}$ with Equation 6;
- 7: Return \mathbf{h}_t^L and $\mathbf{h}_t^{L,i}$.

Let $X \in R^{T \times k}$ denote the observations of variables, where T denotes the number of collections of observations and k denotes the number of unique clinical variables. $Y \in \{0, 1\}^T$ denotes the binary ground truth of whether the patient will progress to sepsis in the coming hours. We represent patients’ data as temporal heterogeneous graphs. Given a loss function \mathcal{L} and a distribution over pairs (X, Y) , the goal is to learn the prediction function f by minimizing the expected loss: $f^* = \arg \min_f E[\mathcal{L}(f(X), Y)]$.

3.2 Static Temporal Graph Construction

We first construct a static temporal graph to represent the observed variables with timestamps in each patient’s EHRs.

3.2.1 Nodes. The temporal graph contains four kinds of nodes. The first is collection nodes which represent a collection of data observed in the same timestamps. The second kind of node is the observed clinical variables with the attributes of the observed values. The third kind of node is the organ nodes that describe the organ function statuses. The fourth kind of node is the calculator node that will be added in subsection 3.4.

3.2.2 Edges. We define three kinds of edges between the nodes. The first is directed edges between successive nodes of the same kind of variable (e.g., black and yellow arrows in Figure 2(C)). The time gaps between two observations are treated as edge attributes. The second kind of edge is the undirected edge between the nodes with the same timestamps. The third kind of edge is the directed relation to describe the interactions between clinical events. Patients might have received treatments (e.g., vasopressor used to avoid low MAP) before the identification of sepsis. Such interaction is important for patients’ EHR modeling and incorporated into the graph (e.g., purple arrows in Figure 2(C)).

3.3 Temporal Heterogeneous Message Passing

We use a graph encoder with temporal heterogeneous message passing to extract the features from the temporal graphs. Algorithm 1 describes the inference process of the module.

3.3.1 Clinical Embedding. We first map all the heterogeneous nodes and edges of the temporal graphs into a same embedding space.

Node Embedding. An embedding layer is used to map each node v into a fixed-sized vector $e^v \in R^d$. We also embed the observed

values as vectors for the nodes with float-value attributes. Following [47], we adopt value embedding to map the observed value as vector $e^{v'}$ and use time embedding to map the time gap δ as e^δ . The concatenation of the node embedding and the value embedding is sent to a linear mapping layer to generate $\mathbf{h}_v \in \mathbb{R}^d$, containing the information of node type, variable name, and observed value:

$$\mathbf{h}_v = L([e^v; e^{v'}]), \quad (1)$$

where $L(\cdot)$ denotes linear mapping functions, with each instance representing a distinct mapping. $[\cdot; \cdot]$ is a concatenation operation. The details of the embedding can be found in subsection A.3.2. **Edge Embedding.** Similarly, we use an embedding layer to map each edge r into a fixed-size vector $e^r \in \mathbb{R}^d$. For the directed edge r with elapsed time, we combine the edge embedding e^r with the time embedding e^δ to generate $\mathbf{h}_r \in \mathbb{R}^d$:

$$\mathbf{h}_r = L([e^r; e^\delta]) \quad (2)$$

3.3.2 Heterogeneous Message Passing. Given the temporal graph and embeddings, we leverage a temporal-aware message-passing mechanism for heterogeneous graphs to effectively gather temporal information, clinical event interaction, and historical observations.

We represent the features of clinical nodes in the (l) -th layer of the network as $\mathbf{h}_v^{(l)} \in \mathbb{R}^d$. They also serve as the input for the subsequent $(l+1)$ -th layer.

The entire message-passing process can be formalized into two stages: aggregation and updating. The first step is to aggregate the information of neighboring nodes. Specifically, we use an attention mechanism to weigh and integrate the features of neighboring nodes and concatenate the output from multi-head attention to obtain the final message. Taking as an example the process of propagating the features of neighbor nodes $v_e \in N(v_t)$ to the collection node v_t , the i -th head attention is as follows:

$$\mathbf{q}_v = L(\mathbf{h}_{v_t}^{(l)}), \quad \mathbf{W}_r = L(\mathbf{h}_r), \quad \mathbf{k}_{v_e} = L(\mathbf{h}_{v_e}^{(l)}), \\ \text{Att}^i(v_e, v_t) = \frac{\mathbf{q}_v \mathbf{W}_r \mathbf{k}_{v_e}^\top}{\sqrt{d}}, \quad (3)$$

where v_e represents the clinical node such as lab test, vital sign, and procedure contained in t -th collection as well as previous collection node v_{t-1} . r denotes the edge between nodes v_t and v_e . The target node for message propagation is t -th collection node v_t . After obtaining the attention scores for different neighboring nodes, we combine them with the mapped neighbor features to complete the entire aggregation process. We formalize it as follows:

$$\tilde{\mathbf{h}}_v^i = \text{Softmax}_{v_e \in N(v_t)} \left(\text{Att}^i(v_e, v_t) \right) \cdot L(\mathbf{h}_{v_e}^{(l)}), \quad (4)$$

where $N(v_t)$ denotes neighbors of v_t , and $\tilde{\mathbf{h}}_v^i$ represents the message passed to the collection node by the i -th head attention. Following the acquisition of aggregated information, the next stage is to combine this aggregated information with the target node's ego information to update the representation of the target node:

$$\mathbf{h}_v^{(l+1),i} = \gamma \cdot L\left(\text{ReLU}\left(\tilde{\mathbf{h}}_v^i\right)\right) + (1-\gamma) \cdot \mathbf{h}_v^{(l)}, \\ \mathbf{h}_v^{(l+1)} = \parallel_{i \in [1,h]} \mathbf{h}_v^{(l+1),i} \quad (5)$$

where $\parallel_{i \in [1,h]}$ denotes concatenating the outputs of multiple heads, γ is learnable coefficient for the skip connection. Leveraging heterogeneous message passing, we integrate the features of medical events in the collection, along with the graphical structure and historical collection information, into the current collection node, thereby updating the patient's representation.

3.3.3 Organ-Specific Node and Collection Node Representation. After L layers of temporal heterogeneous message passing, we obtain the collection node features $\mathbf{h}_{v_t,c}^L$ and six organ node features $\mathbf{h}_{v_t,o,i}^L$ ($1 \leq i \leq 6$). We use linear functions to generate the node representation $\mathbf{h}_t^L, \mathbf{h}_t^{L,i} \in \mathbb{R}^d$ for further dynamic graph construction:

$$\mathbf{h}_t^{L,i} = L(\text{ReLU}(\mathbf{h}_{v_t,o,i}^L)), \quad \mathbf{h}_t^L = L(\text{ReLU}(\mathbf{h}_{v_t,c}^L)) \quad (6)$$

3.4 Dynamic Temporal Graph Construction

As Figure 3(A) shows, we add the accurately estimated calculators to temporal graphs. We estimate organ-specific calculators (e.g., DIC [44]) and multi-organ calculators (e.g., SOFA [43]) with the same structure. In this subsection, we use multi-organ calculator generation with \mathbf{h}_t^L as an example.

3.4.1 Calculator Score Estimation. We estimate the values of calculators $e_t^c \in \mathbb{R}^c$ for all the c calculators:

$$e_t^c = L(\mathbf{h}_t^L) \quad (7)$$

When all the component variables are observed, we have the ground truth for the calculators and use Mean Square Error (MSE) loss to train the calculator estimation module:

$$\mathcal{L}_e = \frac{1}{T} \sum_{t=1}^T \frac{1}{\sum_i M_{t,i}} \sum_{i=1}^c (e_{t,i}^c - \hat{e}_{t,i}^c)^2 M_{t,i}, \quad (8)$$

where $\hat{e}_{t,i}^c$ denotes the ground truth of the clinical calculators. $M_{t,i}$ denotes an indicator variable. $M_{t,i}$ is 1 if the ground truth $\hat{e}_{t,i}^c$ is available, and 0 else. When the ground truth is not available, the module is jointly trained with the sepsis prediction framework.

3.4.2 Dynamic Node Construction. When missing rates are high, and the generated calculators' scores are inaccurate, consideration of such calculators could introduce additional bias and even mislead clinicians when providing clinical decision-making support. We use a confidence score $p_t^c \in \mathbb{R}^c$ to filter out the missing calculators with low confidence (i.e., $p_t^c < 0.5$).

$$p_t^c = \text{Sigmoid}(L(\mathbf{h}_t^L)) \quad (9)$$

We use the following objective to train the dynamic node construction module:

$$\mathcal{L}_d = \frac{1}{T} \sum_{t=1}^T \frac{1}{\sum_i M_{t,i}} \sum_{i=1}^c [-y_{t,i}^c \log p_{t,i}^c - (1 - y_{t,i}^c) \log(1 - p_{t,i}^c)] M_{t,i}, \quad (10)$$

where $y_{t,i}^c = I[(e_{t,i}^c - \hat{e}_{t,i}^c)^2 < 0.01]$ and $I[\cdot]$ is an indicator function that returns 1 if the statement is true; otherwise, 0.

3.4.3 New Edge Generation. To incorporate the clinical calculator mechanism, the model also automatically generates the edges between the generated nodes and their component variables.

Algorithm 2 SepsisCalc

Input: static temporal graph G , calculator ground truth \hat{e}_t , outcome y_t , learning rate lr ;

- 1: **repeat**
- 2: # *Dynamic temporal graph construction*
- 3: Obtain collection node feature \mathbf{h}_t^L and organ node feature $\mathbf{h}_t^{L,i}$ with Algorithm 1 and temporal graph G ;
- 4: Estimate clinical calculators e_t^c with Equation 7;
- 5: Compute calculator estimation loss \mathcal{L}_e with Equation 8;
- 6: Estimate calculator confidence p_t^c with Equation 9;
- 7: Compute calculator confidence loss \mathcal{L}_d with Equation 10;
- 8: Obtain dynamic temporal graph G_d by adding calculators with high confidence to temporal graph G ;
- 9: # *Sepsis risk and organ dysfunction prediction*
- 10: Obtain collection node feature \mathbf{h}_t^D and organ node feature $\mathbf{h}_t^{D,i}$ with Algorithm 1 and dynamic graph G_d ;
- 11: Estimate the sepsis risk p_t and organ dysfunction risk $p_t^{o,i}$ with Equation 11;
- 12: Compute prediction loss \mathcal{L}_c and \mathcal{L}_o with Equation 12;
- 13: Update parameters by minimizing the loss \mathcal{L} in Equation 13;
- 14: **until** Convergence.

3.5 Sepsis and Organ Dysfunction Prediction

After dynamically constructing the temporal graph, we re-extract the features of the new graphs with the same temporal heterogeneous message passing module as in subsection 3.3. We use \mathbf{h}_t^D and $\mathbf{h}_t^{D,i}$ to denote the extracted features for the collection node and i -th organ node at time t from the dynamic temporal graph and continue to predict the clinical risk as shown in Figure 3(B).

3.5.1 Risk Prediction. We use a linear layer followed with a Sigmoid layer to generate the sepsis risk $p_t \in R$ and organ dysfunction risk $p_t^{o,i} \in R$ ($i = 1, 2, \dots, 6$):

$$\begin{aligned} p_t &= \text{Sigmoid}(L(\mathbf{h}_t^D)), \\ p_t^{o,i} &= \text{Sigmoid}(L(\mathbf{h}_t^{D,i})), \end{aligned} \quad (11)$$

3.5.2 Objective Prediction. We use binary-cross-entropy loss to train the framework:

$$\begin{aligned} \mathcal{L}_c &= \frac{1}{T} \sum_{t=1}^T -y_t \log(p_t) - (1 - y_t) \log(1 - p_t), \\ \mathcal{L}_o &= \frac{1}{T} \sum_{t=1}^T \frac{1}{6} \sum_{i=1}^6 -y_t^{o,i} \log(p_t^{o,i}) - (1 - y_t^{o,i}) \log(1 - p_t^{o,i}) \end{aligned} \quad (12)$$

The whole framework is trained with a weighted loss:

$$\mathcal{L} = \mathcal{L}_c + \alpha_o \mathcal{L}_o + \alpha_e \mathcal{L}_e + \alpha_d \mathcal{L}_d, \quad (13)$$

where $\alpha_o, \alpha_e, \alpha_d > 0$ are hyper-parameters. Algorithm 2 describes the training process of the framework.

4 EXPERIMENT

To demonstrate the effectiveness of the proposed SepsisCalc, we conducted extensive experiments on multiple real-world datasets.

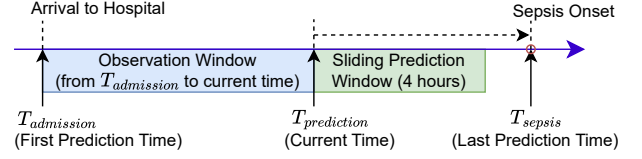


Figure 4: Setting of sepsis onset prediction.

4.1 Datasets and Setup

Datasets. We validated our model on two publicly available datasets (MIMIC-III¹ and AmsterdamUMCdb²) and one proprietary dataset extracted from OSUWMC³. We first extracted all the sepsis patients with suspected infection [36] in the datasets. Patients meeting sepsis-3 criteria [36] are defined as case patients, while the others with only suspected infection are treated as control patients. Following [21, 47], we extracted 30 vital signs and lab tests for sepsis patient status modeling. The statistics of the three datasets are displayed in Table 1. The used single-organ and multi-organ calculators are summarized in Table 5 and Table 7. The used variables and additional details can be found in subsection A.4.1.

Setup. Figure 4 displays the setting of the experiments. After the patients arrive at the hospital, we start to predict whether the patients will suffer from sepsis with a sliding 4-hour prediction window. We run the prediction process hourly until the patients have been diagnosed with sepsis or discharged.

4.2 Comparison Methods

To validate the performance of the proposed SepsisCalc for sepsis prediction, we implemented various models, including clinical calculator-based methods (i.e., NEWS [37], MEWS [38], qSOFA [36], SIRS [2]), RNN-based methods (GRU [3], LSTM [14], DFSP [8]), attention-based methods (RETAIN [5], Dipole [28]), graph-based methods (GTN [6], RGNN [23]). The details of the comparison methods can be found in subsection A.4.3.

We also implemented various versions of the proposed model. **SepsisCalc** is the main version. **SepsisCalc⁻ⁱ** is the simplest version that uses the same graph construction method as [6, 23] (see Figure 2(B)), without any graph interaction and calculators. **SepsisCalc^{imp}** uses an imputation method [27] to replace the dynamic graph construction module. **SepsisCalc^{-d}** removes the dynamic graph construction module. **SepsisCalc^{-o}** removes the organ dysfunction prediction module.

4.3 Implement Details

We implement our proposed model with Python 3.8.10 and PyTorch 1.12.1⁴. For training models, we use Adam optimizer with a mini-batch of 64 patients. The multi-modal data are projected into a 512-d space. We randomly divide the patients in each dataset into 10 sets. All the experiment results are averaged from 10-fold cross-validation, in which 7 sets are used for training, 1 set for validation, and 2 sets for testing. The validation sets are used to determine the best values of parameters in the training iterations.

¹<https://mimic.physionet.org/>

²<https://amsterdammedicaldatascience.nl>

³<https://wexnermedical.osu.edu/>

⁴<https://pytorch.org/>

Table 1: Statistics of MIMIC-III, AmsterdamUMCdb, and OSUWMC datasets.

	MIMIC	AmsterdamUMCdb	OSUWMC
#. of patients	21,686	6,560	85,181
#. of male	11,862	3,412	41,710
#. of female	9,824	3,148	43,471
Age (mean \pm std)	60.7 \pm 11.6	62.1 \pm 12.3	59.3 \pm 16.1
Missing rate	65%	68%	75%
Sepsis rate	32%	35%	29%

For the sepsis prediction tasks, we use Area Under the Receiver Operating Characteristic Curve (AUC), F1 and Recall for evaluation metrics at the collection level (with each collection treated as a separate sample). For the calculator estimation tasks, following [27, 45], we measure the models' performance with normalized Root Mean Square Error (nRMSE). The code and more implementation details can be found in subsection A.4.4 and GitHub⁵.

5 RESULTS

We now report the performance of the proposed model in the three datasets. We focus on answering the following research questions using our experimental results:

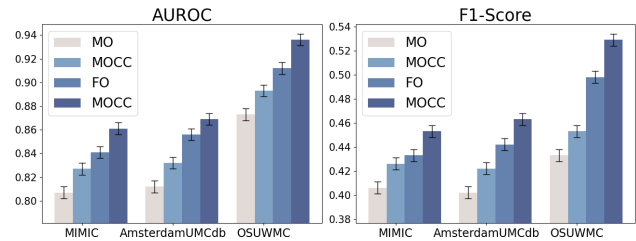
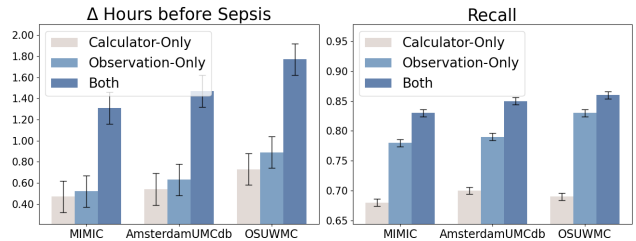
- **Q1: Why must we incorporate the clinical calculator scores?**
- **Q2: Are the estimated clinical calculator scores effective?**
- **Q3: How do estimated calculator scores improve early sepsis prediction system?**

5.1 Q1: Why must we incorporate the clinical calculator scores?

5.1.1 Wide Adoption of Clinical Calculators. Clinical calculators have emerged as indispensable tools within healthcare settings, providing clinicians with evidence-based risk assessments essential for accurate diagnosis and prognostic evaluation [9, 11, 16]. In the context of sepsis, numerous calculators, such as SOFA [36, 43], qSOFA [7, 33], MEWS [38], NEWS [37], and SIRS [2], are extensively studied in existing sepsis-related literature [10, 13, 31, 34], and widely employed as early warning tools in real-world hospital EHR systems for both ICU and hospital wards [1, 36]. Integrating clinical calculators into early sepsis prediction models can align them more closely with clinicians' workflows and enhance comprehensibility.

5.1.2 Improvement from Calculators on Sepsis Prediction Performance. Compared to raw clinical variables, clinical calculators can summarize the patients' health states at a high level (*e.g.*, in single or multiple organ dysfunction levels). We first design experiments to demonstrate the effectiveness of clinical calculators on early sepsis prediction with four settings:

- **Full Observation without Clinical Calculator (FO):** All the component variables of calculators are available while the clinical calculators are not used.

**Figure 5: Sepsis risk prediction performance in both full and missing observation settings.****Figure 6: Average alert time before sepsis and recall.**

- **Full Observation with Clinical Calculator (FOCC):** All the component variables of calculators are available and the clinical calculators are also included when conducting sepsis prediction.
- **Missing Observation without Clinical Calculator (MO):** Only partial component variables of clinical calculators are observed and the calculators are not incorporated.
- **Missing Observation with Clinical Calculator (MOCC):** Partial component variables of clinical calculators are observed and we still use the variables to compute the calculators.

Figure 5 displays the results for sepsis prediction in the four settings. The results show that in both full observation and missing observation settings, clinical calculators (*i.e.*, FOCC and MOCC) can consistently improve sepsis prediction performance, which shows the effectiveness of such domain knowledge in clinical tasks.

5.1.3 Early Sepsis Alerts with Clinical Calculators. Due to the fast-development characteristics of sepsis, delayed identification and treatment will significantly reduce the patients' survival rates [25]. It is critical to early identify the patients with sepsis risks. We also reported the average sepsis alert time before sepsis and recall at the patient level (with each patient treated as a separate sample for evaluation), as shown in Figure 6. The sepsis prediction models that incorporate both raw observations and calculators can predict sepsis approximately 1 hour earlier and achieve higher recall compared to models using only raw observational data, further demonstrating the credibility and effectiveness of the calculators.

5.2 Q2: Are the estimated clinical calculator scores effective?

5.2.1 Missing Rate of Clinical Calculators. Directly incorporating the clinical calculators might not be applicable for two reasons: (i) Lots of risk calculators (*e.g.*, SOFA) aggregate the values of clinical variables in a specific time span and thus are not immediately

⁵<https://github.com/yinchangchang/SepsisCalc>

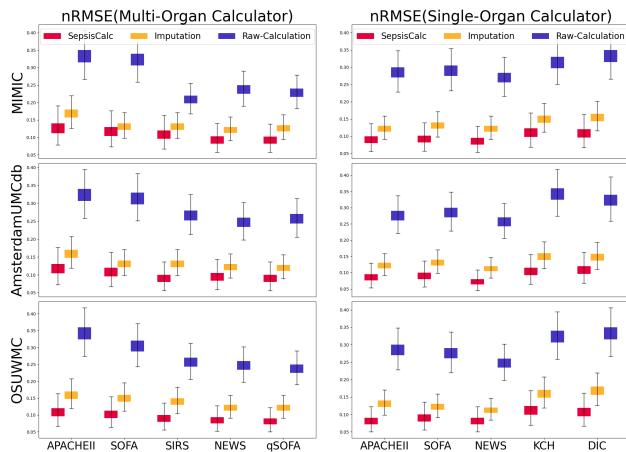


Figure 7: nRMSE of clinical calculator estimation (mask observation setting). All the component variables of the calculators are observed and the ground truths of calculators are available. We randomly mask 70% component variables. Raw-calculation means the original clinical methods that use the latest observed variables to compute the calculators.

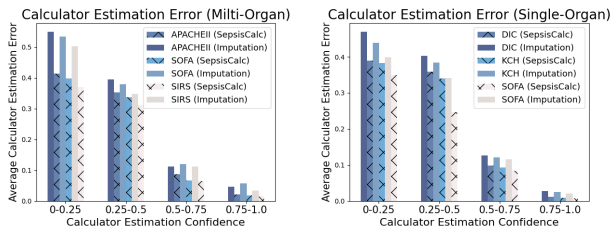


Figure 8: Calculator estimation error over confidence levels.

available, which limits their usage in timely sepsis prediction and detection. (ii) Most risk calculators usually combine the information from multiple variables. In real-world settings, the variables might have high missing rates and not always be available (especially for blood lab tests). Table 8 in Appendix displays the missing rates of the important lab tests and calculators related to sepsis.

5.2.2 Effectiveness of Clinical Calculator Estimation. To address the problem, we propose to estimate the calculators. In this subsection, we evaluated the performance of clinical calculator estimation.

Full Observation Setting. When all the component variables are observed, the ground truths of the clinical calculators are available, we used nRMSE to evaluate the calculator estimation performance. Table 10 in Appendix shows that the nRMSE between the ground truths and the estimated calculators is close to 0, demonstrating that our model can accurately learn the computation mechanisms of clinical calculators.

Mask Observation Setting. To further validate the calculator estimation performance when component variables are missing, we randomly masked 70% component variables to keep them consistent

with real-world missing rates (see Table 8). Figure 7 displays the results of our SepsisCalc, imputation [27] and original calculator computation mechanisms (*i.e.*, Raw-Calculation). SepsisCalc achieved much smaller estimation errors than Raw-Calculation, leading to the improved performance of downstream sepsis prediction tasks.

When component variables are missing, the ground truths of the clinical calculators are not available. Instead, we use the performance of downstream sepsis prediction tasks to validate the effectiveness of calculator estimation tasks. The detailed results are displayed in subsection 5.3.2.

5.2.3 Effectiveness of Clinical Calculator Generation Confidence.

When the missing rates are relatively high, the estimated clinical calculators might be inaccurate. We propose to use the calculator estimation confidence p_i^c in Equation 9 to filter the inaccurately estimated calculator nodes (*i.e.*, $p_i^c < 0.5$). Figure 8 shows that when the confidence of generation is low, calculator estimation performance suffers from a significant decline. Existing imputation models always give imputation results for the missing values, which might introduce more imputation bias to the prediction models and could be harmful for high-stake clinical applications. Moreover, the error-prone imputation in high-missing-rate settings could further mislead clinicians when providing clinical decision-making support. Our dynamic graph construction module only generates the clinical calculators with high confidence, which achieves a good trade-off between introducing more domain-specific knowledge and reducing imputation bias.

5.3 Q3: How do estimated calculator scores improve early sepsis prediction system?

5.3.1 Early Sepsis Prediction. Table 2 displays the sepsis prediction results. All the deep learning methods outperform the early-warning scores (*i.e.*, NEWS, MEWS, qSOFA, SIRS), which shows the promising potential of state-of-the-art deep learning models in real-world clinical applications. Although human-designed calculators are effective, deep-learning methods can capture abnormal values and more complicated temporal patterns inside EHRs.

Compared with attention-based models and graph-based models, the proposed SepsisCalc achieved the best prediction performance. By considering both observations and the estimated clinical calculators, the proposed SepsisCalc can model the organ dysfunctions better, which further improves the performance. The combination of human-designed clinical calculators and end-to-end deep learning methods can not only achieve better performance but also enhance credibility in real-world applications.

5.3.2 Ablation Study. To validate the performance improvement from dynamic calculator generation, we conducted an ablation study with various versions of the proposed model. Table 3 displays the ablation study results. Without the clinical event interaction, SepsisCalc⁻ⁱ performs worse than other versions, demonstrating the effectiveness of the proposed static graph construction method (*i.e.*, with medical event interaction and successive connection of the same variables). SepsisCalc outperforms SepsisCalc^{-d} (without dynamic calculator estimation), demonstrating the effectiveness of clinical calculators in sepsis prediction tasks. SepsisCalc outperforms SepsisCalc^{imp} that use imputation for further calculator

Table 2: Sepsis prediction results.

Method	MIMIC-III			AmsterdamUMCdb			OSUWMC		
	AUC	F1	Recall	AUC	F1	Recall	AUC	F1	Recall
NEWS	0.722±0.012	0.366±0.013	0.620±0.012	0.731±0.013	0.370±0.013	0.627±0.013	0.765±0.012	0.387±0.013	0.656±0.012
MEWS	0.726±0.013	0.372±0.013	0.624±0.012	0.735±0.012	0.376±0.012	0.631±0.013	0.769±0.013	0.393±0.012	0.660±0.013
qSOFA	0.729±0.011	0.374±0.011	0.631±0.011	0.738±0.012	0.379±0.012	0.639±0.012	0.772±0.011	0.396±0.011	0.668±0.011
SIRS	0.733±0.011	0.376±0.012	0.642±0.011	0.742±0.012	0.381±0.012	0.650±0.011	0.776±0.012	0.398±0.012	0.680±0.012
GRU	0.801±0.012	0.397±0.012	0.696±0.013	0.807±0.012	0.400±0.012	0.701±0.012	0.872±0.012	0.432±0.012	0.758±0.012
LSTM	0.807±0.013	0.408±0.012	0.698±0.012	0.813±0.012	0.411±0.012	0.703±0.012	0.879±0.013	0.444±0.012	0.760±0.012
RETAIN	0.814±0.013	0.418±0.013	0.710±0.014	0.820±0.013	0.421±0.013	0.715±0.014	0.886±0.014	0.455±0.013	0.773±0.013
Dipole	0.817±0.014	0.423±0.013	0.703±0.013	0.823±0.014	0.426±0.014	0.709±0.014	0.889±0.013	0.461±0.013	0.766±0.014
DFSP	0.822±0.011	0.424±0.012	0.696±0.011	0.828±0.011	0.425±0.011	0.702±0.011	0.894±0.012	0.465±0.012	0.758±0.012
RGNN	0.819±0.013	0.424±0.012	0.698±0.012	0.825±0.013	0.427±0.013	0.703±0.013	0.892±0.013	0.467±0.012	0.760±0.012
GTN	0.821±0.014	0.423±0.014	0.707±0.013	0.827±0.013	0.426±0.014	0.712±0.014	0.893±0.013	0.465±0.014	0.770±0.014
SepsisCalc	0.839±0.011	0.438±0.012	0.729±0.011	0.848±0.012	0.442±0.011	0.735±0.012	0.918±0.012	0.479±0.011	0.791±0.012

Table 3: Ablation study.

Method	MIMIC-III			AmsterdamUMCdb			OSUWMC		
	AUC	F1	Recall	AUC	F1	Recall	AUC	F1	Recall
SepsisCalc ^{imp}	.825	.426	.713	.835	.433	.719	.908	.466	.775
SepsisCalc ⁻ⁱ	.820	.422	.712	.829	.427	.716	.900	.463	.772
SepsisCalc ^{-d}	.830	.427	.715	.839	.432	.715	.910	.468	.775
SepsisCalc ^{-o}	.831	.429	.718	.841	.431	.716	.911	.465	.781
SepsisCalc	.839	.438	.729	.848	.442	.735	.918	.479	.791

computation, demonstrating the effectiveness of the dynamic graph construction. We speculate the reason is that the imputation results might be not accurate during the high-missing-rate settings and introduce more bias to the downstream sepsis prediction tasks, while the proposed SepsisCalc only includes the accurately estimated calculators with high confidence in the graphs. By adding the multi-task learning for both organ dysfunction and sepsis risk prediction tasks, SepsisCalc performs better than SepsisCalc^{-o}, further showing that the organ dysfunction identification tasks can help the sepsis prediction tasks.

6 DEPLOYMENT

Based on the sepsis prediction model, we implemented a system deployed in the Epic EHR Systems⁶ at OSUWMC (see Figure 9). The system starts to collect patients' data after the patients arrive at hospitals and automatically predicts sepsis risks hourly.

The interactive process with our system is visualized in the provided UI (Figure 9). In this scenario, a clinician is examining high-risk patients. After reviewing the patient list (Figure 9(A)), the clinician selects a patient, prompting the Patient Demographics (Figure 9(B)) section update to display his profile, including age, gender, weight, admission department, and sepsis risk score.

The clinician then focuses on the SOFA Score (Figure 9(C)) to assess the overall organ dysfunctions. Due to the missing variables,

the observed total SOFA score is not applicable, and our SepsisCalc estimates the SOFA score as 12. The radar chart provides a visual summary of scores across different organs. The gray shaded area displays the original SOFA scores at the beginning of the ICU admission, while the blue shaded area represents the observed SOFA (solid blue lines) and estimated SOFA (dashed red lines) scores.

To delve deeper, the clinician clicks the respiration area in the radar chart and reviews how the SOFA score for respiration has evolved in Figure 9(D), and the details of relevant vital signs (*i.e.*, respiratory rate and SpO₂), and lab test results (*e.g.*, PaO₂, SaO₂, HCO₃) in Figure 9(E). We also provide a trend view to display vital signs and lab tests so clinicians can review the history of variables and critical changes in Figure 10(B) in Appendix.

The clinician can also examine other organ SOFA scores and their specific details from the radar chart. By clicking on different areas of the radar chart, the clinician can view the scores and details for the central nervous system (CNS), coagulation, renal, liver, and cardiovascular systems. This allows for a comprehensive assessment of each organ's functional status, thereby facilitating more precise clinical decision-making.

Note that we used OSUWMC data for our algorithm illustration. All patients' names and demographic info in this Figure 9 are randomly generated for illustration purposes. Ongoing deployment also includes recruiting clinicians for usability evaluation to quantitative and qualitatively measure clinical outcomes and user satisfaction of SepsisLab (OSUWMC IRB#: 2020H0018).

7 CONCLUSION

In this work, we aim to develop transparent and convincing models for the real-world sepsis prediction tasks. We propose a novel framework SepsisCalc, which represents patients' EHR data as dynamic temporal graphs and effectively extracts temporal information, clinical event interactions, and organ dysfunction information from the graphs with a temporal heterogeneous message passing module. We introduce a dynamic graph construction module to estimate and integrate clinically widely used calculators into sepsis prediction

⁶<https://www.epic.com/software/>

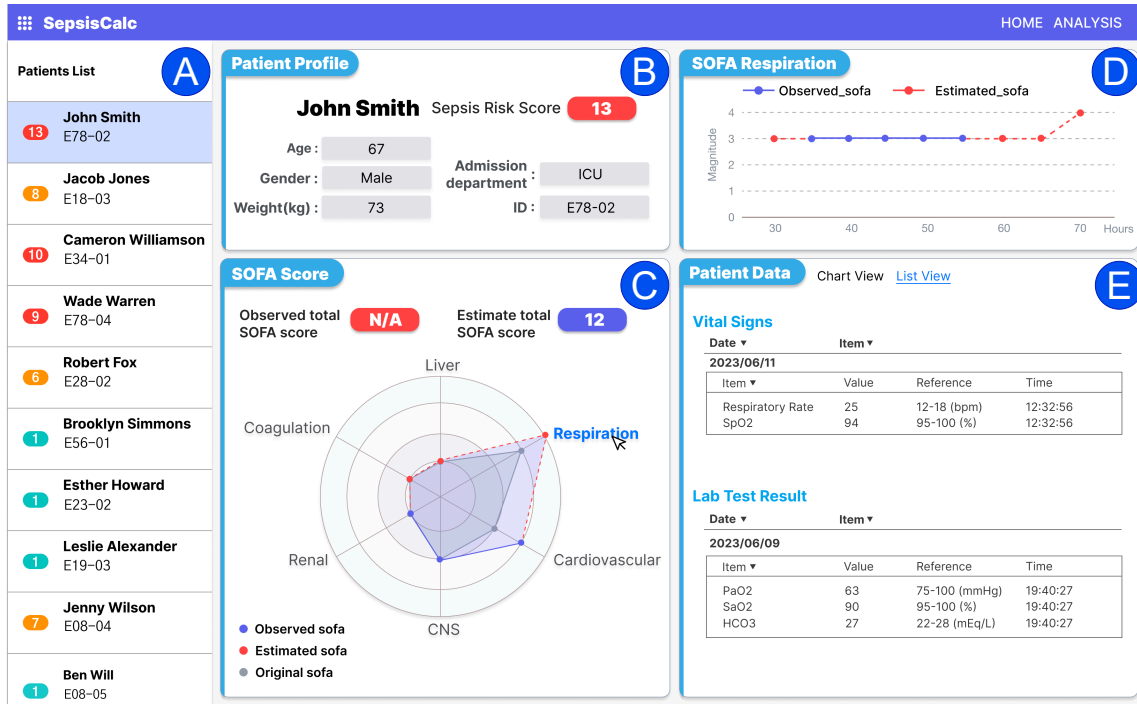


Figure 9: User Interface of SepsisCalc System. (A) Patient list with sepsis risk score. (B) Demographic information. (C) Overall SOFA scores. (D) Organ-specific SOFA score. (E) Vital signs and lab test results related to the specific organ.

models to help assess organ dysfunctions, aligning well with clinicians’ workflows for sepsis identification. Our graph construction method naturally handles the missing values by including only observed variables and high-confidence calculators, thereby avoiding the potential biases introduced by imputation methods that most sepsis prediction models suffer from. Experiments on three real-world datasets show that SepsisCalc can not only accurately estimate the calculators to assess the organ dysfunctions (even with missing values), but also outperform state-of-the-art clinical risk prediction methods, demonstrating the effectiveness of SepsisCalc. Finally, we design a system to display the identified organ dysfunctions and potential sepsis risks, providing a human-AI interaction tool for deployment and paving the way for actionable clinical decision-making support for early intervention.

8 ACKNOWLEDGMENTS

This work was funded in part by the National Science Foundation under award number IIS-2145625, by the National Institutes of Health under award number R01AI188576.

REFERENCES

- [1] Poushali Bhattacharjee, Dana P Edelson, and Matthew M Churpek. 2017. Identifying patients with sepsis on the hospital wards. *Chest* 151, 4 (2017), 898–907.
- [2] Roger C Bone, Robert A Balk, Frank B Cerra, R Phillip Dellinger, Alan M Fein, William A Knaus, Roland MH Schein, and William J Sibbald. 1992. Definitions for sepsis and organ failure and guidelines for the use of innovative therapies in sepsis. *Chest* 101, 6 (1992), 1644–1655.
- [3] Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the Properties of Neural Machine Translation: Encoder-Decoder Approaches.
- [4] Edward Choi, Mohammad Taha Bahadori, Le Song, et al. 2017. GRAM: Graph-based Attention Model for Healthcare Representation Learning. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- [5] Edward Choi, Mohammad Taha Bahadori, Jimeng Sun, et al. 2016. RETAIN: An Interpretable Predictive Model for Healthcare using Reverse Time Attention Mechanism. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems*.
- [6] Shaika Chowdhury, Yongbin Chen, Andrew Wen, Xiao Ma, Qiying Dai, Yue Yu, Sunyang Fu, Xiaoqian Jiang, and Nansu Zong. 2023. Predicting physiological response in heart failure management: A graph representation learning approach using electronic health records. *medRxiv* (2023).
- [7] Maia Dorsett, Melissa Kroll, Clark S Smith, Phillip Asaro, Stephen Y Liang, and Hanwan P Moy. 2017. qSOFA has poor sensitivity for prehospital identification of severe sepsis and septic shock. *Prehospital emergency care* 21, 4 (2017), 489–497.
- [8] Yongrui Duan, Jiazhen Huo, Mingzhou Chen, et al. 2023. Early prediction of sepsis using double fusion of deep features and handcrafted features. *Applied Intelligence* 53, 14 (2023), 17903–17919.
- [9] Mikhail A Dziadzko, Ognjen Gajic, Brian W Pickering, and Vitaly Herasevich. 2016. Clinical calculators in hospital medicine: availability, classification, and needs. *Computer methods and programs in biomedicine* 133 (2016), 1–6.
- [10] Evangelos J Giamarellos-Bourboulis, Thomas Tsaganos, I Tsangaris, M Lada, C Routsis, D Sinapidis, M Koupetori, M Bristianou, G Adamis, K Mandragos, et al. 2017. Validation of the new Sepsis-3 definitions: proposal for improvement in early risk identification. *Clinical Microbiology and Infection* 23, 2 (2017), 104–109.
- [11] Tim A Green, Stevan Whitt, Jeffery L Belden, Sanda Erdelez, and Chi-Ren Shyu. 2019. Medical calculators: prevalence, and barriers to use. *Computer Methods and Programs in Biomedicine* 179 (2019), 105002.
- [12] KDIGO Group et al. 2012. KDIGO clinical practice guideline for acute kidney injury. *Kidney Int. Suppl.* 2 (2012), 1.
- [13] Yuanfang Guan, Hongyang Li, Daiyao Yi, Dongdong Zhang, Changchang Yin, Keyu Li, and Ping Zhang. 2021. A survival model generalized to regression learning algorithms. *Nature computational science* 1, 6 (2021), 433–440.
- [14] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Computation* 8 (1997).
- [15] Md Mohaimenul Islam, Tahmina Nasrin, Bruno Andreas Walther, Chieh-Chen Wu, Hsuan-Chia Yang, and Yu-Chuan Li. 2019. Prediction of sepsis patients using machine learning approach: a meta-analysis. *Computer methods and programs in*

- biomedicine* 170 (2019), 1–9.
- [16] Qiao Jin, Zhizheng Wang, et al. 2024. AgentMD: Empowering Language Agents for Risk Prediction with Large-Scale Clinical Tool Learning. *arXiv preprint arXiv:2402.13225* (2024).
- [17] Alistair E.W. Johnson, Tom J. Pollard, Lu Shen, et al. 2016. MIMIC-III, a freely accessible critical care database. (2016).
- [18] Sundreen Asad Kamal, Changchang Yin, Buyue Qian, and Ping Zhang. 2020. An interpretable risk prediction model for healthcare with pattern attention. *BMC Medical Informatics and Decision Making* 20 (2020), 1–10.
- [19] Patrick S Kamath and W Ray Kim. 2007. The model for end-stage liver disease (MELD). *Hepatology* 45, 3 (2007), 797–805.
- [20] William A Knaus, Elizabeth A Draper, Douglas P Wagner, and Jack E Zimmerman. 1985. APACHE II: a severity of disease classification system. *Critical care medicine* 13, 10 (1985), 818–829.
- [21] Matthieu Komorowski, Leo A Celi, Omar Badawi, Anthony C Gordon, and Aldo Faisal. 2018. The artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care. *Nature medicine* 24, 11 (2018), 1716–1720.
- [22] Mitchell M Levy, Laura E Evans, and Andrew Rhodes. 2018. The surviving sepsis campaign bundle: 2018 update. *Intensive care medicine* 44 (2018), 925–928.
- [23] Sicen Liu, Tao Li, Haoyang Ding, Buzhou Tang, Xiaolong Wang, Qingcai Chen, Jun Yan, and Yi Zhou. 2020. A hybrid method of recurrent neural network and graph neural network for next-period prescription prediction. *International Journal of Machine Learning and Cybernetics* 11 (2020), 2849–2856.
- [24] Vincent Liu, Gabriel J. Escobar, et al. 2014. Hospital Deaths in Patients With Sepsis From 2 Independent Cohorts. *JAMA* 312, 1 (07 2014), 90–92.
- [25] Vincent X Liu, Vikram Fielding-Singh, John D Greene, Jennifer M Baker, Theodore J Iwashyna, Jay Bhattacharya, and Gabriel J Escobar. 2017. The timing of early antibiotics and hospital mortality in sepsis. *American journal of respiratory and critical care medicine* 196, 7 (2017), 856–863.
- [26] Chang Lu, Chandan Reddy, Prithwish Chakraborty, Samantha Kleinberg, and Yue Ning. 2021. Collaborative Graph Learning with Auxiliary Text for Temporal Event Prediction in Healthcare. In *IJCAI*. 3529–3535.
- [27] Yuan Luo, Peter Szolovits, Anand Dighe, and Jason Baron. 2018. 3D-MICE: integration of cross-sectional and longitudinal imputation for multi-analyte longitudinal clinical data. *JAMIA* 25, 6 (2018), 645–653.
- [28] Fenglong Ma, Radha Chitta, Jing Zhou, Quanzeng You, Tong Sun, and Jing Gao. 2017. Dipole: Diagnosis prediction in healthcare via attention-based bidirectional recurrent neural networks. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*. 1903–1911.
- [29] Fenglong Ma, Quanzeng You, et al. 2018. KAME: Knowledge-based Attention Model for Diagnosis Prediction in Healthcare. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management, CIKM*.
- [30] Konstantinos Makris and Loukia Spanou. 2016. Acute kidney injury: diagnostic approaches and controversies. *The Clinical Biochemist Reviews* 37, 4 (2016), 153.
- [31] Paul E Marik and Abdalsamih M Taeb. 2017. SIRS, qSOFA and new sepsis definition. *Journal of thoracic disease* 9, 4 (2017), 943.
- [32] John G O'Grady, Graeme JM Alexander, Karen M Hayllar, and Roger Williams. 1989. Early indicators of prognosis in fulminant hepatic failure. *Gastroenterology* 97, 2 (1989), 439–445.
- [33] Carly J Paoli, Mark A Reynolds, et al. 2018. Epidemiology and costs of sepsis in the United States—an analysis based on timing of diagnosis and severity level. *Critical care medicine* 46, 12 (2018), 1889.
- [34] Xia Qiu, Yu-Peng Lei, and Rui-Xi Zhou. 2023. SIRS, SOFA, qSOFA, and NEWS in the diagnosis of sepsis and prediction of adverse outcomes: a systematic review and meta-analysis. *Expert review of anti-infective therapy* 21, 8 (2023), 891–900.
- [35] Matthew A Reyna, Christopher S Josef, et al. 2019. Early prediction of sepsis from clinical data: the PhysioNet/Computing in Cardiology Challenge 2019. *Critical Care Medicine* (2019).
- [36] Mervyn Singer, Clifford S Deutschman, et al. 2016. The third international consensus definitions for sepsis and septic shock (Sepsis-3). *Jama* 315, 8 (2016), 801–810.
- [37] Gary B Smith, David R Prytherch, Paul Meredith, Paul E Schmidt, and Peter I Featherstone. 2013. The ability of the National Early Warning Score (NEWS) to discriminate patients at risk of early cardiac arrest, unanticipated intensive care unit admission, and death. *Resuscitation* 84, 4 (2013), 465–470.
- [38] Christian P Subbe, M Kruger, et al. 2001. Validation of a modified Early Warning Score in medical admissions. *Qjm* 94, 10 (2001), 521–526.
- [39] Carlos Sánchez, Orlando Pérez-Nieto, and Eder Zamarrón. 2023. Chapter 16 - Mechanical Ventilation in Sepsis. In *The Sepsis Codex*, Marcio Borges, Jorge Hidalgo, and Javier Perez-Fernandez (Eds.). Elsevier, 135–138. <https://doi.org/10.1016/B978-0-323-88271-2.00009-2>
- [40] Patrick Thorat, Jan Peppink, Ronald Driessen, et al. 2020. AmsterdamUMCdb: The First Freely Accessible European Intensive Care Database from the ESICM Data Sharing Initiative. (2020). <https://doi.org/10.1109/JBHI.2020.2995139> access: <https://www.amsterdammedicaldatascience.nl>
- [41] Andrea Tsores and Clinton A Marlar. 2019. Use of the Child Pugh score in liver disease. (2019).
- [42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [43] J L Vincent, Rui Moreno, Jukka Takala, Sheila Willatts, Arnaldo De Mendonça, Hajo Bruining, CK Reinhart, PeterM Suter, and Lambertius G Thijs. 1996. The SOFA (Sepsis-related Organ Failure Assessment) score to describe organ dysfunction/failure: On behalf of the Working Group on Sepsis-Related Problems of the European Society of Intensive Care Medicine (see contributors to the project in the appendix).
- [44] Hideo Wada, Takeshi Matsumoto, and Yoshiaki Yamashita. 2014. Diagnosis and treatment of disseminated intravascular coagulation (DIC) according to four DIC guidelines. *Journal of Intensive Care* 2 (2014), 1–8.
- [45] Chao Yan, Cheng Gao, Ximmeng Zhang, et al. 2019. Deep Imputation of Temporal Data. In *2019 IEEE International Conference on Healthcare Informatics, ICHI 2019, Xi'an, China, June 10-13, 2019*. 1–3.
- [46] Kai Yang, Yongxin Xu, Peinie Zou, Hongxin Ding, Junfeng Zhao, Yasha Wang, and Bing Xie. 2023. KerPrint: Local-Global Knowledge Graph Enhanced Diagnosis Prediction for Retrospective and Prospective Interpretations. *AAAI* 37, 4 (2023).
- [47] Changchang Yin, Ruoqi Liu, Dongdong Zhang, and Ping Zhang. 2020. Identifying sepsis subphenotypes via time-aware multi-modal auto-encoder. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*. 862–872.
- [48] Changchang Yin, Rongjian Zhao, Buyue Qian, Xin Lv, and Ping Zhang. 2019. Domain Knowledge guided deep learning with electronic health records. In *2019 IEEE International Conference on Data Mining (ICDM)*. IEEE, 738–747.
- [49] Dongdong Zhang, Changchang Yin, Katherine M Hunold, Xiaoqian Jiang, Jeffrey M Caterino, and Ping Zhang. 2021. An interpretable deep-learning model for early prediction of sepsis in the emergency department. *Patterns* 2, 2 (2021).

A APPENDIX

A.1 Important Notations

We summarize the important notations in Table 4.

Table 4: Important notations.

Notation	Description
$X \in R^{T \times k}$	The original observed EHRs.
$Y \in \{0, 1\}^T$	The ground truth for sepsis prediction.
T	The number of collections of observations.
t	The t -th collection.
k	The number of unique variables.
c	The number of unique clinical calculators.
$\delta_{i,j} \in R$	The time gap between i -th and j -th collections.
v	A specific node in graph.
r	A specific edge in graph.
$e^v \in R^d$	The embedding vector for node v .
$e^r \in R^d$	The embedding vector for edge r .
d	The dimension of embedding vectors.
$h_v \in R^d$	The node embedding (with observed value).
$h_r \in R^d$	The edge embedding (with time gaps).
$h_v^{(l)} \in R^d$	The l -th layer feature for node v .
h_t^l	t -th collection node’s final features.
$h_t^{l,i}$	i -th organ node’s final features (t -th collection).
$L(\cdot)$	The linear mapping function.
$N(v_t)$	The neighbors of node v_t .
$0 < \gamma < 1$	The coefficient for the skip connection.
$e_t^c \in R^c$	The estimated calculators in t -th collection.
$\hat{e}_t^c \in R^c$	The ground truth for e_t^c .
$M_{t,i}$	Binary variable indicating the availability of $\hat{e}_{t,i}^c$.
p_t^c	The estimated confidence for e_t^c .
y_t^c	The ground truth for p_t^c .
p_t	The predicted sepsis probability in t -th collection.
y_t	The ground truth of p_t .
$p_t^{o,i}$	The organ dysfunction probability for organ i .
$y_t^{o,i}$	The ground truth of $p_t^{o,i}$.
\mathcal{L}_*	Loss functions.

A.2 Clinical Calculators

A.2.1 SOFA Score Calculation. SOFA (Sepsis-related Organ Failure Assessment) score is widely used to describe organ dysfunction for septic patients in real-world clinical settings and is displayed to help clinicians assess the patients’ health states in our system Figure 9. We display the detailed computation method in Table 6. Each organ’s SOFA score ranges from 0 (normal) to 4 (most abnormal). The total SOFA score ranges from 0 (normal) to 24 (most abnormal). Although SOFA is a multi-organ dysfunction calculator, we can use the component scores to assess specific organ dysfunction when partial variables are missing.

A.2.2 Multi-Organ Calculators. In this study, we integrate multiple widely-used and well-validated clinical calculators related to sepsis, including SOFA [43], qSOFA [36], APACHE II [20], SIRS [2], NEWS [37], and MEWS [38]. The calculators can effectively describe the overall health status of critically ill patients (e.g., sepsis

Table 5: Single-organ clinical calculators. PLT: Platelet, Bili: Bilirubin, Enc: Encephalopathy, PT: Prothrombin Time, Fbg: Fibrinogen, DD: D-Dimer, FDPs: Fibrin Degradation Products, KB: Ketone Bodies, LCT: Lactate.

Calculators	Variables	Organ	Range
AKIN [30]	Creatinine, Urine output	Renal	0-3
KDIGO [12]	Creatinine, Urine output	Renal	0-3
KCH [32]	PT, Bili, KB, LCT, Sodium	Liver	0-20
MELD [19]	Bili, INR, Creatinine	Liver	6-40
CPS [41]	Bili, Albumin, PT, Asc, Enc	Liver	5-15
DIC [44]	PLT, PT, APTT, Fbg, DD, FDPs	Coagulation	0-12

patients) by assessing multiple organ dysfunctions. Table 7 displays the component variables for various organ function assessments.

A.2.3 Single-Organ Calculators. When the overall calculator is not applicable due to missing values, we can still use the observed variables to compute the organ-specific calculators (e.g., PaO2 and FiO2 for the respiration system in SOFA as shown in Table 6). In the six organs in Table 7, the variables for coagulation and liver systems (i.e., PLT, PT, and Bili) have relatively higher missing rates. We incorporate multiple organ-specific calculators (e.g., DIC [44] for coagulation, KCH [32], MELD[19], CPS [41] for liver) to assess the corresponding organ status. Table 5 presents the single-organ clinical calculators utilized in this work.

Note that the proposed SepsisCalc can handle various clinical calculators and can be further enhanced with the inclusion of more useful and related calculators.

A.3 Method Details

A.3.1 Model Backbone Selection. Unlike most existing sequence-representation-based clinical prediction studies [5, 28, 49] that treat EHRs as observational sequences, we represent patients’ EHRs as graphs. Modeling EHRs presents several important challenges:

- **Temporal Sequencing of Clinical Events:** The chronological order of clinical events is crucial for accurately describing a patient’s condition.
- **Interaction of Clinical Events:** Clinical events are often closely interconnected, such as the use of vasopressors to treat extremely low mean blood pressure (MBP).
- **High Missing Rate in Clinical Observations:** Many clinical variables, such as lab tests, often have high rates of missing data.

Sequence-based representation can handle the temporal information of EHRs well. However, the models [5, 28] typically require fixed-size vectors as input, which may necessitate additional operations (e.g., imputation) to address the issue of missing data in EHRs. Furthermore, the interaction between clinical events, which plays a significant role in modeling a patient’s health state, is often overlooked by most sequence-based models. Failure to address these last two challenges may result in suboptimal performance of the prediction models.

Table 6: The definition of SOFA score and its components across six organ systems. Each SOFA component score ranges from 0 (normal) to 4 (most abnormal). The total SOFA score ranges from 0 (normal) to 24 (most abnormal).

SOFA score	1	2	3	4
Respiration				
PaO ₂ /FiO ₂ , mmHg	< 400	< 300	< 200	< 100
Coagulation				
Platelets ×10 ³ /mm ³	< 150	< 100	< 50	< 20
Liver				
Bilirubin, mg/dl (μmol/l)	1.2 - 1.9 (20 - 32)	2.0 - 5.9 (33 - 101)	6.0 - 11.9 (102 - 204)	> 12.0 (> 204)
Cardiovascular				
Hypotension	MAP < 70 mmHg	Dopamine ≤ 5 or dobutamine (any dose)	Dopamine > 5 or epinephrine ≤ 0.1 or norepinephrine ≤ 0.1	Dopamine > 15 or epinephrine > 0.1 or norepinephrine > 0.1
Central nervous system (CNS)				
Glasgow Coma Score (GCS)	13 - 14	10 - 12	6 - 9	<6
Renal				
Creatinine, mg/dl (μmol/l) or urine output	1.2 - 1.9 (110 - 170)	2.0 - 3.4 (171 - 299)	3.5-4.9 (300 - 440) or < 500 ml/day	> 5.0 (> 440) or <200 ml/day

Table 7: Multi-organ clinical calculators. Temp: Temperature, RR: Respiratory Rate, HR: Heart Rate, Bili: Bilirubin, DA: Dopamine, DOB: Dobutamine, EPI: Epinephrine, NE: Norepinephrine, SS: Serum Sodium, SP: Serum Potassium, PLT: Platelets.

Calculators	Respiration	Coagulation	Liver	Cardiovascular	CNS	Renal	other	Range
SOFA [43]	PaO ₂ , FiO ₂	PLT	Bili	MAP, DOB, DA, EPI, NE	GCS	Creatinine, Urine output		0-24
qSOFA [36]	RR			SBP	GCS			0-3
SIRS [2]	RR, PaCO ₂			HR			WBC, Temp	0-4
NEWS [37]	RR, SpO ₂			HR, SBP	GCS		Temp	0-20
MEWS [38]	RR			HR, SBP	GCS		Temp	0-15
APACHE II [20]	RR	PLT, PT		MAP, HR	GCS	Creatinine, Urine output	Age, Temp, PH, SS, SP, WBC, LCT	0-71

In this study, we use temporal graphs to represent patients' EHRs. Graphs can naturally model the interaction between clinical variables. Only the observed variables and estimated calculators are included in the temporal graphs, eliminating the need for additional imputation methods and avoiding potential imputation bias. Moreover, we use directed edges between clinical observation nodes to incorporate temporal information. With this graph representation, we employ a graph neural network as the model backbone to represent patients' health states and make sepsis predictions.

A.3.2 Clinical Embedding. Value embedding. For variables, we adopt value embedding [47, 49] to map the values into vectors. Given a variable v and the observed values in the whole dataset, we sort the values and discretize the values into n ($n = 1000$) sub-ranges with equal number of observed values in each sub-range.

The variable v is embedded into a vector $e^v \in R^d$ with an embedding layer. For the the observed value for variable v within sub-range i ($1 \leq i \leq n$), we embed it into a vector $e^{v'} \in R^{2d}$:

$$\begin{aligned} e_j^{v'} &= \sin\left(\frac{i * j}{n * d}\right) \\ e_{d+j}^{v'} &= \cos\left(\frac{i * j}{n * d}\right), \end{aligned} \quad (14)$$

where $0 \leq j < d$. By concatenating e^v and $e^{v'}$, we obtain a vector containing both the variable's and its value's information. A linear layer is followed to map the concatenation vector into a new value embedding vector $\mathbf{h}_v \in R^d$.

$$\mathbf{h}_v = L([e^v; e^{v'}]), \quad (15)$$

where $L(\cdot)$ denotes a linear mapping function.

Time Embedding. In order to incorporate the elapsed time between observed values, we leverage a time embedding [47, 49] for the time gap δ :

$$\begin{aligned} e_j^\delta &= \sin\left(\frac{\delta * j}{T_m * d}\right) \\ e_{d+j}^\delta &= \cos\left(\frac{\delta * j}{T_m * d}\right), \end{aligned} \quad (16)$$

where $0 \leq j < d$, T_m denotes the maximum of time gap ($0 < \delta \leq T_m$). We combine the edge embedding e^r with the time embedding e^δ to generate $\mathbf{h}_r \in R^d$:

$$\mathbf{h}_r = L([e^r; e^\delta]) \quad (17)$$

A.4 Experiment Details

A.4.1 Variables Used for Sepsis Prediction. Following [21, 47], we use following variables to model sepsis patients’ health states: heart rate, Respiratory, Temperature, Spo2, SysBP, DiasBP, MeanBP, Glucose, Bicarbonate, WBC, Bands, C-Reactive, BUN, GCS, Urineoutput, Creatinine, Platelet, Sodium, Hemoglobin, Chloride, Lactate, INR, PTT, Magnesium, Aniongap, Hematocrit, PT, PaO2, SaO2, Bilirubin. The first 8 variables are immediately available vital signs.

A.4.2 Missing Rates of Variables. Table 8 displays the missing rates of the lab tests and calculators. Note that lab tests are usually observed once from several hours to days, so we display the 4-hour missing rates here. A missing value means the variable has not been observed for more than 4 hours.

A.4.3 Methods for Comparison. To validate the performance of the proposed framework for early sepsis risk prediction task, we compare the propose SepsisCalc to the following models:

- **Clinical calculator-based methods:** We use the widely used clinical calculators (i.e., **NEWS** [37], **MEWS** [38], **qSOFA** [36], **SIRS** [2]) to build the sepsis prediction models. A logistic regression is used to predict sepsis with the calculator scores, the component variables, and frequently observed vital signs.
- **GRU and LSTM:** GRU [3] and LSTM [14] are classical RNN based models, which both introduce various gates to improve RNN’s performance.
- **RETAIN:** The REverse Time AttentIoN model (RETAIN) [5] is the first work that tries to interpretate model’s disease risk prediction results with two attention modules. The attention modules generate weights for every medical events. The weights are helpful to analyze different events’ contributions to the output risk.
- **Dipole** [28]: Dipole employs bidirectional recurrent neural networks combined with three distinct attention mechanisms for patient visit information prediction.
- **DFSP** [8]: Double Fusion Sepsis Predictor (DFSP) is an early sepsis prediction model that uses early and late fusion techniques to improve the accuracy and robustness of sepsis prediction.
- **RGNN** [23]: RGNN is a hybrid method of RNN and GNN with RNN to represent patient status sequences and GNN to represent temporal medical event graphs like Figure 2(B).
- **GTN** [6]: GTN also represent EHRs as graphs and adopt a Transformer [42] to make clinical risk predictions.

Table 8: Missing rates of observed lab tests and multi-organ clinical calculators related to sepsis.

variable	AmsterdamUMCdb	OSUWMC	MIMIC-III
GCS	29%	50%	33%
Urine output	23%	39%	33%
Creatinine (CRT)	75%	85%	80%
Platelet (PLT)	76%	88%	82%
PTT	76%	83%	79%
PT	78%	92%	80%
INR	78%	84%	80%
Bilirubin	92%	94%	93%
Glucose (GLC)	34%	49%	36%
PaO2	86%	92%	87%
SaO2	88%	94%	89%
Hemoglobin (HMG)	56%	75%	69%
Bicarbonate (BCB)	69%	74%	67%
Lactate (LCT)	88%	90%	89%
WBC	67%	78%	69%
BUN	63%	76%	66%
Bands	99%	99%	99%
C-reactive	99%	99%	99%
Magnesium	66%	76%	69%
Aniongap (AG)	62%	78%	67%
Hematocrit (HMT)	60%	76%	64%
Chloride (CLR)	62%	70%	66%
Sodium (SDM)	55%	72%	65%
SOFA	94%	95%	94%
APACHE II	85%	92%	88%
SIRS	75%	85%	77%
NEWS	34%	55%	36%
MEWS	33%	54%	36%
qSOFA	33%	53%	35%

A.4.4 Implement Details. We implement our proposed model with Python 3.8.10 and PyTorch 1.12.1⁷. For training models, we use Adam optimizer with a mini-batch of 64 patients. The multi-modal data are projected into a 512-d space ($d = 512$). We train the proposed model on 1 GPU (TITAN RTX 6000), with a learning rate of 0.001. We randomly divide the patients in datasets into 10 sets. All the experiment results are averaged from 10-fold cross-validation, in which 7 sets are used for training every time, 1 set for validation, and 2 sets for testing. The validation sets are used to determine the best values of parameters in the training iterations. We ran the training and test phases 10 times and reported the mean and standard deviation of the metrics in section 5.

For non-graph-based models, we normalize the values of variable i as follows:

$$x^i = \frac{x^i - \text{mean}(x^i)}{\text{std}(x^i)}, \quad (18)$$

⁷<https://pytorch.org/>

Table 9: Organ dysfunction prediction results.

Method	MIMIC-III			AmsterdamUMCdb			OSUWMC		
	AUC	F1	Recall	AUC	F1	Recall	AUC	F1	Recall
SepsisCalc ^{imp}	.835	.440	.761	.840	.443	.778	.912	.470	.839
SepsisCalc ⁻ⁱ	.832	.439	.758	.831	.445	.775	.910	.468	.832
SepsisCalc ^{-d}	.841	.445	.767	.842	.445	.780	.912	.468	.835
SepsisCalc	.862	.458	.787	.865	.453	.790	.925	.485	.856

where *mean* and *std* are the mean value and standard deviation for the variable *i* on the whole dataset. Because the non-graph-based models cannot handle missing variables, we use a popular imputation method 3D-MICE [27] to impute the missing values.

A.4.5 Evaluation Metrics. For the sepsis prediction tasks, we use Area Under the Receiver Operating Characteristic Curve (AUC), F1, and Recall for evaluation metrics. For the calculator estimation tasks, we measure the models' performance with nRMSE. The nRMSE is calculated from the gap between the ground truth and prediction. Given a variable *i*, nRMSE is defined as:

$$nRMSE^i = \sqrt{\frac{\sum_j \sum_t a_t^{(j),i} (\hat{x}_t^{(j),i} - \tilde{x}_t^{(j),i})^2}{\sum_j \sum_t a_t^{(j),i}}}, \quad (19)$$

where $\hat{x}_t^{(j),i}$, $\tilde{x}_t^{(j),i}$, $a_t^{(j),i}$ indicate the ground truth, imputed value, and masking indicator for patient *j*, variable *i* in collection *t*.

A.4.6 Clinical Event Interaction. As Figure 2(C) shows, we incorporate the clinical event interaction to build the temporal graph. Following the surviving sepsis campaign bundle [22], we consider two kinds of important treatments for septic patients: vasopressors and IV fluid to prevent low blood pressure (related variables: SBP, DBP, DBP), antibiotics to treat infections (related variables: WBC, BUN, Bands, C-reactive). We also consider mechanical ventilation as the treatment for acute respiratory distress syndrome (related variables: SpO2, PaO2, SaO2, respiratory rate), which frequently co-occurs with sepsis [39]. We add the interaction relation between the treatments and related variables to the constructed temporal graph for patient health status modeling.

A.5 Additional Experiments

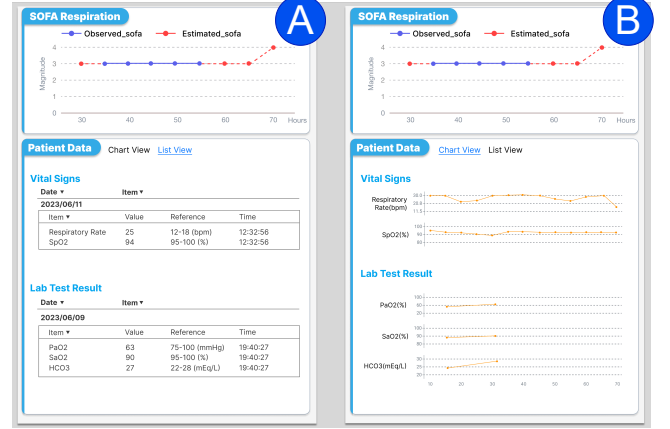
A.5.1 Organ Dysfunction Prediction. This study aims to early identify the patients with potential risk. We adopt additional prediction branches to force the model to learn the organ dysfunction patterns with \mathcal{L}_o in Equation 12.

Organ Dysfunction Prediction Setting. We use the same setting as sepsis prediction (to predict whether the specific organs will suffer from dysfunction with a 4-hour sliding window, similar to Figure 4) to predict organ dysfunction risk.

The results of organ dysfunction predictions are presented in Table 9. The findings show that SepsisCalc outperforms the other versions (*i.e.*, SepsisCalc^{imp}, SepsisCalc⁻ⁱ, SepsisCalc^{-d}). The results demonstrate the effectiveness of the proposed graph construction module in organ dysfunction prediction tasks, which could further improve the sepsis prediction performance of SepsisCalc in Table 3.

Table 10: nRMSE of calculator estimation (full observation setting). All the component variables of the calculators are observed and the ground truths of calculators are available.

Dataset	SOFA	APACHEII	qSOFA	SIRS	MEWS	NEWS
MIMIC-III	0.01	0.01	0.01	0.01	0.01	0.01
AmsterdamUMCdb	0.01	0.02	0.01	0.01	0.02	0.01
OSUWMC	0.01	0.02	0.01	0.01	0.01	0.01

**Figure 10: (A) List view of clinical variables for organ status. (B) Chart view of clinical variables for organ status.**

A.5.2 Effectiveness of Clinical Calculators. When all the component variables are observed, the ground truths of the clinical calculators are available, we use nRMSE to evaluate the calculator estimation performance. Table 10 shows that the nRMSE between the ground truths and the estimated calculators is close to 0, demonstrating that our model can accurately learn the computation mechanisms of clinical calculators.

A.5.3 Hyper-Parameter Optimization. We use the sum of four loss functions to train the model in Equation 13. We also tried to adjust the weights α_o , α_e , α_d for the loss functions. We conducted a grid search to find the best α_* with the following values [0.1, 0.2, 0.3, 0.5, 1, 2, 3, 5, 10]. We found the models achieved the best performance with $0.3 \leq \alpha_* \leq 3$ and are not sensitive to the weights, so we set $\alpha_* = 1$ for the proposed SepsisCalc.

A.6 Additional Details for Deployment

Note that our SepsisCalc system offers both list view (Figure 10(A)) and chart view (Figure 10(B)) to display the sequences of observed variables and the latest observed values with reference ranges, such that the clinicians can track the relative change trends and absolute dysfunction status for various organ systems.

Note that clinicians can interact with the system to display the dysfunction status of various organ systems. Figure 10 just displays the variables related to the respiratory system for illustration purposes.