

# Modeling Long-term Dependencies and Short-term Correlations in Patient Journey Data with Temporal Attention Networks for Health Prediction

Yuxi Liu\*

College of Science and Engineering, Flinders University  
Adelaide, South Australia, Australia  
liu1356@flinders.edu.au

Antonio Jimeno Yepes

School of Computing Technologies, RMIT University  
Melbourne, Victoria, Australia  
antonio.jose.jimeno.yepes@rmit.edu.au

Zhenhao Zhang

College of Life Sciences, Northwest A&F University  
Yangling, Shaanxi, China  
zhangzhenhow@nwfafu.edu.cn

Flora D. Salim

School of Computer Science and Engineering, UNSW  
Sydney, New South Wales, Australia  
flora.salim@unsw.edu.au

## ABSTRACT

Building models for health prediction based on Electronic Health Records (EHR) has become an active research area. EHR patient journey data consists of patient time-ordered clinical events/visits from patients. Most existing studies focus on modeling long-term dependencies between visits, without explicitly taking short-term correlations between consecutive visits into account, where irregular time intervals, incorporated as auxiliary information, are fed into health prediction models to capture latent progressive patterns of patient journeys. We present a novel deep neural network with four modules to take into account the contributions of various variables for health prediction: i) the *Stacked Attention* module strengthens the deep semantics in clinical events within each patient journey and generates visit embeddings, ii) the *Short-Term Temporal Attention* module models short-term correlations between consecutive visit embeddings while capturing the impact of time intervals within those visit embeddings, iii) the *Long-Term Temporal Attention* module models long-term dependencies between visit embeddings while capturing the impact of time intervals within those visit embeddings, iv) and finally, the *Coupled Attention* module adaptively aggregates the outputs of *Short-Term Temporal Attention* and *Long-Term Temporal Attention* modules to make health predictions. Experimental results on MIMIC-III demonstrate superior predictive accuracy of our model compared to existing state-of-the-art methods, as well as the interpretability and robustness of this approach. Furthermore, we found that modeling short-term correlations contributes to local priors generation, leading to improved predictive modeling of patient journeys.

\*Corresponding author

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 ACM 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://permissions.acm.org).

BCB '22, August 7–10, 2022, Northbrook, IL, USA

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9386-7/22/08...\$15.00

<https://doi.org/10.1145/3535508.3545535>

## CCS CONCEPTS

• **Applied computing** → **Health informatics**.

## KEYWORDS

Electronic Health Records, Deep Learning, Natural Language Processing, Patient Representation

## ACM Reference Format:

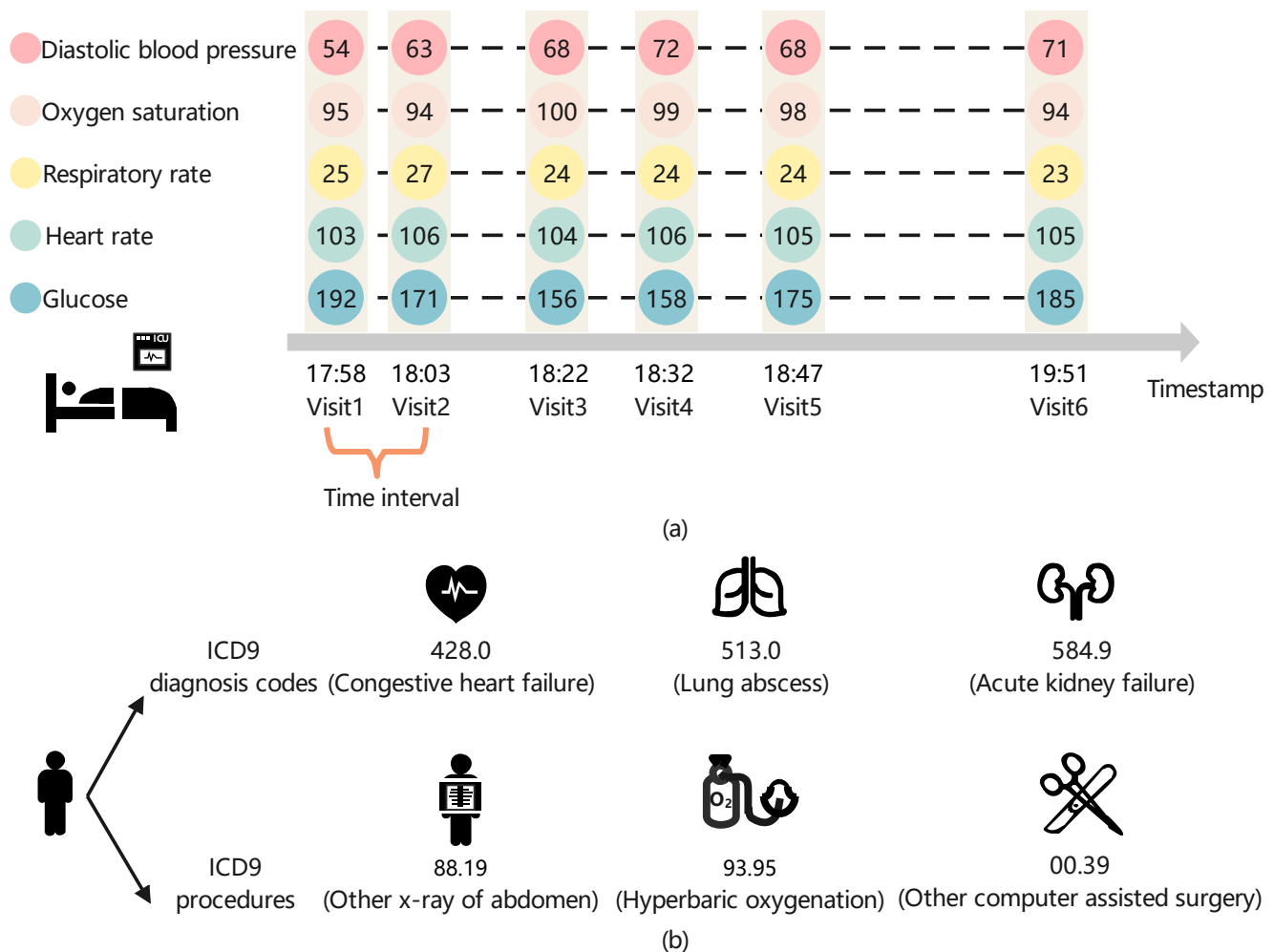
Yuxi Liu, Zhenhao Zhang, Antonio Jimeno Yepes, and Flora D. Salim. 2022. Modeling Long-term Dependencies and Short-term Correlations in Patient Journey Data with Temporal Attention Networks for Health Prediction. In *13th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics (BCB '22)*, August 7–10, 2022, Northbrook, IL, USA. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3535508.3545535>

## 1 INTRODUCTION

Developing deep learning (DL)-based models based on EHR data for predicting future health events or situations such as demands for health services has become an active research area, due to its potential to enable and facilitate care providers to take appropriate mitigating actions to minimize risks and manage demand. Examples of successful applications include early diagnosis prediction [9], disease risk prediction [26], disease progression modeling and intervention recommendation [32].

Although the existing DL-based models have achieved promising results in many health prediction tasks, significant barriers remain when applying them for modeling EHR patient journey data. Raw EHR patient journey data has its own characteristics, such as temporal, multivariate, heterogeneous, irregular, and sparse nature, etc [6]. This paper focuses mainly on the multivariate and temporal nature of EHR patient journey data, which contain inherent relationships at multiple levels and scales, notably, (i) short-term correlations between consecutive clinical visits – how every visit relates to each other in a short period and (ii) long-term dependencies between clinical visits – how each visit relates to the rest visits in the complete EHR patient journey. Integrated modeling of all these relationships is required to achieve more accurate predictions. However, to our knowledge, previous studies on EHR patient journey modeling have largely ignored the effect of (i).

The following example is based on the publicly available MIMIC-III dataset [19].



**Figure 1: Illustration of clinical time series. The time intervals between consecutive visits vary significantly.**

The physician conducts the necessary lab tests for a patient at each visit during his/her intensive care unit (ICU) stays. This usually involves multiple clinical events for a patient that occur at the same point in time or within a short period. These clinical events are strongly related to a patient's health status, particularly in lab tests associated with a series of vital sign measurements (e.g., heart rate, blood glucose, respiratory rate). As shown in Figure 1a). The short-term correlations are how every visit relates to each other in a short period. The benefit of modeling these correlations is that such a consideration can contribute to local priors generation for improving predictive modeling of patient journeys. For example, when a patient is in a certain health status (i.e., 'severe' or 'healthier' status), certain values of a vital sign (e.g., blood glucose) within neighboring visits (e.g., consecutive visits) are likely to correlate to each other, and within those visits, there is a set of vital signs that are most relevant to that health status (i.e., local priors).

The modeling of accumulated patient journey data, i.e., a set of time-ordered clinical visits associated with a series of clinical events, forms a chain of data enabling researchers and health care

providers to capture the long-term trends in patient health status. Besides, a patient's transition to a given health status also depends on his/her personal history of past clinical events, such as the previous disease diagnosis and procedures. Accordingly, both previous ICD9 diagnosis codes and procedures should be taken into consideration (As shown in Figure 1b). For example, when predicting the risk of mortality for patients, we should automatically include learning of the impact of a patient journey over the previous 48 hours on the prognosis. That is, using data from the first 48 hours of an ICU stay to predict in-hospital mortality for patients. In this sense, a predictive model should learn the context of a given patient journey through tracking and capturing the complex dynamic of clinical events and their interactions over time (i.e., long-term dependencies).

Clinical visit timestamps carry important information about the underlying patient journey dynamics, which allow us to talk about the chronologies (timing and order) of visits. In practice, it is important to consider the timing and order of visits to learn the context of a given patient's journey accurately. Usually, existing works

construct sequential models by using recurrent neural networks (RNNs), and mainly model long-term dependencies between visits in patient journey data [9, 11, 18, 32, 35].

In this paper, we propose a novel deep neural network with a modular structure, referred to as *TAttNet* (Temporal Attention Networks) hereafter, to jointly tackle the above issues. The design of *TAttNet* is inspired by ideas in processing sentences in documents from natural language processing (NLP). A patient journey is treated as a document and a clinical visit as a sentence. By learning the context of a given patient journey, *TAttNet* is able to pay attention to both critical indicative previous clinical events, even though they happened a long time ago, and local clinical events between consecutive visits. Meanwhile, *TAttNet* takes the time interval into consideration. In addition, the architecture of *TAttNet* provides interpretability of the model decisions, which is required to support its decisions, which is critical in machine learning EHR prediction systems.

We validate *TAttNet* on the mortality prediction task from a publicly available EHR dataset, on which our method outperforms the baseline attention models by large margins. Further quantitative and qualitative analysis of the learned attentions shows that our method can also provide richer interpretations that align well with the views of relevant literature and medical experts.

## 2 RELATED WORK

Recently, researchers have shown an increased interest in health prediction tasks. Recurrent neural networks such as LSTM [16], GRU [7] and GRU-D [5] have been used widely for EHR patient journey data modeling as representative deep learning models. Examples of successful applications include heart failure prediction [10, 29], and comorbidity prediction and patient similarity analysis [33]. The performance of such sequential models is limited, especially for understanding the timing of visits, because the RNN architecture only contains recurrence (i.e., the order of sequences/visits). In reality, for example, measurements (vital signs) in clinical visits are commonly acquired with irregular time intervals (as shown in Figure 1a). To take the time interval into consideration, T-LSTM [4] is proposed, which assumes that the clinical information may decay if there is a time interval between two consecutive visits. In other words, T-LSTM assumes that the more recent visits are more important than previous visits in general on health prediction tasks. T-LSTM transforms time intervals into weights and uses them to adjust the memory passed from previous moments in the Long Short Term Memory (LSTM). Based on this assumption, RetainEX [21], Timeline [2], and ATTAIN [41] are proposed, which focus on the provision of time-aware mechanisms for improving the predictive strength of RNNs. Their success makes combining RNNs with a time-aware mechanism module become the mainstream method in modeling EHR patient journey data, and they inspire the design principle of later works like ConCare [27] and HiTANet [23].

It is worth noting that the work of HiTANet [23] was with a similar motivation to the above studies but with a different emphasis and method. HiTANet emphasizes that the previous clinical information should not be decayed because latent progressive patterns of patient journeys are non-stationary, where a patient's health condition can

be better or worse at different timestamps. HiTANet designs a Time-aware Transformer to handle irregular time intervals of visits. It encodes time intervals into time vector representations and embeds them into visits first, and the Transformer model [37] then takes time vector representations as a part of the inputs and generates the hidden state representation of those visits. There is still much work to be done to achieve more adaptive time-aware mechanisms and desirable patient journey modeling simultaneously.

In addition, in real hospital scenarios, successful health prediction application requires not only the predictive strength of DL-based models but also model interpretability, which is essential to building trust in predictive models. To take the interpretability of predictive models into consideration, the innovative and seminal work of [8] developed a new DL-based model named RETAIN, which adapts two RNNs into an end-to-end interpretable network for health prediction tasks. The structure of RETAIN includes two RNNs, at visit-level and variable-level of patient contextual information respectively, to implement the attention mechanism, thereby capturing the influential past clinical visits and significant clinical variables within those clinical records. Results obtained from the RETAIN model often take an interpretable representation that is able to retain the semantics of each visit and variable while highlighting their influence on the target prediction by corresponding weights. Motivated by the successful application of RETAIN, attention-based DL-models are widely used in health prediction tasks. Examples of downstream applications include diagnosis prediction [25], mortality prediction [34], and disease progression modeling [2].

Compared to previous research, our work has a similar purpose but with a different emphasis and method. Our purpose is to construct an end-to-end interpretable network by combining deep neural networks and attention mechanisms together and mainly model EHR patient journey data. Our work emphasizes the use of *Short-Term Temporal Attention* and *Long-Term Temporal Attention* modules to separately model the short-term correlations and long-term dependencies in patient journey data. Compared to these prior works [8, 22, 25, 27, 40], this provides richer patient journey representations, which are generated from long-term dependencies between visits and short-term correlations between consecutive visits, both promoting accurate modeling of patient journeys.

By rethinking our proposed method, the intuitions behind *Short-Term Temporal Attention* and *Long-Term Temporal Attention* modules can be seen as capitalizing on the main strength of a convolutional neural network (CNN) and a multilayer perceptron (MLP). On the one hand, *Long-Term Temporal Attention* module builds upon an MLP to effectively model long-term dependencies between clinical visits. On the other hand, *Short-Term Temporal Attention* module utilizes a CNN directly to fill in the gap of MLP in capturing short-term correlations between consecutive clinical visits. Another notable advantage of *Short-Term Temporal Attention* module is that an interpretable attention map is given after the training, which shows the weights learned adaptively for variables that can describe how the importance of previous values of variables varies over time intervals and understand which variables are paid more attention to, during each visit for the target prediction. It is also worth noting that the generated two types of representations based on the above modules are simultaneously tuned together by designing a *Coupled Attention* module to make health predictions.

Furthermore, it is worth bearing in mind that our *TAttNet* starts from a *Stacked Attention* module, and it strengthens the deep semantics in clinical events within each patient journey and generates visit embeddings. One notable advantage of this module is that an interpretable attention map is obtained after the training, which gives valuable information about the target variables on how much they are correlated to each other. Moreover, since it can learn asymmetric relationships, the attention map tells us which variables precede the others in terms of multivariate clinical time series forecasting.

### 3 METHODS

In this section, we introduce the proposed *TAttNet* by discussing the basic notations of health prediction tasks first. Next, we detail four important modules of *TAttNet*. Finally, we present how to use *TAttNet* for health prediction tasks.

#### 3.1 Basic Notations

**3.1.1 EHR Patient Journeys.** Let  $P$  denote the set of all patient journeys, and  $|P|$  is the number of patient journeys. Each patient journey  $p \in P$  can be defined as follows:

$$r = \begin{Bmatrix} r_{11} & \cdots & r_{1T} \\ \vdots & \ddots & \vdots \\ r_{N1} & \cdots & r_{NT} \end{Bmatrix}, r^c = \{0, 1\}^{g_c}, r^d = \{0, 1\}^{g_d}, \quad (1)$$

where  $r$  denotes clinical time series, which consists of the number of  $T$  records and  $N$  is the number of time-variant features in each record. Both  $r^c$  (previous ICD9 diagnosis codes) and  $r^d$  (ICD9 procedures) are time-invariant features.  $g_c$  and  $g_d$  are the number of unique ICD9 diagnosis codes and unique ICD9 procedures, respectively.

**3.1.2 The Temporal Nature of EHR Patient Journeys.** Let  $\mu_i$  and  $\mu_j$  denote the timestamp of records  $r_i$  and  $r_j$ .  $\delta_{ij}$  denotes the time interval between records  $r_i$  and  $r_j$ , i.e.,  $\delta_{ij} = \mu_j - \mu_i$  ( $i < j$ ).

#### 3.2 TAttNet architecture

The architecture of the proposed *TAttNet* is shown in Figure 2. In the following subsections, we provide the inner design of *TAttNet* in greater detail.

**3.2.1 Stacked Attention Module.** In this subsection, we develop a *Stacked Attention* module to strengthen the deep semantics in clinical events within each patient journey and generate visit embeddings (see Figure 2a). This module is mainly composed of the multi-head attention mechanism [37] whereby it contains multiple self-attention layers working in a parallel pattern. A self-attention layer includes three elements: a set of key-query pairs and values. The key-query pairs are used to compute the inner dependency weights, then used to update the values. Compared to a single self-attention layer, one advantage of the multi-head attention mechanism is that it enhances the attention layer with multiple representation subspaces. After obtaining the refined representation of each position by the multi-head attention mechanism, we add a feed-forward network (FFN) sub-layer.

Based on  $r$ , the generated representations  $r^{MHA}$  can be defined as follows:

$$r^{MHA} = \text{MultiHeadAttention}(r) = W_o[head_1(r) \oplus head_2(r) \oplus \cdots \oplus head_m(r)], \quad (2)$$

where  $head_m$  and  $\oplus$  are the  $m$ -th attention head and concatenation operation, respectively.  $W_o \in \mathbb{R}^{N \times mN}$  is a learnable parameter. For simplicity, we use  $r^{MHA}$  to denote  $\text{MultiHeadAttention}(r)$ . In the following, we provide the details of all implementations.

First, we compute the attention weight ( $\xi$ ) in Eq. (3) that determines how much each feature will be expressed at this certain feature.

$$\xi_1, \xi_2, \dots, \xi_N = \text{Softmax}\left(\frac{Q^\top \cdot K_1}{\sqrt{d_K}}, \frac{Q^\top \cdot K_2}{\sqrt{d_K}}, \dots, \frac{Q^\top \cdot K_N}{\sqrt{d_K}}\right), \quad (3)$$

where  $Q$  and  $K_i$  denote query vector and key vector respectively,  $d_K$  is the dimension of  $K_i$ , and the scaled dot product is used as the attention function [37].  $V_i$  denotes value vector. We use the way of [38] to project input vectors into query, key and value spaces to obtain  $Q$ ,  $K_i$  and  $V_i$  separately. The formula can be defined as follows:

$$Q, K_i, V_i = W_Q \cdot r^\top, W_K \cdot r_i^\top, W_V \cdot r_i^\top, \quad (4)$$

where  $W_Q, W_K \in \mathbb{R}^{d_K \times T}$  and  $W_V \in \mathbb{R}^{T \times T}$  are the projection matrices, which are learnable parameters. Each  $head_m$  has its own projection matrix.

Second, we use  $v_i$  and  $\xi_i$  to compute  $head_m$ , which is obtained by using the weighted sum of  $V_i$ . The weights are calculated by applying the attention function to all key-query pairs.

$$head_m(r) = \sum_{i=1}^N \xi_i \cdot V_i^\top. \quad (5)$$

Based on Eq. (3-5), each  $head_m$  is obtained by letting  $r$  attend to all the feature positions so that any feature interdependencies between  $r$  and  $r_i$  can be captured.

After obtaining the refined new representation of each position by the multi-head attention mechanism, we use the way of [37] to add a feed-forward network sub-layer, which is used to transform the features non-linearly. The formula is defined as follows:  $FFN(r^{MHA}) = W_2 \cdot \max(W_1 \cdot r^{MHA} + b_1, 0) + b_2$ . We also employ a residual connection [14] around each of the two sub-layers, followed by layer normalization [1]. For simplicity, we use  $h$  to denote the output of  $FFN(r^{MHA})$ . Last, the output of this module is  $h \in \mathbb{R}^{N \times T}$ .

**3.2.2 Learning Time Feature Embedding.** Before implementing *Short-Term Temporal Attention* and *Long-Term Temporal Attention* modules, we develop a Time-Encoder to embed  $\delta$  and  $h$  into the same feature space. The formula can be defined as below:

$$\delta_{enc} = W_\delta \cdot \delta + b_\delta, \quad (6)$$

where  $W_\delta \in \mathbb{R}^{N \times 1}$  and  $b_\delta \in \mathbb{R}^{N \times 1}$  are learnable parameters.

**3.2.3 Short-Term Temporal Attention Module.** In this subsection, we develop a *Short-Term Temporal Attention* module to model short-term correlations between consecutive visit embeddings while capturing the impact of time intervals within those visit embeddings

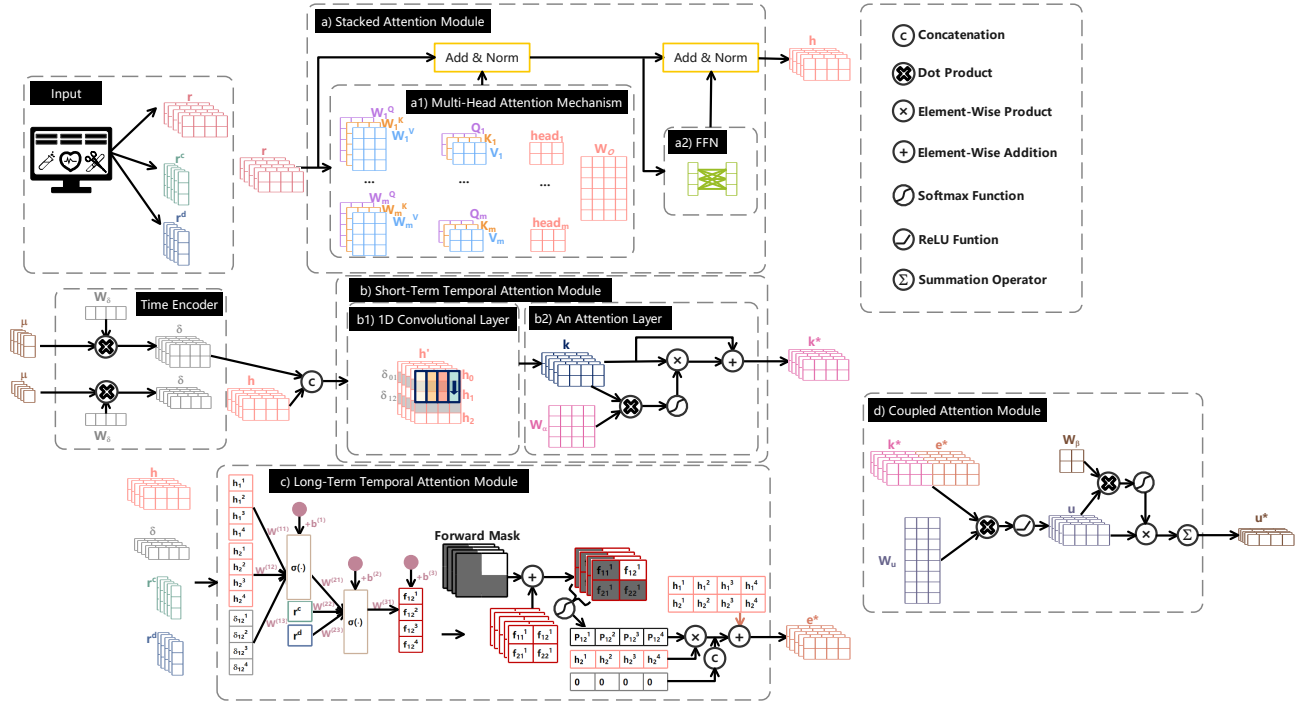


Figure 2: Architecture of the proposed TAttNet.

(see Figure 2b). The module consists of a customized 1D convolutional layer (i.e., kernel size=3, stride=2) and an attention layer. One notable advantage of the module is that an interpretable attention map is given after the training, which shows the weights learned adaptively for variables (clinical events) that can describe how the importance of previous values of variables varies over time intervals.

First, we embed time intervals into visit embeddings. Hence,  $h = [h_1, h_2, \dots, h_T]$  is expanded into  $h' = [h_0, \delta_{01}, h_1, \delta_{12}, h_2, \delta_{23}, h_3, \dots, h_{T-1}, \delta_{(T-1)T}, h_T] \in \mathbb{R}^{N \times (2T+1)}$ . Note that  $h_0$  is a zero vector. Second, we apply 1D convolution operation only over the horizontal dimension. To take time intervals into consideration, we specifically use a combination of  $N$  kernels  $\{W^{(j)}\}_{j=1}^N$ , and each kernel has kernel size 3 and stride 2. E.g.,  $h'_{t:t+2}$  denotes the concatenation of visit embeddings  $h_{t-1}, h_t$  and time intervals  $\delta_{(t-1)t}$  from  $h'_t$  to  $h'_{t+2}$ . A kernel  $W^{(j)} \in \mathbb{R}^{3 \times 1}$  is naturally applied on the window of  $h'_{t:t+2}$  to produce a new feature  $k_{(t+1)/2}^{(j)} \in \mathbb{R}$  with the ReLU activation function as follows:

$$k_{(t+1)/2}^{(j)} = \text{ReLU}(h'_{t:t+2} \cdot W^{(j)} + b^{(j)}), \quad (7)$$

where  $b^{(j)} \in \mathbb{R}$  is a bias term and  $\text{ReLU}(k) = \max(k, 0)$ . This kernel is applied to each possible window of values in the whole description  $\{h'_{1:3}, h'_{3:5}, \dots, h'_{2T-1:2T+1}\}$  to generate a feature map  $k^{(j)} \in \mathbb{R}^T$  as follows:  $k^{(j)} = [k_1^{(j)}, k_2^{(j)}, \dots, k_T^{(j)}]$ ,  $j = 1, 2, \dots, N$ . We can see that each kernel produces a feature. Since we have  $N$  kernels, the final vector representation of a patient journey can be obtained by

concatenating all the extracted features, i.e.,  $k \in \mathbb{R}^{N \times T}$ . Last, an attention layer is applied to produce an attention weight  $\alpha \in \mathbb{R}^{N \times T}$  and the final feature representation  $k^*$  is obtained as follows:

$$\begin{aligned} \alpha &= \text{softmax}(W_\alpha \cdot k + b_\alpha), \\ k^* &= \alpha \odot k + k, \end{aligned} \quad (8)$$

where  $W_\alpha \in \mathbb{R}^{N \times N}$  and  $b_\alpha \in \mathbb{R}^{N \times 1}$  are learnable parameters.  $\odot$  denotes element-wise multiplication. The generated attention weight  $\alpha \in \mathbb{R}^{N \times T}$  encodes the causative and associative relationships of clinical events within each clinical visit embedding.

**3.2.4 Long-Term Temporal Attention Module.** In this subsection, we develop a Long-Term Temporal Attention module to model long-term dependencies between visit embeddings while capturing the impact of time intervals within those visit embeddings (see Figure 2c). This module is specifically designed to learn the contextual information from each patient journey (i.e., the number of  $T$  visit embeddings), whereby it generates an attention weight matrix  $P_{:,j} \in \mathbb{R}^{N \times T}$  ( $j = 1, 2, \dots, T$ ) that encodes the causative and associative relationships between the clinical events of  $j$ -th visit embedding and other visit embeddings. Mathematically, the module can be formalized as follow:

$$\begin{aligned} f(h_i, h_j, \delta_{ij}) &= \sigma(W^{(11)}h_i + W^{(12)}h_j + W^{(13)}\delta_{ij} + b^{(1)}), \\ f(h_i, h_j, \delta_{ij}, r^c, r^d) &= W^{(31)}\sigma(W^{(21)}f(h_i, h_j, \delta_{ij}) + W^{(22)}r^c \\ &\quad + W^{(23)}r^d + b^{(2)}) + b^{(3)}, \\ f'(h_i, h_j, \delta_{ij}, r^c, r^d) &= f(h_i, h_j, \delta_{ij}, r^c, r^d) + M_{ij}^{fw}, \end{aligned} \quad (9)$$

where  $h_i$  and  $h_j$  denote the  $i$ -th and  $j$ -th visit embeddings of  $h$ . All  $W$  and  $b$  are learnable parameters.  $\sigma(\cdot)$  is an activation function.

According to Eq. (9), we can see that the proposed module first computes  $h_i$ ,  $h_j$  and  $\delta_{ij}$  to obtain the  $f(h_i, h_j, \delta_{ij})$ . It then builds the relationship between the obtained  $f(h_i, h_j, \delta_{ij})$  and  $r^c$ ,  $r^d$  and generates the  $f(h_i, h_j, \delta_{ij}, r^c, r^d)$ . Last, it incorporates a forward mask  $M_{ij}^{fw}$  into the  $f(h_i, h_j, \delta_{ij}, r^c, r^d)$  and generates the  $f'(h_i, h_j, \delta_{ij}, r^c, r^d)$ .

The benefit of incorporating the mask operation [31] (i.e.,  $M_{ij}^{fw}$ ) is that such a consideration can help the module learn the context of a patient journey according to the order of their clinical visits. That is, when computing each visit, the module only pays attention to all its previous visits. Specifically, a forward mask  $M^{fw}$  is incorporated as follows:

$$M_{ij}^{fw} = \begin{cases} 0, & i < j \\ -\infty, & \text{otherwise.} \end{cases} \quad (10)$$

In order to generate an attention weight matrix (i.e.,  $P_{ij}$ ), a softmax function is applied to the  $[f'(h_i, h_j, \delta_{ij}, r^c, r^d)]_{i=1}^T$  and thus obtain the probability distribution  $p(g|h, h_j)$ .  $g$  is an indicator of which feature is important to  $h_j$ . A large  $P_{ij}^l$ ,  $P_{ij}^l \triangleq p(g_l = i|h, h_j)$  means that  $h_i$  contributes important information to  $h_j$  on the  $l$ -th feature dimension. For simplicity, we ignore the subscript  $l$  if it does not cause any confusion.

The  $e_j$  is obtained by traversing all  $h_i$  as follows:

$$e_j = \sum_{i=1}^T P_{ij} \odot h_i, j = 1, 2, \dots, T, \quad (11)$$

where  $e_j$  is the weighted average of sampling a visit embedding according to its importance. We obtain  $e = [e_1, e_2, \dots, e_T] \in \mathbb{R}^{N \times T}$  from Eq. (11). We add  $e$  to  $h$  to obtain the final feature representation  $e^*$ , i.e.,  $e^* = e + h$ .

**3.2.5 Coupled Attention Module.** Through the aforementioned design, *TAttNet* has been able to model long-term dependencies and short-term correlations in patient journey data. Meanwhile, it takes the time interval into consideration. The two types of representations ( $k^*$  and  $e^*$ ), generated based on *Short-Term Temporal Attention* and *Long-Term Temporal Attention* modules, should be simultaneously tuned together. In this subsection, we develop a *Coupled Attention* module to adaptively couple both representations (see Figure 2d). Given the obtained representations  $k^*$  and  $e^*$ , an overall representation  $u$  is generated as follows:

$$u = \text{ReLU}(W_u \cdot [k_t^*, e_t^*] + b_u), \quad (12)$$

where  $W_u \in \mathbb{R}^{d_u \times 2N}$  and  $b_u \in \mathbb{R}^{d_u}$  are learnable parameters.  $d_u$  is the dimension of the overall representation  $u$ . The rectified linear unit is defined as  $\text{ReLU}(x) = \max(x, 0)$ . Note that  $\max()$  applies element-wise to vectors.

Next, an attention layer is applied to  $u$  to obtain a final representation  $u^*$  as follows:

$$\beta = \text{softmax}(u \cdot W_\beta + b_\beta), u^* = \sum_{i=1}^T \beta_i \odot u_i, \quad (13)$$

where  $W_\beta \in \mathbb{R}^{T \times T}$  and  $b_\beta \in \mathbb{R}^T$  are learnable parameters.  $\odot$  denotes element-wise multiplication.

### 3.3 Health Prediction

The output of subsection 4.5, a final representation  $\{u_p^*\}_{p=1}^{|P|}$ , is fed into a Softmax output layer to obtain the prediction probability  $\hat{y}_p$ :

$$\hat{y}_p = \text{softmax}(W_y \cdot u_p^* + b_y), \quad (14)$$

where  $W_y$  and  $b_y$  are learnable parameters. The cross-entropy between the ground truth  $y_p$  and the prediction probability  $\hat{y}_p$  is used to calculate the loss. Thus, the objective function of health prediction is the average of cross entropy:

$$\mathcal{L}(\theta) = -\frac{1}{|P|} \sum_{p=1}^{|P|} (y_p^\top \cdot \log(\hat{y}_p) + (1 - y_p)^\top \cdot \log(1 - \hat{y}_p)), \quad (15)$$

where  $\theta$  is the parameter of *TAttNet* and  $|P|$  is the total number of patient journeys.

## 4 EXPERIMENTS

In this section, we report the experimental results on a publicly available EHR dataset to demonstrate the effectiveness of the proposed method. After the discussion of the experimental setting, we first compare our approach with attention-based DL methods. Moreover, we analyze the effectiveness and interpretability of our approach by an ablation study and case studies.

### 4.1 Experimental Setup

**4.1.1 Datasets and Tasks.** We validate the performance of our method on health risk prediction tasks from the publicly available Medical Information Mart for Intensive Care (MIMIC-III) dataset [19], for both the prediction accuracy and prediction robustness. The data were normalized on the basis of the literature [13].

MIMIC-III is one of the largest publicly available ICU datasets, comprising 38,597 distinct patients and a total of 53,423 ICU stays. We utilize clinical times series data (e.g., heart rate, respiration rate) and previous ICD9 diagnosis codes as well as ICD9 procedures as inputs. The prediction tasks here are two binary classification tasks, 1) in-hospital mortality (48 hours after ICU admission): to evaluate ICU mortality based on the data from the first 48 hours after ICU admission, 2) in-hospital mortality (complete hospital stay after ICU admission): to evaluate ICU mortality based on the data during the complete hospital stay.

**4.1.2 Baselines.** To fairly evaluate the effectiveness of the proposed *TAttNet*, we implement the following attention-based DL methods:

- (1)  $\text{GRU}_\alpha$  is the basic GRU [7] with a location-based attention mechanism [25].
- (2) RETAIN [8] (NeurIPS 2016) is an interpretable deep learning model for health prediction. It purposes the use of a two-level attention mechanism (based on two RNNs), which could enhance both the performance and interpretability of the model.
- (3) Dipole [25] (SIGKDD 2017) has a similar purpose with RETAIN, but with a different method. It uses a bidirectional RNN and incorporates a three-level attention mechanism.

- (4) Transformer<sub>e</sub> is the encoder of the Transformer [37] (NeurIPS 2017); in the final step, we use to flatten and feed-forward networks to make health predictions.
- (5) INPREM [40] (SIGKDD 2020) can be seen as a combination of a linear part and a non-linear part. It mainly uses three attention mechanisms to fill in the gap of a linear model in learning the non-linear relationship of features from EHR patient journey data.
- (6) T-LSTM [4] (SIGKDD 2017) is an improved LSTM based approach by modifying gate information to model the time decay.
- (7) RetainEX [21] is built upon RETAIN to learn weights for clinical visits and events while capturing the impact of time intervals within those visits. It mainly handles previous clinical visits by enabling time decay.

**4.1.3 Implementation Details & Evaluation Strategies.** We perform all the baselines and *TAttNet* with Python v3.7.0. For each task, we randomly split the datasets into training, validation, and testing sets in a 75:10:15 ratio. The validation set is used to select the best values of parameters. Binary outcomes were evaluated with the area under the receiver operating characteristic curve (AUROC) and the area under the precision-recall curve (AUPRC). We repeat all the approaches ten times and report the average performance.

We use class weight in CrossEntropyLoss for a highly imbalanced dataset. This is achieved by placing an argument called 'weight' on the CrossEntropyLoss function. This approach offers an effective way of alleviating the problem of being highly imbalanced.

In order to evaluate the interpretability of our *TAttNet*, the following steps were taken: (i) we randomly select a case (i.e., one patient journey) from the MIMIC-III dataset [19]. (ii) we input imputed values for the missing items of this case on the basis of literature [13]. (iii) we use the case as the input of the trained *TAttNet*, and the mask operation is applied to both *Stacked Attention* and *Short-Term Temporal Attention* modules. Altogether, this makes the dimensionality of the data input to *TAttNet* appropriate and avoids taking imputed items into consideration in the target prediction.

## 4.2 Performance Evaluation

Table 1 shows the performance of all methods on the MIMIC-III dataset. The results indicate that *TAttNet* significantly and consistently outperforms other baseline methods. E.g., for the mortality prediction of MIMIC-III (48 hours after ICU admission), *TAttNet* achieves the highest AUROC with 0.9076, one standard deviation of 0.002. Similarly, *TAttNet* achieved the highest AUPRC with 0.6326, one standard deviation of 0.016.

We find that all methods demonstrated good prediction robustness for lengthy EHR patient journeys. E.g., for the mortality prediction of MIMIC-III (complete hospital stay after ICU admission), RETAIN and Dipole achieved AUROCs of 0.8663 and 0.8884, respectively, which is an improvement over the prediction performance of their on 48 hours by roughly 2%. In contrast, *TAttNet* achieves the best AUROC/AUPRC scores amounting to 0.9538 and 0.7999.

## 4.3 Ablation Study

We now need to examine the effectiveness of different modules of our method. To this end, we conduct an ablation study on the

datasets. To determine whether the developed modules improve the prediction performance, we present four variants of *TAttNet* as follows:

*TAttNet<sub>α</sub>*: *TAttNet* without *Stacked Attention*.

*TAttNet<sub>β</sub>*: *TAttNet* without *Short-Term Temporal Attention*.

*TAttNet<sub>γ</sub>*: *TAttNet* without *Long-Term Temporal Attention*.

*TAttNet<sub>δ</sub>*: *TAttNet* without *Coupled Attention*.

We present the ablation study results in Table 2. We find that *TAttNet* outperforms *TAttNet<sub>β</sub>* (i.e., without the *Short-Term Temporal Attention* module). It indicates that modeling of short-term correlations between consecutive visits contributes to local priors generation, for improving predictive modeling of patient journeys. *TAttNet* also outperforms *TAttNet<sub>δ</sub>* (i.e., without the *Coupled Attention* module), which demonstrates that adaptively aggregates two types of representations are generated by *Short-Term Temporal Attention* and *Long-Term Temporal Attention* modules which can produce enhanced prediction performance. The superior performance of *TAttNet* than the *TAttNet<sub>α</sub>* verifies the efficacy of the *Stacked Attention* module which can generate richer patient journey representations and improve the performance.

## 4.4 Case Study: Method Interpretability

We validate the interpretability of our *TAttNet* with random examples selected from the datasets, which are presented in Figures 3 and 4. The proposed *TAttNet* enjoys good interpretability owing to both *Stacked Attention* and *Short-Term Temporal Attention* modules. Specifically, the *Stacked Attention* module is able to provide an interpretable attention map after the training, which gives valuable information about the target variables on how much they are correlated to each other. The *Short-Term Temporal Attention* module is able to provide an interpretable attention map after the training, which shows the weights learned adaptively for variables that can describe how the importance of previous values of variables varies over time intervals and understand which variables are paid more attention to during each visit for the target prediction.

Figure 3 shows the variable/feature importance of patient A who died after 48 hours of an ICU stay. The feature weights ranging from 0.0 ~ 1.0 were calculated by the proposed *Stacked Attention* module. The ordinates y-axis of the figure shows the Query features, and the abscissas x-axis shows the Key features. The boxes intensity of color in the shaded areas in Figure 3 shows how much each Key feature responds to the Query when a Query feature makes a query. Looking at Figure 3, it is apparent that most of the features related to the primary disease are more likely to respond to each other, represented by the feature weights of two matrices (See Figure 3 for details of the primary disease of patient A).

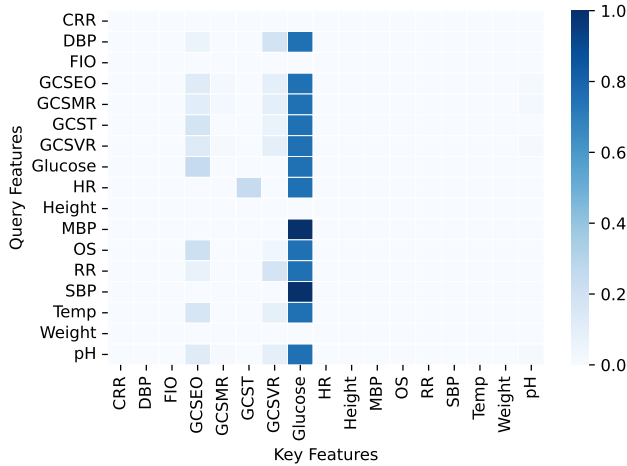
It is common medical knowledge that the glucose of a patient is strongly related to diabetes. As shown in Figure 3, in the box of Glucose-Glucose position, *TAttNet* pays much more attention to the glucose in patient A. *TAttNet* determines that there are relatively high associations between pH and glucose. This is highly consistent with the medical literature [36, 39]. Moreover, *TAttNet* determines that there are relatively high associations between diastolic blood pressure (DBP), systolic blood pressure (SBP), and glucose. This result may be explained by the fact that patient A has both essential hypertension and diabetes. This finding broadly supports the

**Table 1: The AUROC/AUPRC scores of in-hospital mortality prediction with different observation windows for MIMIC-III dataset.**

MIMIC-III/Mortality Prediction		48 hours after ICU admission		complete hospital stay after ICU admission	
Metrics		AUROC	AUPRC	AUROC	AUPRC
Methods	GRU $_{\alpha}$	0.8339(0.005)	0.6248(0.006)	0.8406(0.008)	0.6213(0.010)
	RETAIN	0.8429(0.003)	0.4848(0.008)	0.8663(0.010)	0.6014(0.013)
	Dipole	0.8685(0.002)	0.5357(0.004)	0.8884(0.005)	0.5180(0.007)
	Transformer $_e$	0.8116(0.008)	0.6046(0.010)	0.8260(0.014)	0.5923(0.017)
	INPREM	0.8699(0.008)	0.5207(0.012)	0.8807(0.008)	0.5399(0.018)
	T-LSTM	0.8743(0.002)	0.4354(0.024)	0.9153(0.010)	0.7367(0.031)
	RetainEX	0.8737(0.005)	0.4668(0.016)	0.9240(0.005)	0.7496(0.017)
	<i>TAttNet</i>	<b>0.9076(0.002)</b>	<b>0.6326(0.016)</b>	<b>0.9538(0.003)</b>	<b>0.7999(0.015)</b>

**Table 2: Ablation performance comparison.**

MIMIC-III/Mortality Prediction		48 hours after ICU admission		complete hospital stay after ICU admission	
Metrics		AUROC	AUPRC	AUROC	AUPRC
Methods	<i>TAttNet</i> $_{\alpha}$	0.9005(0.002)	0.6173(0.018)	0.9109(0.002)	0.6406(0.018)
	<i>TAttNet</i> $_{\beta}$	0.9036(0.002)	0.6002(0.017)	0.9360(0.006)	0.7389(0.027)
	<i>TAttNet</i> $_{\gamma}$	0.8886(0.008)	0.6228(0.010)	0.9376(0.003)	0.7511(0.011)
	<i>TAttNet</i> $_{\delta}$	0.9064(0.002)	0.6243(0.010)	0.9491(0.003)	0.7930(0.011)
	<i>TAttNet</i>	0.9076(0.002)	0.6326(0.016)	0.9538(0.003)	0.7999(0.015)

**Figure 3: The variable/feature importance of patient A with respect to the other variables/features. In patient A, diabetes mellitus is the primary disease, accompanied by a series of diabetes-induced diseases and other diseases.**

work of other studies in this area linking essential hypertension with diabetes [24, 28, 42]. This finding also supports evidence from clinical observations [15] that there are relatively high associations between abnormal glucose, blood pressure, and mortality. Furthermore, *TAttNet* determines that there are relatively high associations between four Glasgow Coma Scale scores (i.e., GCSEO, GCSMR,

GCST, GCSVR) and glucose. These are highly consistent with the medical literature [20].

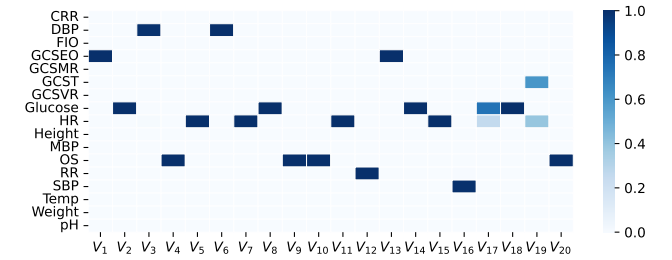
**Figure 4: The variable/feature importance of patient A with respect to each one of the visits (V1 to V20).**

Figure 4 shows the variable/feature importance of patient A who died after 48 hours of an ICU stay (Note that both Figures 3 and 4 have used the same patient A). The feature weights ranging from 0.0 ~ 1.0 were calculated by the proposed *Short-Term Temporal Attention* module. The boxes intensity of color in the shaded areas in Figure 4 shows how much each feature responds to the target prediction when the module calculates a visit. Note that the module considers all visits of patient A (i.e., an entire journey of patient A), but the images are understandably truncated for visibility. We take the first 20 visits of patient A as an example for detailed discussion.

In Figure 4, there are observed changes in the importance of previous values of variables vary over time intervals (i.e., at different visits). E.g., in the box of Glucose- $V_2$  and Glucose- $V_8$  positions,



*TAttNet* pays much more attention to the glucose in patient A. After zooming on Visits 2 and 8 ( $V_2$ ,  $V_8$ ), we find that glucose had greater weight, contributing most to the corresponding visit. When querying the data of patient A from the MIMIC-III dataset, we found that glucose had a larger value of 608 mg/dL within Visit 2. Similarly, glucose had a larger value of 440 mg/dL within Visit 8. In reviewing the literature [12], we found that the commonly-used standard of glucose in critically ill patients is below 180 mg/dL. Overall, these results provide reasonably consistent evidence of an association between abnormal glucose and mortality [3, 17].

## 5 LIMITATIONS AND FUTURE WORKS

Several limitations to this study need to be acknowledged. (i) MIMIC-III database is a large medical dataset comprising all information relating to patients admitted to intensive care units. This work uses clinical time series data and previous ICD9 diagnosis codes as well as ICD9 procedures as inputs. Further investigation and experimentation into more informative details such as admission information and free text diagnosis are strongly recommended. (ii) In our case studies, the scope of the experimental data was chosen based on the literature [30], i.e., the first 48 hours of an ICU stay and complete hospital stay after ICU admission. In terms of directions for future research, further work could explore patient-specific contextual information at various time points, such as 24 or 36 hours. It would be worthwhile to compare patient-specific contextual information at various time points, and the experimental results may offer different prediction probabilities and medical findings. (iii) The lack of model uncertainty analysis in the results adds further caution regarding the generalisability of these findings. This paper argues that successful health prediction requires high predictive accuracy and higher-quality model interpretability. However, DL-based models are usually highly parameterized and thus may increase the uncertainty of their prediction results. In order to achieve high predictive accuracy, researchers often make many efforts to find the best parameters. Hence, the model results often depend on the customization of parameters, but the parameters are uncertain, achieved by analyzing different perspectives and using different optimal algorithms, even human experiences. Future studies should address the questions raised by the model parameter uncertainty (also known as model uncertainty). Altogether, this will build more trust in predictive models and thus enable the transition from academic research to clinical applications.

## 6 CONCLUSIONS

In this paper, we introduce a novel deep neural network named *TAttNet* for health prediction tasks in EHR patient journey data mining. This prediction will enable and facilitate care providers to take appropriate mitigating actions to minimize risks and manage demand. We conduct experiments on a publicly available EHR dataset. *TAttNet* yielded performance improvements in both AUROC and AUPRC over the state-of-the-art prediction methods. Furthermore, our *TAttNet* also provides an interpretation as to the cause of individual patient predictions, allowing the user to take more important features into consideration in future investigations. The authors believe this may be useful for providing personalized estimates of outcome probabilities. A number of possible future studies using

the same experimental setup are apparent such as hospital stays, re-admissions, and future diagnoses.

## REFERENCES

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450* (2016).
- [2] Tian Bai, Shanshan Zhang, Brian L Egleston, and Slobodan Vucetic. 2018. Interpretable representation learning for healthcare via capturing disease progression through time. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 43–51.
- [3] Elizabeth LM Barr, Paul Z Zimmet, Timothy A Welborn, Damien Jolley, Dianna J Magliano, David W Dunstan, Adrian J Cameron, Terry Dwyer, Hugh R Taylor, Andrew M Tonkin, et al. 2007. Risk of cardiovascular and all-cause mortality in individuals with diabetes mellitus, impaired fasting glucose, and impaired glucose tolerance: the Australian Diabetes, Obesity, and Lifestyle Study (AusDiab). *Circulation* 116, 2 (2007), 151–157.
- [4] Inci M Baytas, Cao Xiao, Xi Zhang, Fei Wang, Anil K Jain, and Jiayu Zhou. 2017. Patient subtyping via time-aware LSTM networks. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*. 65–74.
- [5] Zhengping Che, Sanjay Purushotham, Kyunghyun Cho, David Sontag, and Yan Liu. 2018. Recurrent neural networks for multivariate time series with missing values. *Scientific reports* 8, 1 (2018), 1–12.
- [6] Yu Cheng, Fei Wang, Ping Zhang, and Jianying Hu. 2016. Risk prediction with electronic health records: A deep learning approach. In *Proceedings of the 2016 SIAM International Conference on Data Mining*. SIAM, 432–440.
- [7] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078* (2014).
- [8] Edward Choi, Mohammad Taha Bahadori, Joshua A Kulas, Andy Schuetz, Walter F Stewart, and Jimeng Sun. 2016. RETAIN: an interpretable predictive model for healthcare using reverse time attention mechanism. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*. 3512–3520.
- [9] Edward Choi, Mohammad Taha Bahadori, Andy Schuetz, Walter F Stewart, and Jimeng Sun. 2016. Doctor ai: Predicting clinical events via recurrent neural networks. In *Machine learning for healthcare conference*. PMLR, 301–318.
- [10] Edward Choi, Andy Schuetz, Walter F Stewart, and Jimeng Sun. 2017. Using recurrent neural network models for early detection of heart failure onset. *Journal of the American Medical Informatics Association* 24, 2 (2017), 361–370.
- [11] Cristóbal Esteban, Oliver Staack, Stephan Baier, Yinchong Yang, and Volker Tresp. 2016. Predicting clinical events by combining static and dynamic information using recurrent neural networks. In *2016 IEEE International Conference on Healthcare Informatics (ICHI)*. IEEE, 93–101.
- [12] Jan Gunst, Astrid De Bruyn, and Greet Van den Berghe. 2019. Glucose control in the ICU. *Current opinion in anaesthesiology* 32, 2 (2019), 156.
- [13] Hrayr Harutyunyan, Hrant Khachatryan, David C Kale, Greg Ver Steeg, and Aram Galstyan. 2019. Multitask learning and benchmarking with clinical time series data. *Scientific data* 6, 1 (2019), 1–18.
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [15] Patrick Henry, Frédérique Thomas, Athanasios Benetos, and Louis Guize. 2002. Impaired fasting glucose, blood pressure and cardiovascular disease mortality. *Hypertension* 40, 4 (2002), 458–463.
- [16] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [17] Osnat Itzhaki Ben Zadok, Ran Kornowski, Ilan Goldenberg, Robert Klempfner, Yoel Toledano, Yitschak Biton, Enrique Z Fisman, Alexander Tenenbaum, Gregory Golovchiner, Ehud Kadmon, et al. 2017. Admission blood glucose and 10-year mortality among patients with or without pre-existing diabetes mellitus hospitalized with heart failure. *Cardiovascular Diabetology* 16, 1 (2017), 1–9.
- [18] Abhyuday N Jagannatha and Hong Yu. 2016. Structured prediction models for RNN based sequence labeling in clinical text. In *Proceedings of the conference on empirical methods in natural language processing*. conference on empirical methods in natural language processing, Vol. 2016. NIH Public Access, 856.
- [19] Alistair EW Johnson, Tom J Pollard, Lu Shen, H Lehman Li-Wei, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. MIMIC-III, a freely accessible critical care database. *Scientific data* 3, 1 (2016), 1–9.
- [20] Atsushi Kotera, Shinsuke Iwashita, Hiroki Irie, Junichi Taniguchi, Shunji Kasaoka, and Yoshihiro Kinoshita. 2014. An analysis of the relationship between Glasgow Coma Scale score and plasma glucose level according to the severity of hypoglycemia. *Journal of Intensive Care* 2, 1 (2014), 1–6.
- [21] Bum Chul Kwon, Min-Je Choi, Joanne Taery Kim, Edward Choi, Young Bin Kim, Soonwook Kwon, Jimeng Sun, and Jaegul Choo. 2018. Retainvis: Visual analytics

- with interpretable and interactive recurrent neural networks on electronic medical records. *IEEE transactions on visualization and computer graphics* 25, 1 (2018), 299–309.
- [22] Wonsung Lee, Sungrae Park, Weonyoung Joo, and Il-Chul Moon. 2018. Diagnosis prediction via medical context attention networks using deep generative modeling. In *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 1104–1109.
- [23] Junyu Luo, Muchao Ye, Cao Xiao, and Fenglong Ma. 2020. Hitanet: Hierarchical time-aware attention networks for risk prediction on electronic health records. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 647–656.
- [24] Yaogai Lv, Yan Yao, Junsen Ye, Xin Guo, Jing Dou, Li Shen, Anning Zhang, Zhiqiang Xue, Yaqin Yu, and Lina Jin. 2018. Association of blood pressure with fasting blood glucose levels in Northeast China: a cross-sectional study. *Scientific reports* 8, 1 (2018), 1–7.
- [25] 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.
- [26] Fenglong Ma, Jing Gao, Qiuling Suo, Quanzeng You, Jing Zhou, and Aidong Zhang. 2018. Risk prediction on electronic health records with prior medical knowledge. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1910–1919.
- [27] Liantao Ma, Chaohe Zhang, Yasha Wang, Wenjie Ruan, Jiangtao Wang, Wen Tang, Xinyu Ma, Xin Gao, and Junyi Gao. 2020. Concare: Personalized clinical feature embedding via capturing the healthcare context. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 833–840.
- [28] G Mancina. 2005. The association of hypertension and diabetes: prevalence, cardiovascular risk and protection by blood pressure reduction. *Acta Diabetologica* 42, 1 (2005), s17–s25.
- [29] G Maragatham and Shobana Devi. 2019. LSTM model for prediction of heart failure in big data. *Journal of medical systems* 43, 5 (2019), 1–13.
- [30] Yilmazcan Ozyurt, Mathias Kraus, Tobias Hatt, and Stefan Feuerriegel. 2021. AttDMM: an attentive deep Markov model for risk scoring in intensive care units. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 3452–3462.
- [31] Xueping Peng, Guodong Long, Tao Shen, Sen Wang, and Jing Jiang. 2020. Self-attention enhanced patient journey understanding in healthcare system. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 719–735.
- [32] Trang Pham, Truyen Tran, Dinh Phung, and Svetha Venkatesh. 2017. Predicting healthcare trajectories from medical records: A deep learning approach. *Journal of biomedical informatics* 69 (2017), 218–229.
- [33] Tong Ruan, Liqi Lei, Yangming Zhou, Jie Zhai, Le Zhang, Ping He, and Ju Gao. 2019. Representation learning for clinical time series prediction tasks in electronic health records. *BMC medical informatics and decision making* 19, 8 (2019), 1–14.
- [34] Ying Sha and May D Wang. 2017. Interpretable predictions of clinical outcomes with an attention-based recurrent neural network. In *Proceedings of the 8th ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*. 233–240.
- [35] Qiuling Suo, Fenglong Ma, Giovanni Canino, Jing Gao, Aidong Zhang, Pierangelo Veltri, and Gnasso Agostino. 2017. A multi-task framework for monitoring health conditions via attention-based recurrent neural networks. In *AMIA annual symposium proceedings*, Vol. 2017. American Medical Informatics Association, 1665.
- [36] Zuping Tang, Xiaogu Du, Richard F Louie, and Gerald J Kost. 2000. Effects of pH on glucose measurements with handheld glucose meters and a portable glucose analyzer for point-of-care testing. *Archives of pathology & laboratory medicine* 124, 4 (2000), 577–582.
- [37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [38] Zhiwei Wang, Yao Ma, Zitao Liu, and Jiliang Tang. 2019. R-transformer: Recurrent neural network enhanced transformer. *arXiv preprint arXiv:1907.05572* (2019).
- [39] Sakiko Yoshida, Teruki Miyake, Shin Yamamoto, Shinya Furukawa, Tetsuji Niya, Hidenori Senba, Sayaka Kanzaki, Osamu Yoshida, Toru Ishihara, Mitsuhito Koizumi, et al. 2018. Relationship between urine pH and abnormal glucose tolerance in a community-based study. *Journal of diabetes investigation* 9, 4 (2018), 769–775.
- [40] Xianli Zhang, Buyue Qian, Shilei Cao, Yang Li, Hang Chen, Yefeng Zheng, and Ian Davidson. 2020. INPREM: An Interpretable and Trustworthy Predictive Model for Healthcare. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 450–460.
- [41] Yuan Zhang. 2019. ATTAIN: Attention-based Time-Aware LSTM Networks for Disease Progression Modeling.. In *In Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI-2019)*, pp. 4369–4375, Macao, China.
- [42] Honglei Zhao, Fanfang Zeng, Xiang Wang, and Lili Wang. 2017. Prevalence, risk factors, and prognostic significance of masked hypertension in diabetic patients. *Medicine* 96, 43 (2017).