

# MeSIN: Multilevel Selective and Interactive Network for Medication Recommendation

Yang An<sup>a</sup>, Liang Zhang<sup>b</sup>, Mao You<sup>c</sup>, Xueqing Tian<sup>c</sup>, Bo Jin<sup>d,\*</sup> and Xiaopeng Wei<sup>a,\*</sup>

<sup>a</sup>School of Computer Science and Technology, Dalian University of Technology, Dalian, China, 116024

<sup>b</sup>International Business College, Dongbei University of Finance and Economics, Dalian, China, 116025

<sup>c</sup>Department of Health Technology Assessment, China National Health Development Research Center, Beijing, 100001, PR China

<sup>d</sup>School of Innovation and Entrepreneurship, Dalian University of Technology, Dalian, China, 116024

## ARTICLE INFO

### Keywords:

Intelligent healthcare management  
Medication recommendation  
Multilevel interactive learning  
Temporal event modelling

## ABSTRACT

Recommending medications for patients using electronic health records (EHRs) is a crucial data mining task for an intelligent healthcare system. It can assist doctors in making clinical decisions more efficiently. However, the inherent complexity of the EHR data renders it as a challenging task: (1) *Multilevel structures*: the EHR data typically contains multilevel structures which are closely related with the decision-making pathways, e.g., laboratory results lead to disease diagnoses, and then contribute to the prescribed medications; (2) *Multiple sequences interactions*: multiple sequences in EHR data are usually closely correlated with each other; (3) *Abundant noise*: lots of task-unrelated features or noise information within EHR data generally result in suboptimal performance. To tackle the above challenges, we propose a multilevel selective and interactive network (MeSIN) for medication recommendation. Specifically, MeSIN is designed with three components. First, an attentional selective module (ASM) is applied to assign flexible attention scores to different medical codes embeddings by their relevance to the recommended medications in every admission. Second, we incorporate a novel interactive long-short term memory network (InLSTM) to reinforce the interactions of multilevel medical sequences in EHR data with the help of the calibrated memory-augmented cell and an enhanced input gate. Finally, we employ a global selective fusion module (GSFM) to infuse the multi-sourced information embeddings into final patient representations for medications recommendation. To validate our method, extensive experiments have been conducted on a real-world clinical dataset. The results demonstrate a consistent superiority of our framework over several baselines and testify the effectiveness of our proposed approach.

## 1. Introduction

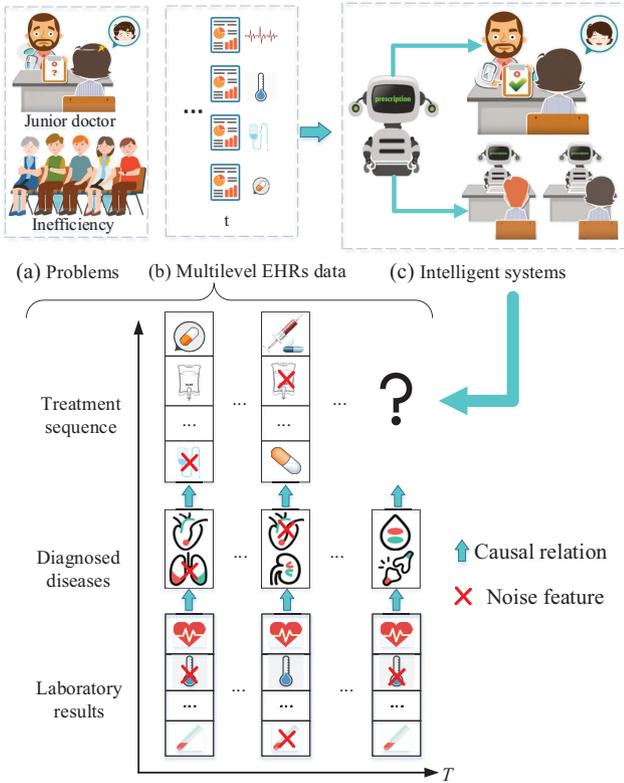
Recently, healthcare intelligence has become a hot research topic, which is mainly due to the following factors: 1) the wide utilization of digital healthcare systems that produce huge valuable data such as electronic health records (EHRs); 2) the tremendous advancements of computational models, in particular the deep learning methods; 3) an urgent need of intelligent healthcare systems to assist the junior doctors and solve the inefficiency of medical resources (Figure 1 (a)), brought by the emergent public health incidents, such as COVID-19 or Coronavirus Pandemic [14]. One of the core EHR-based applications is recommending medications for patients, with the aim to assist or even replace doctors in making effective and safe medication prescriptions for certain patients, as shown in Figure 1 (c).

However, recommending medications for patients is a challenging task due to the complexity of EHR data. As illustrated in Figure 1 (b), this complexity can be attributed to several factors. First, the EHR data typically comprises of multilevel medical records including three key aspects, e.g., laboratory results, diagnosed diseases, and prescribed treatment medications. Within each visit, the multilevel structure is closely related with the decision-making pathway, which is a kind of hierarchical structure. As shown in Figure 1 (b), the hierarchy generally begins with the laboratory re-

sults that precisely record the detailed health progression of a patient, the middle is the diseases diagnosed by doctors according to corresponding laboratory results, and the top is the medications prescribed by doctors after comprehensive decision-making processes. Thus, how to fully leverage the inherent multilevel structural information has become a critical factor for modeling the intelligent medication recommendation systems. Most existing medication recommendation studies [31, 39, 15] put more efforts on modeling the mapping relations between diagnosed diseases and recommended medications. Though these algorithms have achieved early success on the medication recommendation task, they often over-emphasize the visit-level temporal dependency, and overlook the critical influence of such a hierarchy shown in Fig.1.

Second, along with the temporal dependencies of multiple medical sequences, the complex sequential correlations embodied in the multilevel structure of EHR data (Figure 1 (b)) is another challenge for a medication recommendation task. For example, the laboratory results can provide enough hints for certain diseases, e.g., when the anion gap is abnormally high, creatinine is abnormally low, and aspartate is abnormally high, it demonstrates some pulmonary related diseases including respiratory infections, pulmonary emphysema, and the patient needs corresponding treatments such as glucose, sodium bicarbonate, xylitol and budesonide. Such phenomenon clearly indicates multilevel correlations of EHR data. While most existing methods [31, 39, 22] overlook

✉ jinbo@dlut.edu.cn (B. Jin); xpwei@dlut.edu.cn (X. Wei)  
ORCID(s): 0000-0002-4094-7499 (B. Jin)



**Figure 1:** The urgent need of developing intelligent system with multilevel EHRs data, and corresponding complexity.

such important relations of medical sequences and only consider the temporal dependencies. Though LSTM-DE [18] and RAHM [2] model the interactions of two sequences, they only considered the effects of related input sequences on the memory cell state while neglecting the influence on the current input state. Thus, in this paper, we consider infusing the interactions of two sequences both on the memory cell state and input cell state simultaneously into the temporal sequence learning network.

Third, unlike the above discussed structure-related limitations, how to recognize and filter out noisy information existing in EHR data at each timestamp is another important challenge that inhibits the recommendation performance. However, few deep learning studies in health informatics focus on infusing the feature selections into the learning process except LSAN [35] which considers assigning flexible attention weights to different diagnosis codes via their relevance to corresponding diseases for reducing the effect of irrelevant diagnosis codes in EHR data. However, the doctors practically pay more attention on the critical few factors and neglect the irrelevant medical indicators or historical medical codes. In other words, irrelevant features should be deleted in the decision-making process, and unimportant historical medical codes should be given less attention. In this way, the general attention mechanism might not be appropriate in the learning process.

To address the aforementioned challenges, in this paper, we develop a Multilevel Selective and Interactive Network, called MeSIN. The key idea lies in three aspects. First, a

multilevel learning framework is designed to encode the the inherent multilevel structure of EHR data, which imitates the decision-making process of doctors in hospitals. Second, to capture the intra-correlations of multiple visits within each medical sequence and the inter-correlations of multiple sequences of EHR data, we propose a novel interactive temporal sequence learning network. Third, due to the multiple heterogeneous inputs including medical codes embeddings and learned laboratory results embeddings, we introduce multiple attentional selective modules into the framework to make automatic and intelligent selections. Therefore, our developed framework MeSIN consists of three key components including the attentional selective module (ASM), the interactive long-short term memory network (InLSTM), and the global selective fusion module (GSFM). In MeSIN, they tightly work together and significantly enhance each other for medication recommendation.

The main contributions of this study are as follows:

- **Multilevel Selective and Interactive Network (MeSIN).** To the best of our knowledge, MeSIN is the first to formulate medications recommendation task as a multilevel learning framework, which is a challenging process in clinical decision-making systems. It can fully leverage the inherent multilevel structure of EHR data to learn a comprehensive patient representation for reasonable medication recommendation.
- **Interactive Long-Short Term Memory Network (InLSTM).** InLSTM can effectively reinforce the interactions of multiple temporal heterogeneous sequences with the help of a recurrent neural structure, a new calibrated memory-augmented cell and a novel enhanced input gate.
- **Attentional Selective Module (ASM).** We incorporate multiple improved attentional selective modules into MeSIN, which can intelligently assign relevance scores to the learned medical codes embeddings according to their importance with recommended medications.
- **Global Selective Fusion Module (GSFM).** We design a self-attention based global selective fusion module (GSFM) to effectively infuse the obtained heterogeneous embeddings into patient representation according to their respective importance and minimize the adverse effects induced by the irrelevant information.

## 2. Related works

Related studies in healthcare informatics are reviewed from the following three perspectives: medication recommendation, attention mechanism in health informatics, sequence modeling in health informatics.

### 2.1. Medication recommendation

Recently, artificial intelligence, particularly computational intelligence and machine learning methods and algorithms,

has been naturally applied in the development of recommender systems to improve prediction accuracy [36]. Recommending rational and effective medications in time for patients, as a paramount recommendation task in health domain, has attracted great amount of studies. Shang et al. [31] categorized medication recommendation-related tasks into instance-based and longitudinal sequential recommendation methods. Instance-based methods are based on the current disease progression of patients. For example, Zhang et al. [38] formulated the medications recommendation task as a sequential decision-making problem and leveraged a recurrent decoder to model label dependency. Wang et al. [34] addressed the recommendation issues by casting the task as an order-free Markov decision process (MDP) problem. However, they all ignored valuable historical information. Until now, longitudinal sequential recommendation methods mainly consider the impact of historical medical records by modeling their temporal dependencies. For instance, Jin et al. [19] developed three different LSTMs to model heterogeneous data interactions for predicting the next-period prescriptions. Shang et al. [31] incorporated historical diseases and procedure codes, as well as medication records, in their model. Shang et al. [30] considered hierarchical knowledge about diagnoses and medications to enhance the code representation for medication recommendation. An et al. [2] formulated the medication prediction task as hierarchical multi-task learning framework for improving the interpretability of predicted results. However, few of them simultaneously consider all the multiple heterogeneous sequences and the correlations between them in the decision-making of medications recommendation.

## 2.2. Attention mechanism in health informatics

The attention mechanism has been proposed to automatically assign importance scores according to the information relevance. In this case, larger weights indicate that the corresponding vectors are more relevant to generating the output. Due to its powerful ability, the attention mechanism has been widely used in various neural network based applications such as language understanding tasks [11],[28], computer vision problems [32],[17]. Likewise, attention mechanism in health informatics has been prevalent in predictive modelling. For instance, GRAM [7], KAME [25], and G-BERT [30] leveraged the attention mechanism to integrate domain knowledge into disease or medication code representations for better performance. Retain [8], Dipole [24], Timeline [3] and LSAN [35] all introduced attention mechanism to model the disease progression by considering the dependencies among visits and provide some interpretable insights. In addition, GCT [10] were equipped with advanced attention networks, i.e., Transformer [33], to build the correlations between medical codes from every visits based on the automatically learned attention weights. Likely, AMANet [15] utilized multiple attention networks including self-attention and inter-attention to capture the intra-view interaction and inter-view interaction. However, the attention mechanism used in above models all generated the dense attention weights

without zero weights value, which means that they cannot filter out the noise information and attend focus on the critical aspects.

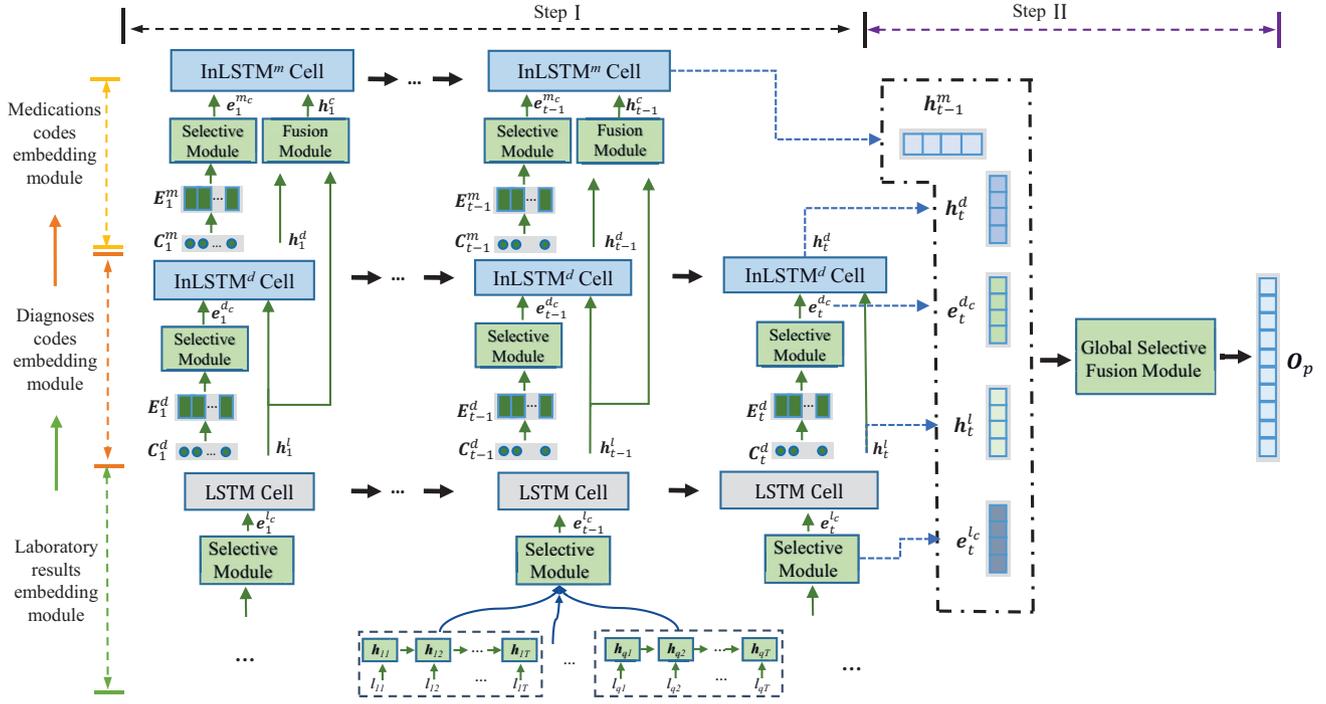
## 2.3. Sequence modeling in health informatics

Due to the complexity of clinical scenarios, as shown in Fig.1, EHR systems in hospitals accumulate complex temporal and heterogeneous sequences. Existing studies in health informatics have widely utilize the temporal sequential records from EHRs to solve healthcare problems such as predicting disease progression [9], [29], [40], [35], medications recommendation [38], [31], [18] and clinical trial recruitment [5], [37]. However, most of the studies such as T-LSTM [4], MNN [29] and LSAN [35] mainly focused on modelling the temporal dependencies of multiple visits of homogeneous sequence such as the history diseases sequence. While the medication recommendation task involves multiple temporal and heterogeneous sequences, not only the temporal intra-dependencies but also the inter-correlations between the sequences should be considered when modelling the sequence learning process. Though GAMENet [31] utilized two medical sequences to model the temporal dependencies for medications recommendation, it didn't consider the correlations of sequences. DMNC [22] presented a two-view sequential learning model to model the complex interactions. However, the complex differentiable neural computer (DNC) blocks used in DMNC [22] do not explicitly model sequential interactions. In contrast, Jin et al. [18] developed three heterogeneous LSTM models to model the correlations of different types of medical sequences by connecting hidden neurons, but neglect the impact on patient's current status. MiME [9] modelled the inherent multilevel structures of medical codes by incorporating the relationships between the diagnoses and their treatments into patient visit representations. AMANet [15] utilized multiple attention networks to capture the intra- and inter- view interactions of heterogeneous and temporal sequences, but overlook the multilevel nature of EHR data.

## 3. Methods

### 3.1. Problem definitions

To facilitate the latter introduction of our computational methods and generalize the applicable dataset, we define the data from electronic health record system (EHR) using the mathematical symbols as follows. The longitudinal EHR data contains a large number of patient records, and each patient can be represented as a sequence of multivariate observations:  $\mathcal{P} = \{\mathcal{X}^1, \mathcal{X}^2, \dots, \mathcal{X}^{t_n}\}$  over time, where  $n \in \{1, 2, \dots, N\}$ ,  $N$  is total number of patients, and  $t_n$  is the number of visits for the  $n$ -th patient. Without loss of generality, we will describe the model for a patient and the subscript (n) will be dropped whenever it is unambiguous. Each visit  $\mathcal{X}^t$  consists of sequential laboratory indicator-wise results  $\mathcal{L}^t = \{[l_{11}, \dots, l_{1T}], \dots, [l_{q1}, \dots, l_{qT}]\}$ , where  $l_{qT}$  denotes the  $q$ -th indicator result at  $T$ -th timestamp within  $t$ -th visit, and categorized data including  $C_d^t \subset C_d$  (a union set of diagnoses codes) and  $C_m^t \subset C_m$  (a union set of medications



**Figure 2:** The architecture of MeSIN. Overall, from the bottom up, MeSIN comprises of three hierarchically correlative modules and a global fusion module. The laboratory results embedding module first projects the temporal sequence of each laboratory indicator  $\{l_{q1}, \dots, l_{qT}\}$  into embedded vectors  $\{h_{1T}^l, \dots, h_{qT}^l\}$  via a multi-channel GRU, and then uses the attentional selective module to compute the enhanced embedding  $e_t^{lc}$  which would be input into LSTM to obtain the visit-level embedding  $h_t^l$ ; The diagnoses codes embedding module and medications codes embedding module respectively contain three substructures: the medical codes embedding layer for mapping the medical codes set  $C_*^t$  into dense embeddings set  $E_*^t$ , the attentional selective module for computing the enhanced embedding  $e_*^{t,c}$ , and the interactive LSTM (InLSTM) for calculating the visit-level embedding  $h_t^*$ . Finally, the global selective fusion module is used to infuse the learned multi-sourced embeddings into patient representation  $O_p$  for medication recommendation.

codes). For simplicity, we use  $C_*^t$  to represent the unified definition of medical codes.  $C_*$  denotes the medical code set and  $|C_*|$  denotes the size of medical code set.  $c_*^j$  is the  $j^{\text{th}}$  medical code in  $C_*$ .

**Problem Definition 1 (Medication recommendation).** Given the historical visit records of a patient  $\{\mathcal{X}^1, \mathcal{X}^2, \dots, \mathcal{X}^{t-1}\}$ , the current laboratory results  $\mathcal{L}^t$  and the diagnosed diseases  $C_d^t$ , our goal is to recommend reasonable medications by generating the multi-label output  $\hat{y}_t^m \in \{0, 1\}^{|C_m^t|}$ :

$$\hat{y}_t^m = f(\{\mathcal{X}^1, \mathcal{X}^2, \dots, \mathcal{X}^{t-1}\}, \mathcal{L}^t, C_d^t). \quad (1)$$

### 3.2. Multilevel selective and interactive network

We propose a novel architecture, MeSIN, to implement the medications recommendation task. As shown in Fig.2, MeSIN is a multilevel learning framework, which mainly consists of two steps. In step I, the hierarchical historical information embeddings learning process begins with the laboratory results embedding learning module, followed by the diagnoses codes embedding module, and then the medications codes embedding module. In step II, the global selective fusion module is utilized to infuse the learned heterogeneous embeddings into the patient representation according

to the selective weights.

#### 3.2.1. Laboratory sequence embedding module

As shown in Fig.2, the module mainly consists of three key parts: a multi-channel time-series embedding layer, an attentional selective layer, and a temporal sequence learning network.

**Multi-channel time-series embedding layer.** As the meanings of particular clinical features for patients in diverse medical conditions are different, the progression of laboratory indicators are distinct accordingly. Thus, the embedding of sequential feature representing the indicator changing progress is distinct from each other. Here, inspired by ConCare [26], we employ the multi-channel time-series embedding layer to embed the sequence of each laboratory indicator feature separately by multi-channel GRUs:

$$h_{qT}^l = \text{GRU}_q(l_{q1}, l_{q2}, \dots, l_{qT}), \quad (2)$$

where  $\{l_{q1}, l_{q2}, \dots, l_{qT}\}$  denotes the time series and  $h_{qT}^l$  represents the embedded vectors of feature q. Therefore, all the embedded vectors of time series of indicator features  $\{h_{1T}^l, \dots, h_{qT}^l\}$  can be acquired in the same way. To reduce

clutter, the superscript (t) representing the results generated at  $t$ -th visit will be dropped whenever it is unambiguous.

**ASM for laboratory results embeddings selection.** For each sequence of laboratory results, we gain the corresponding embeddings  $\mathbf{h}_{1T}^l, \dots, \mathbf{h}_{qT}^l$ . However, in clinical scenarios, doctors pay attention to only a few paramount indicators according to clinical experience, which can effectively improve work efficiency.

In this case, the attention mechanism that computes the attention weights using **softmax function** [6] might be inappropriate, because it results in dense attention alignments that is wasteful and making models less interpretable. Therefore, we introduce a sparse attention using **entmax** [27] computing attention weights for ASM of MeSIN to increase focus on relevant source medical codes embeddings and make the model more interpretable. Here, we employ the specially proposed ASM to compute the enhanced laboratory results embedding  $\mathbf{e}_i^{lc}$ :

$$\begin{aligned} \mathbf{e}_i^{lc} &= \sum_{i=1}^q \alpha_i^l \mathbf{h}_{iT}^l (\alpha_i^l \in \boldsymbol{\alpha}^l), \\ \boldsymbol{\alpha}^l &= \alpha - \text{entmax}([\alpha_1^l, \alpha_2^l, \dots, \alpha_q^l], \gamma_l), \\ \alpha_i^l &= \tanh(\mathbf{W}_{l_a}^\top \mathbf{h}_{iT}^l + \mathbf{b}_{l_a}), \end{aligned} \quad (3)$$

where  $\mathbf{W}_{l_a} \in \mathbb{R}^d$  and  $\mathbf{b}_{l_a} \in \mathbb{R}$  are the parameters of ASM to be learned.  $\alpha - \text{entmax}$  [27] is a special method of entmax, by which we can find the optimal equilibrium point by controlling the value of  $\gamma_l$ . For  $\gamma_l > 1$ , as the value increases, entmax tends to produce sparse probability distributions, yielding a function family interpolating between softmax and sparsimax. In this way, we can compute all the enhanced laboratory results embeddings  $\{\mathbf{e}_1^{lc}, \dots, \mathbf{e}_t^{lc}\}$  at each timestamp.

**Temporal sequence learning network.** To further capture the temporal dependency of multi-visit laboratory results, the enhanced laboratory results embeddings  $\{\mathbf{e}_1^{lc}, \dots, \mathbf{e}_t^{lc}\}$  will be input into the temporal sequence learning network for combining with the historical laboratory results:

$$\mathbf{h}_t^l = \text{LSTM}_L(\mathbf{h}_{t-1}^l, \mathbf{e}_t^{lc}), \quad (4)$$

where  $\text{LSTM}_L$  represents the long-short temporal neural network (LSTM) for capturing the temporal dependency of laboratory examination sequence,  $\mathbf{h}_t^l$  denotes the obtained visit-level laboratory results embedding containing the history information at  $t$ -th visit. With the same calculations in the remaining timestamps, we can finally have all the history laboratory results embeddings  $\{\mathbf{h}_1^l, \dots, \mathbf{h}_t^l\}$  which will be input into the following embedding modules.

### 3.2.2. Diagnoses codes embedding module

After checking the laboratory results, the doctors tend to retrieve the history diagnosed diseases and combines them with current disease condition for comprehensive decision-making. Likely, as shown in Fig.2, we design a module that

contains three critical parts: a diagnosis code embedding layer, an attentional selective module and a novel temporal sequence interactive learning network.

**Diagnosis code embedding layer.** Taking the timestamp  $t$  as an example, MeSIN first encodes each diagnosis code  $\mathbf{c}_i^d$  into a dense representation vector  $\mathbf{e}_i^d \in \mathbb{R}^d$  as:

$$\mathbf{e}_i^d = \mathbf{W}_e^d \mathbf{c}_i^d, \quad (5)$$

where  $\mathbf{W}_e^d \in \mathbb{R}^{d \times |C_*|}$  is the embedding matrix of medical codes that needs to be learned,  $d$  is the size of embedding dimension, and  $|C_*|$  is the size of medical code set. Thus, for the diagnosis code set  $C_d$ , we can represent it by a collection of dense representation vectors  $\mathbf{E}^d = [\mathbf{e}_1^d, \dots, \mathbf{e}_{|C_d|}^d] \in \mathbb{R}^{d \times |C_d|}$ . Then, for the  $i$ -th visit, we can obtain dense embedding set  $\mathbf{E}_t^d = [\mathbf{e}_1^d, \dots, \mathbf{e}_m^d] \in \mathbb{R}^{d \times m}$ , in which each embedding vector is extracted from  $\mathbf{E}^d$  if it exists in  $i$ -th visit.

**ASM for diagnoses codes embeddings selection.** However, as discussed before, not every historical disease has impact on the future disease risk, we should assign different relevance scores on each code embedding according to their importance degree. Here we also leverage ASM for diagnoses codes embeddings selection. The enhanced diagnosis code embedding  $\mathbf{e}_t^{dc}$  can be calculated as:

$$\begin{aligned} \mathbf{e}_t^{dc} &= \sum_{i=1}^m \alpha_i^d \mathbf{e}_i^d (\alpha_i^d \in \boldsymbol{\alpha}^d), \\ \boldsymbol{\alpha}^d &= \alpha - \text{entmax}([\alpha_1^d, \alpha_2^d, \dots, \alpha_m^d], \gamma_d), \\ \alpha_i^d &= \tanh(\mathbf{W}_{d_a}^\top \mathbf{e}_i^d + \mathbf{b}_{d_a}), \end{aligned} \quad (6)$$

where  $\mathbf{W}_{d_a} \in \mathbb{R}^d$  and  $\mathbf{b}_{d_a} \in \mathbb{R}$  are the parameters of ASM to be learned.  $\gamma_d$  is the hyper-parameter of  $\alpha - \text{entmax}$  [27] in this module. In this way, we can compute all the enhanced diagnoses codes embeddings  $\mathbf{E}^{dc} = \{\mathbf{e}_1^{dc}, \dots, \mathbf{e}_t^{dc}\}$ .

**InLSTM in diagnosis code sequence learning.** In addition to modeling of the temporal dependency of a single sequence, we should also consider the interactions of multiple sequences in the sequence learning network. As before, the laboratory results could be regarded as critical references when making the diagnoses by doctors. Hence, the sequence of laboratory results should be used to control the diagnosed disease sequence learning process. Therefore, such kind of network will adopt two input sequences: one is the primary input  $\mathbf{x}_t^p$  of sequence learning network such as the gained diagnosis code embedding  $\mathbf{e}_t^{dc}$ , another is the auxiliary input  $\mathbf{x}_t^a$  for assisting in controlling the primary sequence learning process such as the learned visit-level laboratory results embedding  $\mathbf{h}_t^l$ .

Therefore, the basic LSTM model [16] is not appropriate under such circumstances. Inspired by LSTM-DE [19], we propose a novel interactive long-short term memory network (**InLSTM**), as shown in Fig.3, to reinforce the interaction process of two associated sequences, which brings

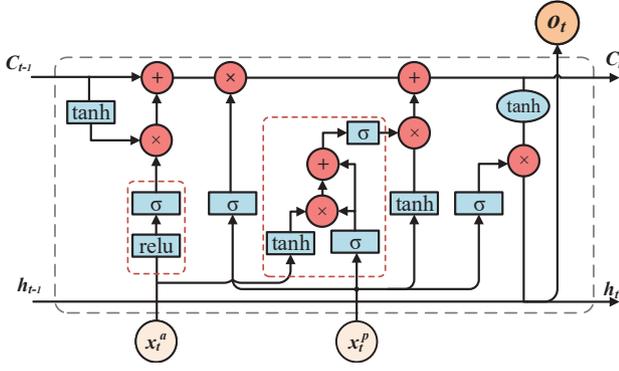


Figure 3: The structure of InLSTM.

in two novel components including a calibrated memory-augmented cell and an enhanced input gate. It can be defined as:

$$\mathbf{h}_t = \text{InLSTM}(\mathbf{h}_{t-1}, \mathbf{x}_t^p, \mathbf{x}_t^a), \quad (7)$$

where the detailed mathematical expression of InLSTM is:

$$\begin{bmatrix} \tilde{\mathbf{C}}_t \\ \mathbf{o}_t \\ \mathbf{i}_t \\ \mathbf{f}_t \end{bmatrix} = \begin{bmatrix} \tanh \\ \sigma \\ \sigma \\ \sigma \end{bmatrix} \left( \mathbf{W} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{h}_{t-1} \end{bmatrix} + \mathbf{b} \right), \quad (8)$$

$$\mathbf{C}_t = \mathbf{f}_t * (\mathbf{C}_{t-1} + \tilde{\mathbf{C}}_t^e) + \mathbf{i}_t^e * \tilde{\mathbf{C}}_t,$$

$$\mathbf{h}_t = \mathbf{o}_t * \tanh(\mathbf{C}_t),$$

where  $\tilde{\mathbf{C}}_t^e$  denotes the calibrated memory-augmented cell state calculated through Eq. (9), and  $\mathbf{i}_t^e$  represents the enhanced input gate which can be computed through Eq. (10):

$$\begin{aligned} \mathbf{d}_t &= \tanh(\mathbf{W}_{enh} \mathbf{C}_{t-1} + \mathbf{b}_{enh}), \\ \mathbf{d}_t^e &= \sigma(\mathbf{U}_d^r \text{ReLU}(\mathbf{W}_d^r \mathbf{x}_t^a)), \\ \tilde{\mathbf{C}}_t^e &= \mathbf{d}_t * \mathbf{d}_t^e, \end{aligned} \quad (9)$$

$$\begin{aligned} \hat{\mathbf{i}}_t^e &= \tanh(\mathbf{W}_i^e \mathbf{x}_t^a + \mathbf{b}_i^e), \\ \tilde{\mathbf{i}}_t^e &= \mathbf{i}_t * \hat{\mathbf{i}}_t^e, \\ \mathbf{i}_t^e &= \sigma(\tilde{\mathbf{i}}_t^e + \mathbf{i}_t) \\ &= \sigma((\tilde{\mathbf{i}}_t^e + 1)\mathbf{i}_t), \end{aligned} \quad (10)$$

where  $\mathbf{d}_t$  indicates the obtained history memory state,  $\mathbf{d}_t^e$  denotes the calibrated gate calculated with auxiliary input  $\mathbf{x}_t^a$ . Afterwards,  $\mathbf{d}_t^e$  will be used to obtain the calibrated memory-augmented cell state value  $\tilde{\mathbf{C}}_t^e$  by multiplying the  $\mathbf{d}_t$ . By this way, the calibrated gate can selectively assign more weights to the representative and predictive memory neurons while suppressing the unimportant neurons. For the input cell, besides the primary input  $\mathbf{x}_t^p$  itself, it is also influenced by the auxiliary input  $\mathbf{x}_t^a$ . Under such a circumstance, the calibrated auxiliary input  $\hat{\mathbf{i}}_t^e$  is introduced to calculate the auxiliary influence score  $\tilde{\mathbf{i}}_t^e$  by multiplying the normal input gate  $\mathbf{i}_t$ . Finally, the enhanced input gate  $\mathbf{i}_t^e$  is computed by the addition of the auxiliary influence score  $\tilde{\mathbf{i}}_t^e$  and the normal input

gate  $\mathbf{i}_t$ , and then adjusted to the value between 0 and 1 via a sigmoid function.

Therefore, for the diagnoses codes embedding module, we can calculate the final visit-level diagnoses codes embedding  $\mathbf{h}_t^d$  by fusing with the historical diagnosed disease as:

$$\mathbf{h}_t^d = \text{InLSTM}_d(\mathbf{h}_{t-1}^d, \mathbf{e}_t^{dc}, \mathbf{h}_t^l), \quad (11)$$

where  $\text{InLSTM}_d$  denotes the proposed InLSTM (Eq.(7)) in this module, which is used to capture the correlations between the primary input of diagnosis code embedding sequence and the auxiliary input of laboratory results embedding sequence.

### 3.2.3. Medications codes embedding module

Similar to the previous hierarchy, diagnoses codes embedding module, the medications codes embedding module still comprise three main parts: a code embedding layer, an attentional selective module, and the temporal sequence interactive learning network.

**Medication code embedding layer.** In this module, MeSIN still first encodes each medication code  $\mathbf{c}_z^m$  into a dense embedding vector  $\mathbf{e}_z^m \in \mathbb{R}^d$  as:

$$\mathbf{e}_z^m = \mathbf{W}_e^m \mathbf{c}_z^m, \quad (12)$$

where  $\mathbf{W}_e^m \in \mathbb{R}^{d \times |C^*|}$  is the embedding matrix of medicinal codes that needs to be learned. Then we can gain the dense medication code embedding set  $\mathbf{E}_z^m = [\mathbf{e}_1^m, \dots, \mathbf{e}_n^m] \in \mathbb{R}^{d \times |n|}$ , in which each embedding is extracted from  $\mathbf{E}^m = [\mathbf{e}_1^m, \dots, \mathbf{e}_{|C_m^*|}^m] \in \mathbb{R}^{d \times |C_m^*|}$  if it existed in  $i$ -th visit.

**ASM for medications codes embeddings selection.** As mentioned before, MeSIN needs to filter out the noise coming from irrelevant historical medication codes sets medications at each timestamp. In this case, we should assign different relevance scores on different codes embeddings using the attentional selective module for computing the enhanced medications set embedding  $\mathbf{e}_{t-1}^{m_c}$ :

$$\begin{aligned} \mathbf{e}_{t-1}^{m_c} &= \sum_{i=1}^n \alpha_i^m \mathbf{e}_i^m (\alpha_i^d \in \alpha^d), \\ \alpha^m &= \alpha - \text{entmax}([\alpha_1^m, \alpha_2^m, \dots, \alpha_n^m], \gamma_m), \\ \alpha_i^m &= \tanh(\mathbf{W}_{m_a}^T \mathbf{e}_i^m + \mathbf{b}_{m_a}), \end{aligned} \quad (13)$$

where  $\mathbf{W}_{m_a} \in \mathbb{R}^d$  and  $\mathbf{b}_{m_a} \in \mathbb{R}$  are the parameters of ASM to be learned.  $\gamma_m$  is the hyper-parameter of  $\alpha - \text{entmax}$  [27] in this module. Likely, we can compute the enhanced medications codes embeddings sequence  $\mathbf{E}^{m_c} = \{\mathbf{e}_1^{m_c}, \dots, \mathbf{e}_{t-1}^{m_c}\}$  at historical timestamps.

**InLSTM in medications codes sequence learning.** Finally, for capturing the temporal dependency of historical medications, the gained enhanced medication code embedding sequence  $\mathbf{E}^{m_c}$  is treated as the primary input of sequence learning network. In addition, recommending medications is essentially a comprehensive decision-making process, the historical prescribed medications must be affected

by the laboratory results and diagnosed diseases. Therefore, the sequences of laboratory results and diagnosed diseases are taken as the auxiliary input assisting in controlling the sequence learning process. Then, we can calculate the final visit-level medication code embedding  $\mathbf{h}_{t-1}^m$  by fusing the historical disease progression using InLSTM Eq.(7) as:

$$\mathbf{h}_{t-1}^m = \text{InLSTM}_m(\mathbf{h}_{t-2}^m, \mathbf{e}_{t-1}^{m_c}, \mathbf{h}_{t-1}^c), \quad (14)$$

where  $\text{InLSTM}_m$  denotes the proposed InLSTM in this module, which is used to capture the correlations between the primary input of diagnosis code embedding and the auxiliary input  $\mathbf{h}_{t-1}^c$ . Further, the auxiliary input  $\mathbf{h}_{t-1}^c$  is calculated via a fusion module:

$$\mathbf{h}_{t-1}^c = \sigma(\mathbf{W}_c^d \mathbf{h}_{t-1}^d + \mathbf{W}_c^l \mathbf{h}_{t-1}^l), \quad (15)$$

where  $\mathbf{W}_c^d, \mathbf{W}_c^l \in \mathbb{R}^{d \times d}$ ,  $\sigma$  denotes the activation function  $\tanh$ .  $\mathbf{h}_{t-1}^l$  and  $\mathbf{h}_{t-1}^d$  respectively represent the obtained visit-level laboratory results embedding using Eq. (4) and diagnosed diseases embedding using Eq. (11).

### 3.2.4. Global selective fusion module

In step **I**, by modeling the hierarchically interactive temporal sequence learning process, we can obtain the visit-level laboratory results embedding  $\mathbf{h}_t^l$ , diagnoses codes embedding  $\mathbf{h}_t^d$  and the medications codes embedding  $\mathbf{h}_{t-1}^m$ . All above three kinds of embeddings have incorporated corresponding historical information. For recommending medications at current timestamp, the current enhanced laboratory results embedding  $\mathbf{e}_t^{l_c}$  and diagnosis code embedding  $\mathbf{e}_t^{d_c}$  should be given more attention when making final decisions.

To effectively fuse above five heterogeneous embeddings according to their importance scores and minimize the effect introduced by irrelevant information as much as possible, in step **II**, we design a global selective fusion module, which is realized by a self-attention mechanism. Since the five types of embeddings are heterogeneous, here, we first calculate the information importance scores by themselves as:

$$\begin{bmatrix} \alpha_m \\ \alpha_d \\ \alpha_l \\ \alpha_{dc} \\ \alpha_{lc} \end{bmatrix} = \begin{bmatrix} \mathbf{W}_m^a \\ \mathbf{W}_d^a \\ \mathbf{W}_l^a \\ \mathbf{W}_{dc}^a \\ \mathbf{W}_{lc}^a \end{bmatrix} * \begin{bmatrix} \mathbf{h}_{t-1}^m \\ \mathbf{h}_t^d \\ \mathbf{h}_t^l \\ \mathbf{e}_t^{d_c} \\ \mathbf{e}_t^{l_c} \end{bmatrix} + \begin{bmatrix} \mathbf{b}_m^a \\ \mathbf{b}_d^a \\ \mathbf{b}_l^a \\ \mathbf{b}_{dc}^a \\ \mathbf{b}_{lc}^a \end{bmatrix}, \quad (16)$$

where  $\mathbf{W}_m^a, \mathbf{W}_d^a, \mathbf{W}_l^a, \mathbf{W}_{dc}^a, \mathbf{W}_{lc}^a$  and  $\mathbf{b}_m^a, \mathbf{b}_d^a, \mathbf{b}_l^a, \mathbf{b}_{dc}^a, \mathbf{b}_{lc}^a$  are the parameters to be learned.  $\alpha_m, \alpha_d, \alpha_l, \alpha_{dc}, \alpha_{lc}$  are the information importance scores, by which we can calculate the final information importance scores  $\mathbf{a}$  as:

$$\mathbf{a} = \text{softmax}([\alpha_m, \alpha_d, \alpha_l, \alpha_{dc}, \alpha_{lc}]), \quad (17)$$

where  $\mathbf{a} = [\alpha'_m, \alpha'_d, \alpha'_l, \alpha'_{dc}, \alpha'_{lc}]$ .

Finally, we obtain the ultimate patient representation vector by summing up the heterogeneous information vectors according to importance scores from Eq.(17) as:

$$\mathbf{O}_p = \alpha'_m \mathbf{h}_{t-1}^m + \alpha'_d \mathbf{h}_t^d + \alpha'_l \mathbf{h}_t^l + \alpha'_{dc} \mathbf{e}_t^{d_c} + \alpha'_{lc} \mathbf{e}_t^{l_c}, \quad (18)$$

where  $\mathbf{O}_p$  is the patient representation vector to be used to recommend reasonable medications in the next subsection.

### 3.2.5. Medication recommendation

Doctors make decisions about recommending reasonable medications for patients after comprehensive consideration. Likewise, the learned patient representation  $\mathbf{O}_p$  is employed in this study to recommend reasonable medications as:

$$\hat{\mathbf{y}}_t = \text{softmax}(\mathbf{W}_o \cdot \mathbf{O}_p + \mathbf{b}_o), \quad (19)$$

where  $\hat{\mathbf{y}}_t$  denotes the set of recommended multi-label medications,  $\mathbf{W}_o \in \mathbb{R}^{d_m \times r}$  and  $\mathbf{b}_o \in \mathbb{R}^{d_m}$  are parameters to be learned.

### 3.3. Model training

Since medication recommendation task belongs to the domain of sequential multi-label prediction task, we utilize the binary cross-entropy loss  $\mathcal{L}_{ce}$  and multi-label margin loss  $\mathcal{L}_{mg}$  as the objective functions. The prediction objective function binary cross-entropy loss  $\mathcal{L}_{ce}$  is formulated as:

$$\mathcal{L}_{ce} = - \sum_{t=1}^T \mathbf{y}_t \log \sigma(\hat{\mathbf{y}}_t) + (1 - \mathbf{y}_t) \log(1 - \sigma(\hat{\mathbf{y}}_t)), \quad (20)$$

The corresponding objective function multi-label margin loss  $\mathcal{L}_{mg}$  is:

$$\mathcal{L}_{mg} = \sum_t \sum_i |C| \sum_j \frac{\max(0, 1 - (\hat{\mathbf{y}}_t[Y_j^t] - \hat{\mathbf{y}}_t[i]))}{L}, \quad (21)$$

where  $\hat{\mathbf{y}}_t[i]$  is the value of  $i^{th}$  coordinate at  $t^{th}$  visit and  $\hat{\mathbf{y}}_t[Y_j^t]$  denotes the predicted label value indexed by  $Y_j^t$ , the  $j^{th}$  value in the ground truth label set  $Y^t$  at  $t$ -th visit for a patient. For function (20) and (21). So we can get two binary cross entropy loss functions  $\mathcal{L}_{ce}^d, \mathcal{L}_{ce}^m$ , and two multi-label margin loss functions  $\mathcal{L}_{mg}^d, \mathcal{L}_{mg}^m$ .

To facilitate the joint optimization process of two tasks, we combine the aforementioned loss functions to build a joint loss function  $\mathcal{L}$ :

$$\mathcal{L} = \eta \mathcal{L}_{ce} + \varepsilon \mathcal{L}_{mg}, \quad (22)$$

where  $\eta, \varepsilon \geq 0$  are the mixture weights, and  $\eta + \varepsilon = 1$ . The training algorithm is detailed in Algorithm 1.

## 4. Experiments and discussion

### 4.1. Datasets description

As is analyzed in Section 1, the aim of study is to recommend medications for patients based on the heterogeneous multilevel EHR data. Hence, we should conduct experiments on a cohort where patients have at least two visits and their EHRs are complete. Here, we choose a real-world publicly available dataset MIMIC-III [20]<sup>1</sup>, in which patients stayed within the intensive care units (ICU) at Beth Israel Deaconess Medical Center and had relatively complete

<sup>1</sup><https://mimic.physionet.org>

**Algorithm 1** Model training for MeSIN.

---

**Require:** Training set  $\mathbf{R}$ , training epochs  $N$ , batch size  $BS$ , mixture weights  $\eta, \varepsilon$  in Eq. (22);  
Use uniform distribution to initialize the model parameters  $\theta \sim U(-1, 1)$ ;

- 1: **for**  $i = 1$  to  $N * |\mathbf{R}|$  **do**
- 2:   Sample  $BS$  patients ( $\mathbf{P} = \{\mathcal{X}^1, \dots, \mathcal{X}^{T_n}\}$ ) from  $\mathbf{R}$ ;
- 3:   **for**  $t = 1$  to  $T_i$  **do**
- 4:     `/**Laboratory results embedding module**/`
- 5:     Obtain multi-channel time-series embeddings  $\{h_{1T}^l, \dots, h_{qT}^l\}$  using Eq. (2);
- 6:     Obtain enhanced laboratory results embedding  $e_i^{lc}$  using ASM using Eq. (3);
- 7:     Compute the visit-level laboratory results embedding  $h_t^l$  using Eq. (4);
- 8:     `/**Diagnoses codes embedding module**/`
- 9:     Obtain diagnoses codes embeddings  $E_t^d = [e_1^d, \dots, e_m^d] \in \mathbb{R}^{d \times m}$  using Eq. (5);
- 10:     Obtain the enhanced diagnosis code embedding  $e_t^{dc}$  using ASM as Eq. (6);
- 11:     Compute the visit-level diagnoses codes embedding  $h_t^d$  using InLSTM (Eq. (7-11));
- 12:     `/**Medications codes embedding module**/`
- 13:     Obtain the medications codes embeddings  $E_z^m = [e_1^m, \dots, e_n^m] \in \mathbb{R}^{d \times n}$  using Eq. (12);
- 14:     Obtain the enhanced medications codes embedding  $e_{t-1}^{mc}$  using ASM as Eq. (13);
- 15:     Compute the visit-level medications codes embedding  $h_{t-1}^m$  using InLSTM (Eq. (7-11));
- 16:     `/**Global selective fusion module**/`
- 17:     Incorporate the multi-source embeddings into patient representation  $\mathbf{O}_p$  using Eq. (16-18);
- 18:     Compute recommended medications  $\hat{\mathbf{y}}^t$  in Eq.(19);
- 19:    **end for**
- 20:    Update  $\theta$  by optimizing the total loss  $\mathcal{L}$  in Eq. 22;
- 21: **end for**

---

health records with multilevel heterogeneous data. Though MIMIC-III belongs to the ICU data, there are certain patients with multiple visits. Hence, we utilize it as our experimental dataset. Similar to [31], we choose the medications prescribed by doctors for each patient within the first 24 hours as medicine set since it is usually a critical period for each patient to get rapid and accurate treatment [12]. Besides, the medicine codes form NDC are transformed to ATC Level 3 for integrating with MIMIC-III. Meanwhile, we employ the second hierarchy codes of the ICD9 codes<sup>2</sup> as the disease category labels, since predicting category information not only guarantees the sufficient granularity of all the diagnoses but also improves the training speed and predictive performance [24, 7]. For considering the laboratory results into decision-making process, we follow the feature extraction method used in [13]. Here, the time-window of each laboratory indicator is 24 hours. More information about

<sup>2</sup><http://www.icd9data.com>

**Table 1**

Statistics of the MIMIC-III datasets

MIMIC III	Quantity
# of patients	4631
# of unique diagnosis	1879
# of unique medication	143
# of unique laboratory indicators	17
avg # of visits	2.55
avg # of diagnoses	10.16
avg # of medications	7.33

the patients cohort from the dataset is listed in Table 1.

**4.2. Evaluation metrics**

To evaluate the performance, we adopt the Jaccard Similarity Score (Jaccard), Precision Recall AUC (PR-AUC), Average Recall (Recall) and Average F1 (F1) as the evaluation metrics. Jaccard is defined as the size of the intersection divided by the size of the union of predicted set  $\hat{Y}_t^i$  and ground truth set  $Y_t^i$ . Precision is used to measure the correctness of predicted medicines and Recall is used to measure the completeness of predicted medicines. F1 is often used as the comprehensive evaluation metric of prediction model.

$$\text{Jaccard} = \frac{1}{\sum_i^N \sum_t^{T_i} 1} \sum_i^N \sum_t^{T_i} \frac{|Y_t^i \cap \hat{Y}_t^i|}{|Y_t^i \cup \hat{Y}_t^i|}, \quad (23)$$

where  $N$  denotes the number of patients in test set and  $T_i$  is the number of visits for the  $i^{\text{th}}$  patient. Given

$$\text{Recall} = \frac{1}{\sum_i^N \sum_t^{T_i} 1} \sum_i^N \sum_t^{T_i} \frac{|Y_t^i \cap \hat{Y}_t^i|}{|Y_t^i|}, \quad (24)$$

$$\text{Precision} = \frac{1}{\sum_i^N \sum_t^{T_i} 1} \sum_i^N \sum_t^{T_i} \frac{|Y_t^i \cap \hat{Y}_t^i|}{|\hat{Y}_t^i|}, \quad (25)$$

the valuation metric F1 can be calculated as:

$$\text{Jaccard} = \frac{1}{\sum_i^N \sum_t^{T_i} 1} \sum_i^N \sum_t^{T_i} \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}. \quad (26)$$

**4.3. Benchmark methods**

To evaluate the effectiveness of the proposed model, it was compared to the following baseline methods:

- *Nearest*. To predict treatment medicines for a patient  $p_i$ , Nearest was proposed to choose the treatment medications prescribed for patient  $p_j$ , who has the most similar historical laboratory indicators and medications with  $p_i$ .
- *LR* [23]. It is a logistic regression with L1/L2 regularization. We sum the multi-hot vector of each visit together and apply the binary relevance technique [23] to handle multi-label output.

- *NBN* [1]. Here, this method mainly utilizes the prior knowledge and employ the statistical methods to recommend corresponding medications for patients.
- *Retain* [8]. RETAIN is an interpretable model with a two-level reverse time attention mechanism to predict diagnoses, which can detect significant past visits and associated clinical variables. It can be used for similar sequential prediction tasks, such as predicting treatment medicines.
- *DELSTM* [18]. This model utilizes additional input sequence as the input of a decomposed gate to control the memory cell state, which indirectly interacts with primary input sequence.
- *PCLSTM* [18]. This structure takes all heterogeneous sequences as input. In other words, multiple sequences interact with each other via the neuron interactions in the way of concatenating both hidden states.
- *RAHM* [2]. It builds a relation augmented hierarchical multi-task learning framework for learning multi-level relation aware patient representation for medication prediction.
- *LEAP* [39]. Leap formulates the medicine prediction problem as a multi-instance multi-label learning problem, which mainly uses a recurrent neural network (RNN) to recommend medicines.
- *DMNC* [22]. DMNC uses a memory augmented neural network to model the interaction of two asynchronous sequences for treatment prediction task [22].
- *GAMENet* [31]. It employs a dynamic memory network to save encoded historical medication information, and further utilizes a query representation formed by encoding sequential diagnosis and procedure codes to retrieve medications from the memory bank.
- *AMANet* [15]. AMANet leverages self-attention and inter-attention to capture the intra-view and inter-view interactions. Then it concatenates the information from history attention and dynamic external memory to predict the medications.

#### 4.4. Experimental settings

We randomly split the patients in MIMIC-III dataset into training, validation and test sets with 2/3 : 1/6 : 1/6 ratios. The random splitting and training processes were performed five times. Table 2 lists the results as averages across the five runs obtained for all the compared models in terms of the four evaluation metrics described in Section 4.2. Specifically, the embedding size and the hidden layer dimension for LSTM and GRU are all set as 128 and 128, respectively. The dropout rate is set as 0.4, batch size is set as 10, and the mixture weights of objective function are set as  $\eta = 0.99, \epsilon = 0.01$ . The values of attention sparse degree controlling parameters  $\gamma_l$  in ASM are set as 1.5, 1.5, 1.3, respectively. Training is done through Adam [21] at learning rate

**Table 2**

Performance Comparison of Benchmark Methods on MIMIC-III Dataset

Methods	Jaccard	PR-AUC	Recall	F1
Nearest	0.2019	0.2227	0.3099	0.2899
LR [23]	0.3401	0.5549	0.4549	0.4901
NBN [1]	0.3341	0.5479	0.5081	0.4839
Retain [8]	0.3267	0.5364	0.4841	0.4905
DELSTM [18]	0.3357	0.5371	0.5144	0.5005
PCLSTM [18]	0.3301	0.5109	0.4952	0.4940
RAHM [2]	0.3558	0.5393	0.5357	0.5122
LEAP[39]	0.3111	0.4212	0.4541	0.4627
DMNC [22]	0.3272	0.4476	0.5136	0.5021
GAMENet [31]	0.3452	0.4828	0.5246	0.5115
AMANet [15]	0.3641	0.5595	0.5547	0.5381
MeSIN <sub>DE</sub>	0.3929	0.5633	0.5717	0.5624
MeSIN <sub>Soft</sub>	0.3941	0.5663	0.5798	0.5636
MeSIN	<b>0.3975</b>	<b>0.5684</b>	<b>0.5934</b>	<b>0.5670</b>

2e-4, and we report the model performance in test set within 40 epochs. All methods are trained on an Ubuntu 16.04 with 64GB memory and Nvidia TAITAN XP GPU using the Pytorch 1.0 framework.

#### 4.5. Performance comparison

As demonstrated in Table 2, the benchmark models used in health informatics are divided into three categories: Shallow methods, including Nearest, LR and NBN; Predictive models, including Retain, DELSTM, PCLSTM and RAHM; Recommendation models, including LEAP, DMNC, GAMENet and AMANet. From the table, we have an impressive observations that the proposed MeSIN achieves the superior performance over the listed benchmark models. Through detailed comparison with the benchmark models, several interesting observations can be made as follows.

First, as for the shallow methods, Nearest, LR and NBN achieved about at least 5.74%, 1.35%, 8.53% and 7.69% lower scores on medication recommendation task with respect to Jaccard, PR-AUC, Recall and F1 score, respectively, than MeSIN. On the one hand, this kind of method does not consider the temporality and heterogeneity of EHR data, and also overlook the relations between multiple sequences. On the other hand, Nearest achieves the worst performance which indicates that most of the patients possess distinct disease conditions within continuous admissions to hospital.

Second, as for predictive models in health informatics, Retain, DELSTM, PCLSTM and RAHM also achieve poor performance on medication recommendation task than MeSIN. We think that the main reason might be that they can not consider the current medical records including laboratory indicators and diagnosed diseases status, and only take the historical records into accounts. In addition, Retain is a two-

level attention based model, that can capture temporal correlations and identify influential past visits. However, it can not consider all heterogeneous sequences respectively and can just concatenate the heterogeneous embeddings into one embedding which would confuse the embeddings obtained from different hierarchies of EHR data. DELSTM and PCLSTM mainly pay more attention on the sequences interaction in the temporal sequence learning process while overlook the inherent hierarchy structure of EHR data. In MeSIN, the multilevel learning framework is incorporated to extract useful information from such kind of inherent structure. RAHM partially employs the hierarchy nature of EHR data to acquire better performance via multi-task learning framework. However, it still employ single sequence learning methods to integrate with the historical information which might cause the confusion of different historical medical sequences.

Third, our proposed MeSIN outperforms all state-of-the-art methods used for medication recommendation such as LEAP, DMNC, GAMENet and AMANet about at most 8.64%, 14.72%, 13.93% and 10.43%, and at least 3.34%, 0.92%, 3.87% and 2.89% with respect to Jaccard, PR-AUC, Recall and F1 score. In practice, the medication recommendation problem within an admission might not be the pure sequential recommendation process and it also refers to the diverse correlations. Their poor performance could be attributable to the poor ability to capture such complicated correlations. Particularly, LEAP cannot capture the inherent multiple relations among heterogeneous sequences. While DMNC realizes the interactions of two sequences through attention based DNC blocks but neglecting the utilization of medications in the history visits. Similarly, AMANet also does not consider historical medications prescribed for patients, but it achieves relatively better performance through multiple attention networks for capturing the inter- and intra- correlations of heterogeneous sequences. However, AMANet neglects the captured evolution information such as disease progression through temporal sequence learning network, which is still a kind of important information in decision-making process.

Finally, we can observe the MeSIN also outperforms two special variants, MeSIN<sub>DE</sub> and MeSIN<sub>Soft</sub>. For the former variant, we replace the developed interactive LSTM network (InLSTM) in MeSIN with DELSTM network [18]. In this variant, the interaction processes just consider the impact of auxiliary input on the memory cell state of primary input sequence learning network, and overlook the impact on the current input cell state. While the InLSTM network in MeSIN simultaneously considers the sequential interactions from above two aspects. For the latter variant, we replace the incorporated attention weights computation method *Entmax* in attentional selective module (ASM) of MeSIN with *Softmax*. In this variant, unlike *Entmax*, *Softmax* will generate the dense attention alignments that is wasteful, and can not pay more focus on the really important feature embeddings.

Therefore, the critical reasons that MeSIN achieves the best performance compared with all benchmark models can be summarized as follows: (1) The multilevel learning frame-

work can help capture the inherent causal relations of adjacent hierarchies; (2) The incorporated multiple attentional selective modules in the framework realizes the effective embeddings selection and make the learned patient representation be more expressive; (3) The designed InLSTM further reinforces the sequences interactions from both the historical memory cell and the input cell, which can further optimize the temporal sequence learning process by incorporating more useful calibrated information.

#### 4.6. Ablation study

We now need to examine the effectiveness of different components in MeSIN and evaluate the contribution of different source data. Hence, we conduct two kinds of ablation studies respectively on model's components and multi-sourced EHR data.

##### 4.6.1. Model components

This ablation study is conducted to verify the effectiveness of different MeSIN components to its overall performance. To determine whether the incorporated components improve the performance, we add them one by one from scratch and verify their performance by all evaluation metrics including Jaccard, PRAUC, Recall and F1 score. Table 3 presents the recommendation results of distinct MeSIN variants on the MIMIC-III dataset. One of the basic baseline models, Vanilla, the medical codes are respectively added together as the enhanced embeddings in every module. Besides, the standard LSTM networks are also respectively employed as the temporal sequence learning networks in three distinct modules, and we employ the concatenation-based fusion method to replace the proposed global selective fusion module. However, Vanilla still achieves relatively better performance compared with benchmark models, which can attribute to the incorporation of multi-sources data and the integration of current laboratory results and diagnosed disease by concatenation-based fusion method.

**Attentional selective module (ASM).** As explained in Section 3.2, ASM is introduced to automatically select the useful information and filter out noise information as much as possible by assigning corresponding attention weights to embeddings according to their respective importance. The following variants are tested to evaluate the contribution of ASMs from different modules to the overall performance of MeSIN:

- ASM<sub>L</sub>. In this variant, we incorporate an attentional selective module for laboratory results embeddings selection. The overall performance is slightly improved by 0.17% on Jaccard in this case compared with the Vanilla model. This testifies that the introduced ASM module can help focus on the useful laboratory results embeddings by controlling the value of  $\gamma_l$ , by which we can obtain relatively better enhanced embedding as the input of temporal sequence learning network.
- ASM<sub>LD</sub>. Similar to ASM<sub>L</sub>, in this variant, we introduce an attentional selective module to replace the addition operation for diagnoses codes embeddings selection. In this

**Table 3**  
Performance Comparison of the variants of MeSIN on MIMIC-III Dataset

Model	ASM			InLSTM		GSFM	Recommendation performance			
	Lab	Diag	Med	Diag	Med	Fusion	Jaccard	PR-AUC	Recall	F1
Vanilla	✗	✗	✗	✗	✗	✗	0.3832	0.5579	0.5451	0.5482
ASM <sub>L</sub>	✓	✗	✗	✗	✗	✗	0.3849	0.5592	0.5522	0.5501
ASM <sub>LD</sub>	✓	✓	✗	✗	✗	✗	0.3865	0.5595	0.5549	0.5546
ASM <sub>LDM</sub>	✓	✓	✓	✗	✗	✗	0.3887	0.5592	0.5684	0.5583
ASM_InLSTM <sub>D</sub>	✓	✓	✓	✓	✗	✗	0.3916	0.5643	0.5673	0.5619
ASM_InLSTM <sub>DM</sub>	✓	✓	✓	✓	✓	✗	0.3935	0.5651	0.5768	0.5639
MeSIN	✓	✓	✓	✓	✓	✓	<b>0.3975</b>	<b>0.5684</b>	<b>0.5934</b>	<b>0.5670</b>

way, the diagnoses codes embeddings that are irrelevant with the recommendation task would be discarded by sparse attention under the value control of  $\gamma_d$ . As a result, the performance of ASM<sub>LD</sub> is improved by 0.16% on Jaccard, which indicates the importance of the ASM in MeSIN in selecting the useful information from numerous medical codes embeddings.

- ASM<sub>LDM</sub>. In this variant, we further incorporate the third ASM module into the prescribed medications embedding module for selecting the most relevant historical medication codes embeddings to build the patient representation. As a result, ASM<sub>LDM</sub> makes relatively better improvement compared with ASM<sub>L</sub> and ASM<sub>LD</sub>. We think that it can attribute that the medication embedding module has direct relevance with the medication recommendation task. In the end, the incorporation of above three attentional selective modules bring about 0.55% on Jaccard, 0.13% on PR-AUC, 2.33% on Recall, 1.01% on F1 in total compared with Vanilla. But the important is that ASM makes MeSIN more interpretable by focusing on the really important features.

**Interactive Long-Short Term Memory network (InLSTM).** InLSTM is developed for reinforcing the interaction process of heterogeneous sequences, which is beneficial to capture the correlations of sequences. The following variants are tested to evaluate the contribution of InLSTM to the overall performance of MeSIN:

- ASM\_InLSTM<sub>D</sub>. In this variant, we incorporate a novel InLSTM to replace the standard LSTM in ASM<sub>LDM</sub> in diagnoses codes embedding module for enhancing the interaction process of disease progression and changing laboratory results. It achieves by 0.29% on Jaccard compared with ASM<sub>LDM</sub>, which verifies the importance of considering the correlations of sequences such as between laboratory results and diagnosed diseases into the temporal sequence learning process.
- ASM\_InLSTM<sub>DM</sub>. Here, the interactive LSTM is further introduced to the top hierarchy, prescribed medications embedding module for facilitating the medication

**Table 4**  
Contribution of different data sources to MeSIN performance

Methods	Jaccard	PR-AUC	Recall	F1
NoLab	0.3861	0.5612	0.5526	0.5567
NoDiag	0.3552	0.5548	0.5158	0.5236
NoMed	0.3859	0.5587	0.5514	0.5537
AllData	<b>0.3975</b>	<b>0.5684</b>	<b>0.5934</b>	<b>0.5670</b>

codes sequential learning process. In this sequence learning network, the diagnosed diseases are utilized to enhance the interaction process with prescribed medications for providing complimentary useful information. Thus, the performance of ASM\_InLSTM<sub>DM</sub> outperforms the fifth variant ASM\_InLSTM<sub>D</sub> by 1.9% on Jaccard, which further indicates the superiority of InLSTM in MeSIN than the standard LSTM in ASM<sub>LDM</sub>.

**Global selective fusion module (GSFM).** After step I, the hierarchically interactive temporal sequence learning procedure, the obtained multi-sourced embeddings are integrated together via proposed global selective fusion module (GSFM) for obtaining the patient representation. In this way, MeSIN can automatically learn the contribution scores of distinct embeddings to the medication recommendation task. As a result, it improves by 0.4% on Jaccard compared with the sixth variant ASM\_InLSTM<sub>DM</sub>. This also indicates the advantage of GSFM than the concatenation-based method used in above six variant models. However, owe to the utilization of concatenation-based fusion method, Vanilla gains relatively better performance than the benchmark methods.

#### 4.6.2. Heterogeneous Data

According to the proposed method, multilevel EHR data need to be input to MeSIN for obtaining the final patient representation. Though each of them plays a paramount role in the clinical decision-making scenario, here, we would build the following MeSIN variants to evaluate the impact of different heterogeneous data on medication recommendation results (Table 4). In *NoLab*, the laboratory results embed-

ding module in MeSIN is removed, and the introduced InLSTM<sub>d</sub> in diagnoses codes embedding module needs to be replaced by the standard LSTM network. In this case, patient's detailed health status is unknown. In NoDiag, the diagnoses codes embedding module is removed from MeSIN, and just retains the remain two modules. Under such circumstance, the learned patient representation will lose the key disease progression information. In NoMed, the medications codes embedding module is removed from MeSIN. In this way, the learned patient representation will lose the historical medications information.

Clearly, it can be noticed from Table 4 that the performance of variants all drop owing to some apparent reasons. First, in practice, most medications are prescribed conditioned on the diagnosed diseases. Therefore, the performance of *NoDiag* drops dramatically, which validates the crucial role of diagnosed diseases and disease progression in medication recommendation task. Second, though historical medications in ICU are not so much valuable for most patients, but are still a kind of important information in understanding patient's history diseases which can help know some detailed information such as allergic condition. Thus, the performance of *NoMed* also drops significantly on medications recommendation task. Third, the performance of *NoLab* also drops significantly in this case but slightly compared with *NoDiag*. The main reason is that the current patient health status including diagnosed diseases and key laboratory indicators results is still a paramount indicator indicating patient's health status. Thus, though the importance of each hierarchy data within EHRs are diverse from each other, all of them play important roles in medication recommendation task.

#### 4.7. Attention analysis in selective module

As discussed above, our newly developed MeSIN outperforms all benchmark models on medication recommendation for patients. Among the constituent components of MeSIN, the attentional selective module (ASM) plays a great role in the model, which has been testified through ablation studies about MeSIN in section 4.6.1. Actually, the positive influence of ASM should attribute to the selective ability of *entmax*, which can increase focus on important medical codes embeddings and make the process more interpretable. Hence, we perform attention analysis to explore the crucially attentive process shown in Figure 4, visualize the difference of softmax and entmax shown in Figure 5, and investigate the importance of multi-source embeddings shown in Figure 6.

**The attentive process.** To clearly interpret the attentive process, as shown in Figure 4, we just consider the relations between the second and third hierarchies (diagnoses codes embedding module and prescribed medications embedding module) within our multilevel learning framework. In addition, the quantitative value in column DA and column MA respectively denotes the attention weights calculated by Eq. (6) and Eq. (13). As shown in Figure 4, the attentive process can be categorized into four distinct but correlated processes. In this case, we have four interesting observations.

Visit 1				Visit 2				Visit 3				Medication Recommendation	
DC	DA	MC	MA	DC	DA	MC	MA	DC	DA	DC	DA	Recommend	Label
40301	0.243	N02B	0.022	99673	0.112	N02B	0.034	40301	0.107			A06A	A06A
5856	0.219	A02B	0.013	5856	0.114	B01A	0.025	5856	0.094			A02B	A02B
2724	0.182	B01A	0.009	4538	0.012	C09A	0.184	41071	0.103			C10A	C10A
41401	0.208	A01A	0.024	40301	0.130	C01D	0.179	99673	0.092			C07A	C07A
4168	0.148	C09A	0.142	41071	0.117	A04A	0.151	44489	0.050			N02B	N02B
28521	0.000	C10A	0.161	7863	0.084	A06A	0.052	4254	0.044			B01A	B01A
		A06A	0.025	4280	0.075	V03A	0.155	4592	0.096			A04A	N02A
		A04A	0.124	E8791	0	C07A	0.021	E8791	0			A01A	A04A
		C01D	0.145	30560	0.105	C08C	0.199	2767	0.073			C09A	A01A
		D04A	0	V4511	0.121			4240	0.061			C01D	C09A
		V03A	0.109	V1581	0.050			V4511	0.101			C08C	C01D
		C08C	0.137	V1251	0.079			V1581	0.039			V03A	C08C
		C07A	0					V5861	0.077				A03B
		N02A	0.089					V1251	0.063				C02A
		B03A	0										V03A

DC: Diagnoses codes DA: Attention in ASM<sub>d</sub> Disease codes: ICD-10  
MC: Medication codes MA: Attention in ASM<sub>m</sub> Medication codes: ATC

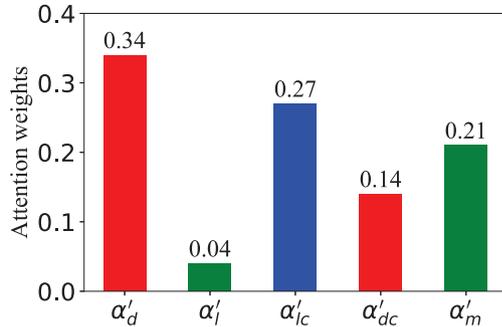
**Figure 4:** The visualization of attentive process between the second and third hierarchy within our multilevel learning framework.

First, in the attentive process (1), the learned visit-level diagnoses codes embedding will be input into the medication codes embedding module for interacting with the medications codes embedding within historical visits. Thus, we can observe that there exists strong causal relations between DC column (diagnoses codes) and MC (medication codes) within each visit. In this way, the diagnoses codes that corresponds to the medications codes existing in the recommendation label column will be assigned more attention weights within each visit. Second, owing to the temporal dependencies of EHR data, the recommended medications in *Recommend* column not only depends on the diagnosed diseases in DC column in the third visit, but also relies on the historical prescribed medications and diseases progression. As for this, as shown in the attentive process (2), the medication codes embeddings in historical visits but existing in *label* column will be assigned more attention weights. Similarly, in the attentive process (3), as for the inherent causal relations between diseases and medications, the corresponding diagnoses codes embeddings will be also assigned more attention weights. In the end, as shown in attentive process (4), for capturing the temporal dependency of EHR data, the medications codes sequence learning process will be influenced by the diagnoses codes sequence learning process under the help of proposed InLSTM in MeSIN.

**The calculation methods of attention weights: Softmax and Entmax.** In MeSIN, *Entmax* has been incorporated into the attentional selective module (ASM) to make more intelligent selections: assist to filter out noisy information and pay more focus on the important feature embeddings. In figure 5, we provide a hot map which demonstrates the difference of attention weights computed by *Entmax* [27] and *Softmax* in laboratory results embeddings selection module. As shown in figure, we observe that the calculation method *Entmax* can generate sparse attention weights within each visit, in other words, it can make the attention scores of some unimportant indicator results embeddings to zero such as

Visit 1	Softmax	0.017	0.011	0.047	0.140	0.126	0.127	0.203	0.045	0.044	0.045	0.011	0.011	0.025	0.046	0.046	0.046	0.011
	Entmax	0	0	0.080	0.159	0.080	0.080	0.239	0.048	0.036	0.039	0	0	0	0.080	0.080	0.080	0
Visit 2	Softmax	0.017	0.012	0.048	0.121	0.096	0.130	0.191	0.026	0.044	0.045	0.041	0.033	0.045	0.047	0.047	0.047	0.011
	Entmax	0	0	0.095	0.131	0	0.095	0.212	0.057	0.022	0.073	0	0	0.060	0.085	0.085	0.085	0
Visit 3	Softmax	0.018	0.012	0.050	0.113	0.100	0.137	0.201	0.048	0.049	0.048	0.034	0.012	0.018	0.050	0.049	0.049	0.012
	Entmax	0.000	0.000	0.098	0.099	0.000	0.099	0.197	0.070	0.083	0.078	0.000	0.000	0.000	0.099	0.089	0.089	0.000
		CRR	DBP	FIO	GCSE	GCSM	GCS	GCSV	GLU	HR	HT	MBP	SAT	RR	SBP	T	WT	pH

**Figure 5:** The comparison of attention weights computed by Entmax and Softmax in laboratory results embeddings selection module.



**Figure 6:** Attention weights distribution in global selective module.

CRR, DBP, MBP, OS, RR and PH. In this way, the MeSIN can intelligently make selections about which features embeddings are more important to be focused on and which are unnecessary to be paid so much attention on when making decisions in clinical decision-making process. Therefore, such attention weights computation method in ASMs can help MeSIN increase focus on the really important features embeddings and make the model more interpretable.

**The importance of multi-source embeddings.** As mentioned in section 3.2.4, the global selective fusion module, to fuse the obtained five heterogeneous embeddings, we introduce a global selective fusion module, which can integrate them into patient representation according to respective importance score and minimize the adverse effect caused by noisy information. In figure 6, we can observe that  $\alpha'_d > \alpha'_{lc} > \alpha'_m > \alpha'_{dc} > \alpha'_i$  (see details in Eq.(16-18)), which indicates the importance ranking of multi-source embeddings. Such a phenomenon further testifies that diagnosed diseases especially the disease progression with historical disease information is the most important information for the medication recommendation task, which have been proved in the ablation study shown in Table 4. The current laboratory result is the second critical factor when making decisions about the recommended medications. In addition, the historical prescribed medications are also taken into account. Finally, the historical laboratory results might be not so important in the intensive care unit (ICU). However, owing to that different patients might have different diseases status, the learned attention weights are also dynamically changing, which makes the computed relevance scores distribution are also diverse. For example, the historical medications might

be more important than the diagnosed diseases when the diagnosis is adverse drug reaction. Through the above analysis, we can see that MeSIN can provide some insightful and interpretable recommendation results.

## 5. Conclusion

In this paper, we propose a novel multilevel selective and interactive network for medication recommendation task with clinical EHR data. In our model, the inherent causal relations and temporal dependencies of EHR data are effectively captured via proposed multilevel learning framework and a novel interactive LSTM cell. Considering the inevitable noise within EHR data, multiple attentional selective modules are incorporated into model for paying more focus on the really important feature embeddings and meanwhile provide insightful and interpretable recommendation results. Finally, we evaluate our model on a real world and public clinical dataset. The experimental results show that our model achieves the best recommendation performance against eleven baselines in terms of Jaccard, PR-AUC, Recall and F1 score. In the future, we plan to adapt the proposed approach for more healthcare prediction tasks based on sequential data and explore its usage in domains other than healthcare.

## Acknowledgements

Funding: This research was partially supported by the National Key R&D Program of China (2018YFC0116800), National Natural Science Foundation of China (No. 61772110 and 71901011).

## Declaration of Competing Interest

Authors declare that there is no conflict of interest.

## References

- [1] Alexiou, A., Mantzavinos, V.D., Greig, N.H., Kamal, M.A., 2017. A bayesian model for the prediction and early diagnosis of alzheimer's disease. *Frontiers in Aging Neuroscience* 9. doi:https://doi.org/10.3389/fnagi.2017.00077.
- [2] An, Y., Mao, Y., Zhang, L., Jin, B., Xiao, K., Wei, X., Yan, J., 2020. Rahm: Relation augmented hierarchical multi-task learning framework for reasonable medication stocking. *Journal of Biomedical Informatics* 108, 103502. doi:https://doi.org/10.1016/j.jbi.2020.103502.

- [3] Bai, T., Zhang, S., Egleston, B.L., Vucetic, S., 2018. Interpretable representation learning for healthcare via capturing disease progression through time. *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 43–51 doi:https://doi.org/10.1145/3219819.3219904.
- [4] Baytas, I.M., Xiao, C., Zhang, X., Wang, F., Jain, A.K., Zhou, J., 2017. Patient subtyping via time-aware lstm networks, in: *Proceedings of the 23th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, p. 65–74. doi:https://doi.org/10.1145/3097983.3097997.
- [5] Biswal, S., Xiao, C., Glass, L., Milkovits, E., Sun, J., 2020. Doctor2vec: Dynamic doctor representation learning for clinical trial recruitment, in: *AAAI*, pp. 557–564.
- [6] Bridle, J.S., 1990. Probabilistic interpretation of feedforward classification network outputs, with relationships to statistical pattern recognition, in: *Neurocomputing*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 227–236. doi:10.1007/978-3-642-76153-9\_28, .
- [7] Choi, E., Bahadori, M.T., Song, L., Stewart, W.F., Sun, J., 2017. Gram: Graph-based attention model for healthcare representation learning. *Proceedings of the 23th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 787–795 doi:https://doi.org/10.1145/3097983.3098126.
- [8] Choi, E., Bahadori, M.T., Sun, J., Kulas, J.A., Schuetz, A., Stewart, W.F., 2016. Retain: An interpretable predictive model for healthcare using reverse time attention mechanism, in: *Advances in Neural Information Processing Systems 29*, Curran Associates, Inc.. pp. 3504–3512.
- [9] Choi, E., Xiao, C., Stewart, W.F., Sun, J., 2018. Mime: Multilevel medical embedding of electronic health records for predictive healthcare, in: *Advances in Neural Information Processing Systems 31*, Curran Associates, Inc.. pp. 4547–4557.
- [10] Choi, E., Xu, Z., Li, Y., Dusenberry, M.W., Flores, G., Xue, Y., Dai, A.M., 2020. Learning the graphical structure of electronic health records with graph convolutional transformer, in: *AAAI*, pp. 606–613. doi:https://doi.org/10.1609/aaai.v34i01.5400.
- [11] Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2019. BERT: Pre-training of deep bidirectional transformers for language understanding, in: *ACL, Association for Computational Linguistics*. pp. 4171–4186. doi:10.18653/v1/N19-1423.
- [12] Fonarow, G.C., Wright, R.S., Spencer, F.A., Fredrick, P.D., Dong, W., Every, N., French, W.J., 2005. Effect of statin use within the first 24 hours of admission for acute myocardial infarction on early morbidity and mortality. *American Journal of Cardiology* 96, 611–616.
- [13] Harutyunyan, H., Khachatryan, H., Kale, D.C., Galstyan, A., 2019. Multitask learning and benchmarking with clinical time series data. *Scientific Data* 6.
- [14] He, X., Lau, E.H.Y., Wu, P., Deng, X., Wang, J., Hao, X., Lau, Y.C.L., Wong, J.Y., Guan, Y., Tan, X., Mo, X., Chen, Y., Liao, B., Chen, W., Hu, F., Zhang, Q., Zhong, M., Wu, Y., Zhao, L., Zhang, F., Cowling, B.J., Li, F., Leung, G., 2020a. Temporal dynamics in viral shedding and transmissibility of covid-19. *Nature Medicine* 26, 672–675.
- [15] He, Y., Wang, C., Li, N., Zeng, Z., 2020b. Attention and memory-augmented networks for dual-view sequential learning. *SIGKDD*, 125–134 doi:https://doi.org/10.1145/3394486.3403055.
- [16] Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural computation* 9, 1735–1780.
- [17] Hu, Y., Yang, Y., Zhang, J., Cao, X.B., Zhen, X., 2021. Attentional kernel encoding networks for fine-grained visual categorization. *IEEE Transactions on Circuits and Systems for Video Technology* 31, 301–314.
- [18] Jin, B., Yang, H., Sun, L., Liu, C., Qu, Y., Tong, J., 2018a. A treatment engine by predicting next-period prescriptions, in: *SIGKDD, ACM*. pp. 1608–1616.
- [19] Jin, B., Yang, H., Sun, L., Liu, C., Qu, Y., Tong, J., 2018b. A treatment engine by predicting next-period prescriptions, in: *SIGKDD*, p. 1608–1616. doi:https://doi.org/10.1145/3219819.3220095.
- [20] Johnson, A.E.W., Pollard, T.J., Shen, L., Lehman, L.W.H., Feng, M., Ghassemi, M., Moody, B., Szolovits, P., Anthony Celi, L., Mark, R.G., 2016. MIMIC-III, a freely accessible critical care database. *Scientific Data* .
- [21] Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization. *CoRR abs/1412.6980*.
- [22] Le, H., Tran, T., Venkatesh, S., 2018. Dual memory neural computer for asynchronous two-view sequential learning, in: *SIGKDD, ACM*. pp. 1637–1645.
- [23] Luaces, O., Díez, J., Barranquero, J., del Coz, J.J., Bahamonde, A., 2012. Binary relevance efficacy for multilabel classification. *Progress in Artificial Intelligence* 1, 303–313.
- [24] Ma, F., Chitta, R., Zhou, J., You, Q., Sun, T., Gao, J., 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*, p. 1903–1911. doi:https://doi.org/10.1145/3097983.3098088.
- [25] Ma, F., You, Q., Xiao, H., Chitta, R., Zhou, J., Gao, J., 2018. KAME: Knowledge-based attention model for diagnosis prediction in healthcare, in: *CIKM*, pp. 743–752. doi:https://doi.org/10.1145/3269206.3271701.
- [26] Ma, L., Zhang, C., Wang, Y., Ruan, W., Wang, J., Tang, W., Ma, X., Gao, X., Gao, J., 2020. Concare: Personalized clinical feature embedding via capturing the healthcare context. *Proceedings of the AAAI Conference on Artificial Intelligence* 34, 833–840. doi:10.1609/aaai.v34i01.5428.
- [27] Peters, B., Niculae, V., Martins, A.F.T., 2019. Sparse sequence-to-sequence models, in: *ACL*, pp. 1504–1519.
- [28] Pruthi, D., Gupta, M., Dhingra, B., Neubig, G., Lipton, Z.C., 2020. Learning to deceive with attention-based explanations, in: *ACL*, pp. 4782–4793.
- [29] Qiao, Z., Wu, X., Ge, S., Fan, W., 2019. Mnn: Multimodal attentional neural networks for diagnosis prediction, in: *IJCAI*, pp. 5937–5943.
- [30] Shang, J., Ma, T., Xiao, C., Sun, J., 2019a. Pre-training of graph augmented transformers for medication recommendation, in: *IJCAI*, pp. 5953–5959.
- [31] Shang, J., Xiao, C., Ma, T., Li, H., Sun, J., 2019b. Gamenet: Graph augmented memory networks for recommending medication combination, in: *AAAI*, pp. 1126–1133. doi:https://doi.org/10.1609/aaai.v33i01.33011126.
- [32] Shen, C., Qi, G.J., Jiang, R., Jin, Z., Yong, H., Chen, Y., Hua, X., 2019. Sharp attention network via adaptive sampling for person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology* 29, 3016–3027.
- [33] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I., 2017. Attention is all you need, in: *NIPS*, pp. 5998–6008.
- [34] Wang, S., Ren, P., Chen, Z., Ren, Z., Ma, J., de Rijke, M., 2019. Order-free medicine combination prediction with graph convolutional reinforcement learning, in: *CIKM*, pp. 1623–1632.
- [35] Ye, M., Luo, J., Xiao, C., Ma, F., 2020. Lsan: Modeling long-term dependencies and short-term correlations with hierarchical attention for risk prediction. *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 1753–1762 doi:https://doi.org/10.1145/3340531.3411864.
- [36] Zhang, Q., Lu, J., Jin, Y., 2020a. Artificial intelligence in recommender systems. *Complex & Intelligent Systems* 7, 439–457.
- [37] Zhang, X., Xiao, C., Glass, L., Sun, J., 2020b. Deepenroll: Patient-trial matching with deep embedding and entailment prediction, in: *WWW*, pp. 1029–1037.
- [38] Zhang, Y., Chen, R., Jie, T., Stewart, W.F., Sun, J., 2017a. Leap: Learning to prescribe effective and safe treatment combinations for multimorbidity, in: *SIGKDD, ACM*. pp. 1315–1324.
- [39] Zhang, Y., Chen, R., Tang, J., Stewart, W.F., Sun, J., 2017b. Leap: Learning to prescribe effective and safe treatment combinations for multimorbidity, in: *Proceedings of the 23th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, p. 1315–1324. doi:https://doi.org/10.1145/3097983.3098109.
- [40] Zhang, Y., Yang, X., Ivy, J.S., Chi, M., 2019. Attain: Attention-

based time-aware lstm networks for disease progression modeling, in:  
IJCAI, pp. 4369–4375.