# Features Fusion Framework for Multimodal Irregular Time-series Events

Peiwang Tang[1,2] and Xianchao Zhang[3,4(✉)]

[1] Institute of Advanced Technology, University of Science and Technology of China, Hefei 230026, China
`tpw@mail.ustc.edu.cn`
[2] G60 STI Valley Industry & Innovation Institute, Jiaxing University, Jiaxing 314001, China
[3] Key Laboratory of Medical Electronics and Digital Health of Zhejiang Province, Jiaxing University, Jiaxing 314001, China
`zhangxianchao@zjxu.edu.cn`
[4] Engineering Research Center of Intelligent Human Health Situation Awareness of Zhejiang Province, Jiaxing University, Jiaxing 314001, China

**Abstract.** Some data from multiple sources can be modeled as multimodal time-series events which have different sampling frequencies, data compositions, temporal relations and characteristics. Different types of events have complex nonlinear relationships, and the time of each event is irregular. Neither the classical Recurrent Neural Network (RNN) model nor the current state-of-the-art Transformer model can deal with these features well. In this paper, a features fusion framework for multimodal irregular time-series events is proposed based on the Long Short-Term Memory networks (LSTM). Firstly, the complex features are extracted according to the irregular patterns of different events. Secondly, the nonlinear correlation and complex temporal dependencies relationship between complex features are captured and fused into a tensor. Finally, a feature gate are used to control the access frequency of different tensors. Extensive experiments on MIMIC-III dataset demonstrate that the proposed framework significantly outperforms to the existing methods in terms of AUC (the area under Receiver Operating Characteristic curve) and AP (Average Precision).

**Keywords:** Features Fusion · LSTM · Multimodal · Time-series.

## 1 Introduction

In general terms, a modality refers to the way in which something happens or is experienced [2]. To our best knowledge, many existing works have demonstrated that Neural Network can achieve an excellent result in single modality processing such as image classification [23], speech synthesis [13], natural language processing [26]. In the field of data, multimodal is used to represent different forms of data, or different formats of the same form, which generally represents text, picture, audio and video [11,22]. Hence, multimodal data processing have attracted

a wide attention from the academia, especially for multimodal fusion which is one of the original topics in multimodal machine learning [2]. Neural Networks is expected to tackle the multimodal fusion problem [18] and has been used extensively to fuse information for text, image and audio [14, 15], gesture recognition [17], and video or image description generation [21, 27], since the earliest investigation of AVSR [20]. However, almost all these studies focus on text, images or speech modes rather than multimodal time-series which is a critical ingredient across many domains, so how to effectively process multimodal data still need further study. Many methods have been launched to process simple single mode time-series data [1, 29], which have achieved the best result in their respective field. But they have no way to directly use multimodal time-series data, for example multisensor data, medical time-series data.

The problem of features fusion is challenging in multimodal irregular time-series data processing [4]. For example, for clinical data, patient's electronic health records can be abstracted into thousands of interrelated medical events with temporal information, including complex allergy history, family genetic history, drug list, hospitalization records and other historical records. Different event has almost absolutely different frequency of recording. E.g. patient's hospitalization records may be only once a few years, but medication records could be many times a year. Not only different events have different recording frequencies, but also the same type events have significant differences in their different nature. For example, attributes of drug taking events such as drug type, dose and test events include specific indicators and comparison results with normal range values. In order to integrate the features of these events, we must describe these dependencies.

In order to solve the above problems in multimodal irregular time-series events, in this paper, the following contributions is presented in this paper: (1) We propose a new features fusion method to deal with multimodal data, where the features of complex data are fused into a common feature subspace. This method can be applied to different multimodal data. (2) We explore different encoding methods for temporary features, and found a method to embed the temporary features into the non-temporary features, which allows us to better deal with time-series data (3) We propose a model called FG-LSTM which developed from the Recurrent Neural Network such as Phased LSTM [16] to deal with the problem of irregular time-series data. Our proposed model filters the input features by feature gate while recording the complex temporal relationship between different features. (4) We compare with other models, and the experiment results based on the real data demonstrate that the prediction performance of our model is significantly improved.

## 2   Related Work

### 2.1   Multimodal Fusion problem

Multimodal fusion mainly refers to the comprehensive processing of multimodal data by computer, which is responsible for fusing the information of each mode

to perform target prediction [22]. Tensor Fusion Network (TFN) [28]is a multimodal network for features fusion through matrix operation to directly fuse the three features vectors of the data with three modes (such as text, image and audio). However, since TFN calculates the correlation between the elements of different modes through the tensor outer product between modes, it will greatly increase the dimension of features tensor and result in a too large model that is difficult to train. Low-rank Multimodal Fusion [14] uses a low rank matrix to decompose the weight, and hence the TFN process is changed into a single linear transformation of each mode. Then the received multi-dimensional point by Low-rank Multimodal Fusion can be regarded as the sum of multiple low rank vectors, and thus the number of parameters in the model is reduced. Although Low-rank Multimodal Fusion is an upgrade of TFN, once the features are too long, it is still easy to explode parameters. Multimodal Adversarial Representation Network [10] adds a dual discriminator countermeasure network based on multimodal fusion (ordinary attention fusion), which captures dynamic commonness and invariance respectively. Multimodal Bottleneck Transformer [15] uses a shared token between two Transformer, so that this token becomes a communication bottleneck of different modes to save computational attention. In this way, multimodal interaction can be limited to several shared tokens. Compared with the above researches, we pay more attention to multimodal time-series events, and the above researches can also be regarded as special cases of multimodal time-series events.

## 2.2    Time-series Forecasting

Recurrent neural network (RNN) is a neural network used to process sequence data. Theoretically, RNN can store long-term memory and update the previous state according to the current input at any time, but in fact, it is very difficult. In another word, RNN is difficult to solve the problem of long-term dependence [5]. LSTM [6] is a special RNN, which is mainly to solve the problems of gradient disappearance and gradient explosion in the process of long sequence training. Compared with ordinary RNN, LSTM has better performance in long sequences, but LSTM can only maintain a long-term dependence within about 50 time steps. Phased LSTM [16] can solve the problem that LSTM can not process irregular input sequences. By integrating different sampling frequencies or irregularly sampled data on phase gate, Phase LSTM can remember signals with different periods, and the state can propagate for a long time. When the processing sequence reaches thousands of steps, LSTM is almost unavailable, while Phased LSTM performs well. But Phased LSTM is not suitable for modeling the complex event sequence with thousands of event types. HE-LSTM [12] is proposed to deal with heterogeneous temporal events in long-term dependence, but it can only extract event types while the features relationship of events can not be obtained. Transformer [26] is a powerful architecture that can achieve excellent performance on a variety of sequential learning tasks, which does not perform recursion on the sequence, but processes the feedforward model of the whole sequence simultaneously. Recent research shows that transformer has the
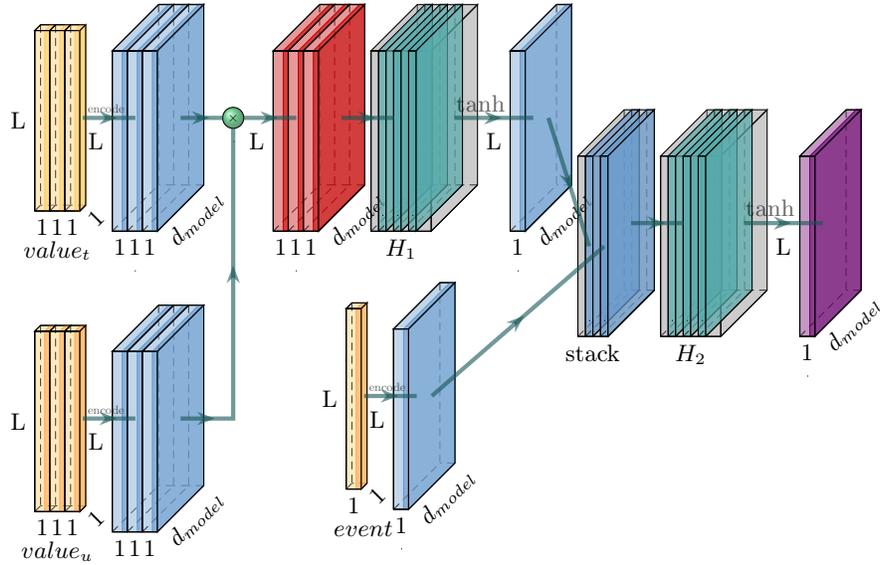
Fig. 1: The non-temporal features fusion method.

potential to improve the prediction ability [24]. However, transformer has some serious problems that make it unable to be directly applied to multimodal irregular time-series data, such as quadratic time complexity, high memory utilization and the inherent limitations of encoder-decoder architecture [29]. In addition to the above problems, the biggest problem of Transformer is that the model contains no recurrence and no convolution, which results in the input tensor can not contain the time relationship of the input sequence effectively [3].

## 3    Methodology

For time-series $\mathbf{S}$ in a given scene, the features of the sequence are consist of dynamic events $\{\mathbf{E}_t\}_{1 \leq t \leq L}$ with length $L$, and each event occurs at the same or different time. We arrange it according to the chronological order of events. Among the events that occur at the same time, the events recorded earlier are arranged in front. Each time-series $\mathbf{S}$ corresponds to a discrete label, which indicates the state of the object at a certain time in the future. For example, in Clinical Endpoint Prediction Task, 0 or 1 indicates the patient's status (death or not) at a certain point in time in the future. So prediction of the results of time-series $\mathbf{S}$ with this corresponding discrete label is defined as the classification of time-series $\mathbf{S}$.

### 3.1   Features fusion

Each time-series $\mathbf{S}$ contains many types of events, and each event has its own time of occurrence. We use $\mathcal{S}$ to represent the feature space where $\mathbf{S}$ is located, $\mathcal{E}$ to represent non-temporal information, that is, the feature space where the events is located, and $\mathcal{T}$ to represent the feature space where temporal information is located. Formally, a sequence $S = [E_1, E_2, \ldots, E_L]$, defines each element as $E_i = (e_i, t_i)$, with $e_i \in \mathcal{E}$ being the non-temporal features at time $i$ and $t_i \in \mathcal{T}$ as an temporal features, and $t_i$ is the interval between the occurrence time of this event and the time when the first event of this time-series occurs. The features vector are defined over a joint space : $\mathcal{S} := (\mathcal{E} \times \mathcal{T})$. The resulting permutation-invariant set is: $\mathbf{S}_E = \{E_1, E_2, \ldots, E_L\} = \{(e_1, t_1), (e_2, t_2), \ldots, (e_L, t_L)\}$. For each event we define $e_i = (type, attribute)$, where $type$ is the type of event, we use $\mathcal{F}_t$ to represent the feature space where type is located; $attribute$ is the attribute of the event, we use $\mathcal{F}_a$ to represent the feature space where attribute is located. So the feature space of event : $\mathcal{E} := (\mathcal{F}_t \times \mathcal{F}_a)$, $\mathcal{E}$ is obviously a joint space, where $\mathcal{F}_t$ is the discrete feature space and $\mathcal{F}_a$ is the continuous feature space. Similarly, attribute consists of two parts: $attribute = (value_t, value_u)$, where $value_t$ is the type of attribute, and $value_u$ is the specific value of attribute. The feature vector of attribute are defined over a joint space : $\mathcal{F}_a := (\mathcal{V}_t \times \mathcal{V}_u)$, where $\mathcal{V}_t$ is the discrete feature space of $value_t$ and $\mathcal{V}_u$ is the continuous feature space of $value_u$.

For each type of event, it can contain multiple types of attributes. While for different types of events, it may contain the same type of attributes or different types of attributes. Therefore, it is difficult to find the feature space of events directly, and we need to characterize the complex relationship between different events. We demonstrate the non-temporal features fusion method as shown in Fig. 1, where $d_{model}$ is the encoded dimension: (1) Select the first three-dimensional feature of attribute, fill up the deficiencies with 0, encode $value_t$ and $value_u$ as $V_t$, $V_u$ respectively, and then use $V_t \times V_u$ to get a new three-dimensional feature; (2) Use $1 \times 1$ convolution kernel to increase the dimension of the features obtained in the previous step, and then use $1 \times 1$ convolution kernel to reduce the dimension to one-dimensional features after being processed by the $tanh$ activation function; (3) Stack the features obtained in the previous step with the features encoded by event, then use $1 \times 1$ convolution kernel to increase the dimension, after processing by the $tanh$ activation function, use $1 \times 1$ convolution kernel to reduce the dimension to obtain the one-dimensional non-temporal features. For $\mathcal{V}_u$ of continuous feature space, we do not simply encode $value_u$ with convolution or fully connected layers, instead encode $value_u$ with the help of $\mathcal{V}_t$ of discrete feature space, as shown in the formula:

$$V_u = W_{v_t} \times value_u + B_{v_t} \tag{1}$$

Where $W_{v_t} \in \mathbb{R}^{L_W \times d_{model}}$ and $B_{v_t} \in \mathbb{R}^{L_B \times d_{model}}$ is the tensor after embedding $value_t$, and $V_u \in \mathbb{R}^{L_V \times d_{model}}$ is the result after encoding $value_u$.

For the fusion of temporal and non-temporal features, many studies directly adopt the additive method, such as the most famous Transformer architecture
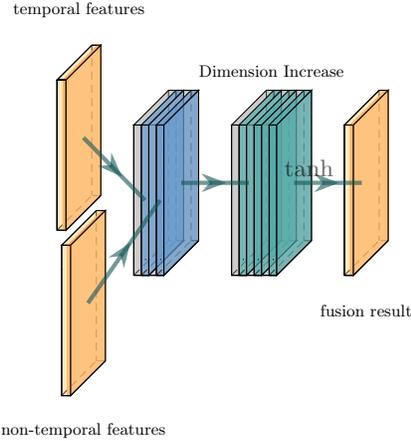
temporal features



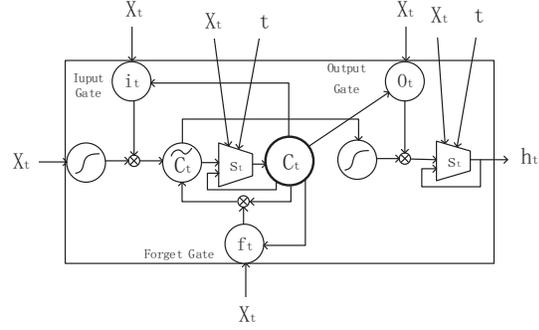Fig. 2: Temporal and non-temporal features fusion methods

Fig. 3: FG-LSTM Model

[26]. The fusion of temporal and non-temporal features is not a simple additive relationship, so the method shown in the Fig. 2 is proposed. Firstly, stack the temporal and non-temporal features, then increase the dimension of the two features used $1 \times 1$ convolution structure. After processing by the *tanh* activation function, we eventually fuse features into one-dimensional tensor on another $1 \times 1$ convolution structure.

Because the time interval between events is not equal, and the time of each event is a very important feature that can not be ignored. We add "time encoding" to the input embeddings and use two methods to encode time:

**Function Encoding** We use sine and cosine functions of different frequencies just as "positional encoding" [26]:

$$FE_{(time,2i)} = sin\left(time/10000^{2i/d_{model}}\right) \tag{2}$$

$$FE_{(time,2i+1)} = cos\left(time/10000^{2i/d_{model}}\right) \tag{3}$$

where $i$ is the dimension, $i \in \{1, \ldots, d_{model}/2\}$. That is, each dimension of the time encoding corresponds to a sinusoid. We chose this function because for any fixed time offset $k$, $FE_{time+k}$ can be represented as a linear function of $FE_{time}$. The time encoding have the same dimension $d_{model}$ as the embeddings, so that the two can be summed.

**Convolution Encoding** We use the convolution structure to learn time encoding:

$$\mathbf{H} = tanh(Conv1d(T^L)) \in \mathbb{R}^{L \times d_{model/2}} \tag{4}$$

$$\mathbf{T} = tanh(Conv1d(\mathbf{H})) \in \mathbb{R}^{L \times d_{model}} \tag{5}$$

For the temporal features with length $L$ and dimension 1, that is, the size $L \times 1$, use the convolution kernel of $1 \times 1$ to learn the matrix with size $1 \times d_{model}/2$, change the temporal features into the matrix with size $L \times d_{model}/2$. After the $tanh$ activation function, use the convolution kernel of $1 \times 1$ to learn the matrix with size $d_{model}/2 \times d_{model}$ again, and change the temporal features into a matrix with the size of $L \times d_{model}$. Finally get the temporal features with the size of $L \times d_{model}$ after $tanh$ activation function.

### 3.2  Model Architecture

Long short-term memory (LSTM) [6] is an important ingredient for modern deep RNN architectures. The FG-LSTM extends the LSTM model by adding a new feature gate $s_t$, and the Fig. 3 shows the FG-LSTM model. The $x_t$ is the input features at time $t$, and others are basically consistent with ordinary LSTM. The feature gate has two factors: a feature filter and a time gate.

The combination of features and time gates only allows the features of certain kinds of features to be input into the neuron, and makes the neuron open only in a specific cycle. This ensures that each neuron will only capture the features of specific types of events and sample them, which solves the problem of poor training effect caused by the complexity and diversity of time and long event sequence.

The opening and closing of this feature gate is controlled by the features and time. Updates to the cell state $c_t$ and $h_t$ are permitted only when the gate is open. We proposed a particularly successful formulation of the feature gate as following:

$$s_t = ReLU(W_{hs}\tanh(W_{xh}x_t + b_h) + b_s) \odot k_t \tag{6}$$

where $W_{xh} \in \mathbb{R}^{d_{model} \times h}, W_{hs} \in \mathbb{R}^{h \times s}$, $b_h \in \mathbb{R}^{1 \times h}$ and $b_s \in \mathbb{R}^{1 \times s}$ are the parameters to be learned, $h$ is hidden size, $s$ is output size. $ReLU$ and $tanh$ is the activation function, $x_t$ is the tensor input at time $t$, and $k_t$ is the time gate [16].

Compared with traditional RNN and other excellent variants of RNN [9], FG-LSTM can choose to update the learned parameters at the time point $t$ of irregular sampling. This allows the FG-LSTM to work with asynchronously sampled irregular time-series data. We can then rewrite the regular LSTM cell update equations for $c_t$ and $h_t$, using proposed cell updates $\tilde{c_t}$ and $\tilde{h_t}$ mediated by the feature gate $s_t$ :

$$i_t = \sigma_i(x_t W_{xi} + h_{t-1} W_{hi} + w_{ci} \odot c_{t-1} + b_i) \tag{7}$$

$$f_t = \sigma_f(x_t W_{xf} + h_{t-1} W_{hf} + w_{cf} \odot c_{t-1} + b_f) \tag{8}$$

$$\tilde{c_t} = f_t \odot c_{t-1} + i_t \odot \tanh_c(x_t W_{xc} + h_{t-1} W_{hc} + b_c) \tag{9}$$

$$c_t = s_t \odot \tilde{c_t} + (1 - s_t) \odot c_{t-1} \tag{10}$$

$$o_t = \sigma_o(x_t W_{xo} + h_{t-1} W_{ho} + w_{co} \odot c_t + b_o) \tag{11}$$

$$\tilde{h_t} = o_t \odot \tanh_h(\tilde{c_t}) \tag{12}$$

$$h_t = s_t \odot \tilde{h_t} + (1 - s_t) \odot h_{t-1} \tag{13}$$

To sum up, for a neuron, only when it meets the type conditions of the corresponding feature gate, and the features information in its sampling period, neron will be updated. Therefore, it can be considered that this neuron represents the state of a certain type of features in a certain sampling period. This is because the feature gate $s_t$, can be seen as a binary classifier to chose the cluster of features types responsible for each neuron. In addition, neurons do not update any information in the closing stage and maintain a perfect memory of past information, i.e. $c_j = c_{j-\Delta}$ if $k_t = 0$ for $t_{j-\Delta} \leq t \leq t_j$. Therefore, other neurons that track other features can directly use the information of this set of features, even if they are far away from each other in sequence indexing. Because of this special mechanism, FG-LSTM can have much diverse and longer memory for modeling the dependency of multiple features.

We use a Softmax layer to predict the true label $\widehat{y_t}$ of the learned features tensor of sequence in the given decision times. This consists of two linear transformations with a $ReLU$ activation in the middle.

$$y_t = softmax(max(0, h_t W_1 + b_1)W_2 + b_2) \tag{14}$$

We use cross-entropy to calculate the classification loss of the prediction $y_t$ and true label $\widehat{y_t}$ of each sample as follows:

$$Loss(\widehat{y_t}, y_t) = \frac{1}{L} \sum_{1 \leq t \leq L} (\widehat{y_t} \times \ln y_t + (1 - \widehat{y_t}) \times \ln(1 - y_t)) \tag{15}$$

We can sum up the losses of all the samples in one minibatch to get the total loss for back propagation.

## 4   Experiments

The dataset used in this experiment is generated by Intensive Care Unit patient medical record data (MIMIC-III) of Beth Israel Deaconess Medical Center in the United States [7]. More than 20000 patient samples in MIMIC-III were ex-

Table 1: The Dataset Distribution

| Dataset | Target 0 | Target 1 | Total |
|---|---|---|---|
| training set | 475291 | 59328 | 534619 |
| validation set | 61698 | 6540 | 68238 |
| evaluation set | 143373 | 19622 | 162995 |

tracted from the dataset, covering more than 4000 kinds and a total of more than 20 million multimodal irregular time-series data. In the experiment, the dataset is divided into training set,validation set and evaluation set, with a ratio of $7 : 1 : 2$. Table 1 shows the data distribution of the dataset, which is divided into two classes. All experiments were implemented by Pytorch [19], optimized by Adam optimization algorithm [8], with the learning rate of 0.0001 and the other parameters are selected as default parameters. We set the random number

Table 2: Results of different non-temporal features fusion methods on different models, among them, different non-temporal feature fusion methods perform the best results, we use bold numbers in black, and underlined numbers are the best results in different models of the same fusion method.

| Model | | LSTM | Bi-LSTM | Phased LSTM | HE-LSTM | Transformer | Informer | FG-LSTM | Count |
|---|---|---|---|---|---|---|---|---|---|
| Our Method | AUC | **75.63** | **75.59** | **72.52** | 74.21 | **75.69** | **76.05** | <u>**78.85**</u> | 12 |
| | AP | **34.96** | **34.92** | **30.45** | 32.44 | **34.31** | **34.93** | <u>**38.90**</u> | |
| Other Method | AUC | 68.59 | 70.36 | 68.95 | **76.35** | 64.23 | 75.91 | <u>76.37</u> | 2 |
| | AP | 25.94 | 26.85 | 26.63 | **34.80** | 21.71 | 32.85 | <u>36.42</u> | |
| Count | | 0 | 0 | 0 | 0 | 0 | 0 | <u>4</u> | - |

seed to 1 to ensure the repeatability of the experimental results. Unless otherwise specified, $d_{model}$ (the dimension after features coding) is 256, the batchsize is 128. The detailed parameter settings of different experiments are described below. All the experiments are conducted on a single Nvidia RTX 3090 GPU (24GB memory), which is sufficient for all the baselines.

### 4.1 Evaluating Metrics

AUC (the area under Receiver Operating Characteristic curve) and AP (Average Precision) [25] are uesd in this paper. AUC is the area of ROC curve and the x-axis, and AP is the area of PRC (precision recall curve) and the x-axis, both of which are robust to the imbalanced data of positive and negative samples.

### 4.2 Comparing Methods

Because the proposed FG-LSTM is a variant based on the classical LSTM [6], we choose the classical LSTM and three other excellent variants including BI-LSTM, Phase LSTM [16] and HE-LSTM [12]. Recently, Transformer architecture has achieved the best performance in many problems, so we discuss the ability of Transformer related architecture to deal with multimodal irregular time-series. We chose the vanilla Transformer [26] and further select one of excellent variants in it called Informer [29]. Because our experiment does not involve the generation process, therefore, only the encoder part of the Transformer architecture is used, and get the final output directly through a fully connected feed-forward network. For LSTM related architectures, only use one layer. For Informer, the number of layers in the original author's open source code is selected, that is, $n = 2$. For Transformer, in order to better compare with Informer, we selecte the same encoder layers as Informer. In addition, $d_{model}$ is changed to 256, which is consistent with LSTM architecture, and there is no change in the parameter settings of Transformer related architecture.

### 4.3 Experimental Result

**Non-temporal Features fusion methods** In many previous studies, the processing methods of features from different feature spaces are only simple addition.

According to this idea, a method is proposed as a comparative experiment, as shown below:

$$x = V_e + sum\left(V_t \times V_u\right) \tag{16}$$

Where $V_e \in \mathbb{R}^{L_e \times d_{model}}$, $V_t \in \mathbb{R}^{L_t \times d_{model}}$ and $V_u \in \mathbb{R}^{L_u \times d_{model}}$ are the tensor encoded by $event$, $value_t$ and $value_u$ respectively. In this experiment, the coding method without considering the temporary features. We uniformly choose the (2) (3) proposed above. For the fusion method of temporal features and non-temporal features, the addition method is directly selected, and the rest are discussed in detail below.
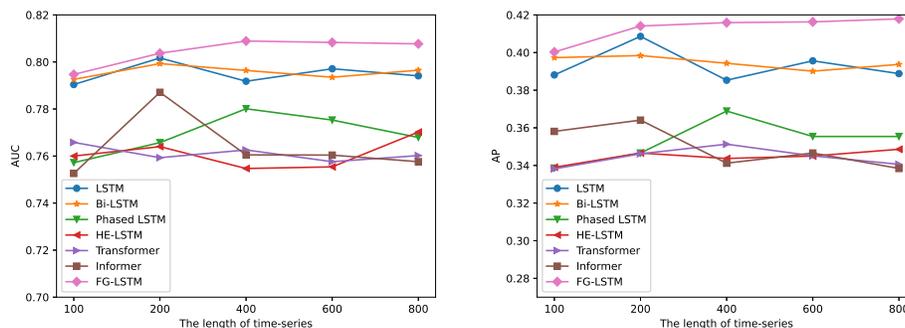
Table 2 shows the experimental results of AUC and AP on Table 1 dataset with different model architectures and different non-temporal features fusion methods. It is obvious that, compared with the common methods, the proposed method of non-temporary features fusion has better performance. Except for the best performance in HE-LSTM framework, our proposed method has advantages in all other frameworks. The most obvious improvement is the Transformer framework, which has increased by 17.84% in AUC and 58.03% in AP. However, for the excellent Informer framework proposed for single-modal time-series, the improvement is not very obvious. The AUC and AP have only increased by 0.18% and 6.33% respectively, which shows that the Informer framework is not very sensitive to feature fusion methods. If we do not pay much attention to features fusion methods, Informer framework is indeed a good choice. For our proposed model FG-LSTM, the best performance of all models is obtained in different non-temporal feature fusion methods, and the AUC and AP are also improved by 3.24% and 6.80% respectively. Although the improvement is not very obvious, it also proves the superiority of our proposed model itself. In general, different feature fusion methods have great impact on the performance of different models, but excellent models are not particularly sensitive to feature fusion methods.

**Temporal and non-temporal features fusion methods** The advanced of the proposed non-temporal features fusion method has been proved. Therefore, in this experiment, we verify the progressiveness of our proposed temporal and non-temporal features fusion method. In order to explore whether it is necessary to upgrade the dimension, we set up a group of control experiments to fuse the two features after stack directly with the help of $1 \times 1$ convolution kernel.

Table 3 shows our experimental results. Where **FE** is function encoding, **CE** is convolution encoding, **add** is a direct addition method, **conv-add** is our method, and $conv-add$ is a comparative method without dimension upgrading. For the **FE** method without learning parameters, it can be seen that the **conv-add** method has achieved the most advanced experimental results in different models, while the $conv-add$ method without dimension upgrading is not as good as the **conv-add** method. But it is still better than the direct addition method in many models. For the **CE** method of learning parameters, it can be seen that no matter what kind of temporal and non-temporal feature fusion method, **CE** is better than **FE**, but none of the three feature fusion methods

Table 3: Results of different temporal features fusion methods on different models.

| Model | FE | | | | | | CE | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | add | | conv-add | | $conv-add$ | | add | | conv-add | | $conv-add$ | |
| | AUC | AP | AUC | AP | AUC | AP | AUC | AP | AUC | AP | AUC | AP |
| LSTM | 75.63 | 34.96 | *76.88* | *36.38* | 76.87 | 36.35 | **79.47** | **39.50** | 79.09 | 38.56 | 79.18 | 38.53 |
| Bi-LSTM | 75.59 | 34.92 | *78.28* | *37.67* | 76.98 | 37.46 | 79.04 | 38.56 | 77.95 | 36.06 | **79.64** | **39.43** |
| Phased LSTM | 72.52 | 30.45 | *74.82* | *33.57* | 72.83 | 30.69 | 76.34 | 34.03 | 74.54 | 32.90 | **78.01** | **36.89** |
| HE-LSTM | 74.21 | 32.44 | *76.78* | *34.73* | 75.46 | 33.24 | **77.06** | **35.86** | 76.79 | 34.72 | 75.47 | 34.37 |
| Transformer | 75.69 | 34.31 | *76.14* | *34.88* | 75.85 | 34.87 | **76.41** | **35.19** | 75.65 | 34.22 | 76.26 | 35.13 |
| Informer | 76.05 | 34.93 | *76.19* | *35.66* | 75.96 | 34.86 | **78.11** | **35.29** | 75.99 | 34.90 | 76.05 | 34.12 |
| FG-LSTM | 78.85 | 38.90 | *80.67* | *41.94* | 78.47 | 38.01 | 79.33 | 39.09 | **81.20** | **42.69** | 80.89 | 41.59 |
| Count | 0 | | 0 | | 0 | | 8 | | 2 | | 4 | |



(a) The AUC of different models          (b) The AP of different models

Fig. 4: The performance of different models on different time-series length.

always has best performance in all models. Because our upgraded **conv-add** method also has parameters to learn, we believe that as long as the dataset is larger, the upgraded **conv-add** method can still be better than other methods in different models. Finally, for different models, different time coding methods and different feature fusion methods are used. Our FG-LSTM model is better than other models, which is enough to prove the robustness of our FG-LSTM. It also shows that the variants of LSTM are not necessarily inferior to the models of Transformer series.

**Experimental comparison of different length time-series** In order to verify the proposed model in this paper has stronger ability to capture the temporal dependence between features than other models, in this experiment, different models are input with different lengths of time-series data, ranging from 100 to 800. For the temporal feature coding method, choose **CE**. For the non-temporal

feature fusion method, use our own method. For the temporal and non-temporal feature fusion method, choose the $conv-add$ method without dimension upgrading. Fig. 4(a) and Fig. 4(b) show the results of this experiment, for Transformer, when the length of time-series is 600 and 800, the Transformer failure for the out-of-memory, so we set the $batchsize$ is 64 to make 24GB memory enough. From the experimental results, we can draw the following conclusions:

Firstly, the time-series information is effective for the prediction results. When the input length is less than 400, most models will be improved with the increase of the length of the input sequence. Secondly, compared with other models, FG-LSTM is better at capturing the timing dependency in time-series. When the sequence length exceeds 400 and becomes longer and longer, the performance of the model is not improved much in AUC, but the AP is still improved steadily. However, other models can not capture the timing dependence under ultra long sequences, so they have not been greatly improved, and even the effect has become worse. Finally, we can see that the classical LSTM model is superior to Transformer and its variant model Informer, which shows that the time-series information extraction of Transformer series models is still slightly insufficient.

## 5    Conclusion

This paper proposes a features fusion framework and FG-LSTM model updated on the basis of LSTM. The model can well deal with multimodal irregular time-series data. At the same time, we also explore how to better encode time features and how to better integrate temporal features and non-temporal features, which is particularly important for irregular time-series data. Firstly, through the temporal features coding method and features fusion framework, the representation tensor obtained by the model can fuse the features and temporal dependency between different non-temporal information, effectively capture the temporal dependency under ultra long sequences and the feature information of a minority events. Then, input the representation tensor of the obtained time-series into the FG-LSTM, due to the existence of feature gates, the model can automatically adapt to the multi-scale sampling frequency of multi-source complex data, asynchronously track the temporal information and feature information of different events. Finally, the experiments demonstrate that the method proposed in this paper has better performance than other typical methods on real datasets. The method in this paper is promising to expand and popularize, and can be further migrated to diverse fields, especially for multi-source asynchronous sampling sensor data and behavior recording data.

## References

1. Armandpour, M., Kidd, B., Du, Y., Huang, J.Z.: Deep personalized glucose level forecasting using attention-based recurrent neural networks. arXiv preprint arXiv:2106.00884 (2021)

2. Baltrušaitis, T., Ahuja, C., Morency, L.P.: Multimodal machine learning: A survey and taxonomy. IEEE transactions on pattern analysis and machine intelligence **41**(2), 423–443 (2018)
3. Dai, Z., Yang, Z., Yang, Y., Carbonell, J.G., Le, Q.V., Salakhutdinov, R.: Transformer-xl: Attentive language models beyond a fixed-length context. In: ACL (1) (2019)
4. Fu, Y., Cao, L., Guo, G., Huang, T.S.: Multiple feature fusion by subspace learning. In: Proceedings of the 2008 international conference on Content-based image and video retrieval. pp. 127–134 (2008)
5. Hochreiter, S., Bengio, Y., Frasconi, P., Schmidhuber, J., et al.: Gradient flow in recurrent nets: the difficulty of learning long-term dependencies (2001)
6. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural computation **9**(8), 1735–1780 (1997)
7. Johnson, A.E., Pollard, T.J., Shen, L., Li-Wei, H.L., Feng, M., Ghassemi, M., Moody, B., Szolovits, P., Celi, L.A., Mark, R.G.: Mimic-iii, a freely accessible critical care database. Scientific data **3**(1), 1–9 (2016)
8. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. Computer Science (2014)
9. Koutnik, J., Greff, K., Gomez, F., Schmidhuber, J.: A clockwork rnn. In: International Conference on Machine Learning. pp. 1863–1871. PMLR (2014)
10. Li, X., Wang, C., Tan, J., Zeng, X., Ou, D., Ou, D., Zheng, B.: Adversarial multi-modal representation learning for click-through rate prediction. In: Proceedings of The Web Conference 2020. pp. 827–836 (2020)
11. Liu, J., Li, T., Xie, P., Du, S., Teng, F., Yang, X.: Urban big data fusion based on deep learning: An overview. Information Fusion **53**, 123–133 (2020)
12. Liu, L., Shen, J., Zhang, M., Wang, Z., Liu, Z.: Deep learning based patient representation learning framework of heterogeneous temporal events data. Big Data Research **5**(1), 2019003 (2019)
13. Liu, P., Cao, Y., Liu, S., Hu, N., Li, G., Weng, C., Su, D.: Vara-tts: Non-autoregressive text-to-speech synthesis based on very deep vae with residual attention. arXiv preprint arXiv:2102.06431 (2021)
14. Liu, Z., Shen, Y., Lakshminarasimhan, V.B., Liang, P.P., Zadeh, A.B., Morency, L.P.: Efficient low-rank multimodal fusion with modality-specific factors. In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (2018)
15. Nagrani, A., Yang, S., Arnab, A., Jansen, A., Schmid, C., Sun, C.: Attention bottlenecks for multimodal fusion. Advances in Neural Information Processing Systems **34** (2021)
16. Neil, D., Pfeiffer, M., Liu, S.C.: Phased lstm: Accelerating recurrent network training for long or event-based sequences. In: NIPS (2016)
17. Neverova, N., Wolf, C., Taylor, G., Nebout, F.: Moddrop: adaptive multi-modal gesture recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence **38**(8), 1692–1706 (2015)
18. Ngiam, J., Khosla, A., Kim, M., Nam, J., Lee, H., Ng, A.Y.: Multimodal deep learning. In: ICML (2011)
19. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems **32**, 8026–8037 (2019)

20. Potamianos, G., Neti, C., Gravier, G., Garg, A., Senior, A.W.: Recent advances in the automatic recognition of audiovisual speech. Proceedings of the IEEE **91**(9), 1306–1326 (2003)
21. Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., Chen, M., Sutskever, I.: Zero-shot text-to-image generation. arXiv preprint arXiv:2102.12092 (2021)
22. REN, Z., WANG, Z., KE, Z., LI, Z., Wushour·Silamu: Survey of multimodal data fusion. Computer Engineering and Applications **57**(18),  16 (2021)
23. Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: International Conference on Machine Learning. pp. 6105–6114. PMLR (2019)
24. Tsai, Y.H.H., Bai, S., Yamada, M., Morency, L.P., Salakhutdinov, R.: Transformer dissection: An unified understanding for transformer's attention via the lens of kernel. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 4344–4353 (2019)
25. Turpin, A., Scholer, F.: User performance versus precision measures for simple search tasks. In: Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval. pp. 11–18 (2006)
26. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: Advances in neural information processing systems. pp. 5998–6008 (2017)
27. Wu, C., Liang, J., Ji, L., Yang, F., Fang, Y., Jiang, D., Duan, N.: Nüwa: Visual synthesis pre-training for neural visual world creation. arXiv preprint arXiv:2111.12417 (2021)
28. Zadeh, A., Chen, M., Poria, S., Cambria, E., Morency, L.P.: Tensor fusion network for multimodal sentiment analysis. arXiv preprint arXiv:1707.07250 (2017)
29. Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., Zhang, W.: Informer: Beyond efficient transformer for long sequence time-series forecasting. In: Proceedings of AAAI (2021)