

Modelling Irregularly Sampled Time Series Without Imputation

Rohit Agarwal¹, Aman Sinha^{2,3}, Dilip K. Prasad¹, Marianne Clausel², Alexander Horsch¹, Mathieu Constant², Xavier Coubez³

¹Bio-AI Lab, UiT The Arctic University of Norway, Tromsø

²Université de Lorraine, Nancy, France

³Institut de Cancérologie Strasbourg, France
{agarwal.102497, amansinha091}@gmail.com

Abstract

Modelling irregularly-sampled time series (ISTS) is challenging because of missing values. Most existing methods focus on handling ISTS by converting irregularly sampled data into regularly sampled data via imputation. These models assume an underlying missing mechanism leading to unwanted bias and sub-optimal performance. We present SLAN (Switch LSTM Aggregate Network), which utilizes a pack of LSTMs to model ISTS without imputation, eliminating the assumption of any underlying process. It dynamically adapts its architecture on the fly based on the measured sensors. SLAN exploits the irregularity information to capture each sensor’s local summary explicitly and maintains a global summary state throughout the observational period. We demonstrate the efficacy of SLAN on publicly available datasets, namely, MIMIC-III, Physionet 2012 and Physionet 2019. The code is available at <https://github.com/Rohit102497/SLAN>.

Introduction

An irregularly sampled time series (ISTS) is a multivariate time series recorded at inconsistent or non-uniform time intervals. Such data can be found in various fields dealing with complex generative processes, such as meteorology (Mudelsee 2002), seismology (Ravuri et al. 2021), user social-media activity logs (Zeng and Gao 2022), e-commerce transactions (Wu, Hernández-Lobato, and Ghahramani 2013), epidemiological and clinical research (Shrive et al. 2006; Yadav et al. 2018). The missingness in ISTS data can be categorized (Rubin 1976) as missing completely at random, missing at random and missing not at random (MNAR) (Ma and Zhang 2021). In contrast to the first two categories, the cause of the missingness in MNAR relates to unobserved data. Thus, modelling applications concerning ISTS with MNAR data is challenging. In this paper, we model ISTS with MNAR. We refer to ISTS with MNAR as only ISTS for ease of reading. Fig. 2b shows a pictorial representation of the ISTS data.

Many methods handle ISTS by filling missingness via imputation converting ISTS to regularly sampled time series (Fig. 2a), assuming an underlying missing mechanism (Ipsen, Mattei, and Frellsen 2020). Imputation is the process of filling up missing values with estimated values whenever input is not observed. The imputation can be performed via forward filling, mean (Che et al. 2016), interpolation (Shukla

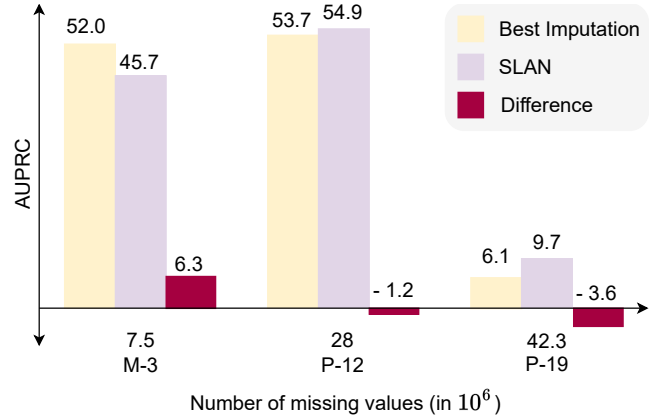


Figure 1: Best imputation model (GRU-D for M-3 and P-12 and GRU-Simple for P-19) vs our model (SLAN) performance. As the number of missing values increases (from M-3 to P-19), the performance of the best imputation model compared to SLAN deteriorates. The negative sign in the difference bar denotes the performance of SLAN is better.

and Marlin 2018), model-based technique (Lim et al. 2021), etc. However, any form of imputation alters the original nature of the data with artificial and sub-optimal approximation leading to unwanted distribution shifts (Zhang et al. 2023). This is evident from Fig. 1 where the best imputation-based ISTS model’s performance deteriorates compared to our model (non-imputation based) as the number of imputed observations increases.

We argue that learning the imputation task is challenging because of the underlying missing mechanism and is not required for the downstream task. Some non-imputation methods (Horn et al. 2020; Vaswani et al. 2017) exist in the literature but do not properly exploit the temporal structure of missing values. It is worth noting that MNAR datasets are not simply incomplete but rather contain informative missingness (Rubin 1976). Therefore, specialized methods are required to handle such missingness in a meaningful way.

We present *Switch LSTM Aggregate Network* (SLAN), which utilizes a pack of LSTM to handle multivariate ISTS. Our proposed model exploits the irregularity information of the ISTS data to maintain the local summary state for each

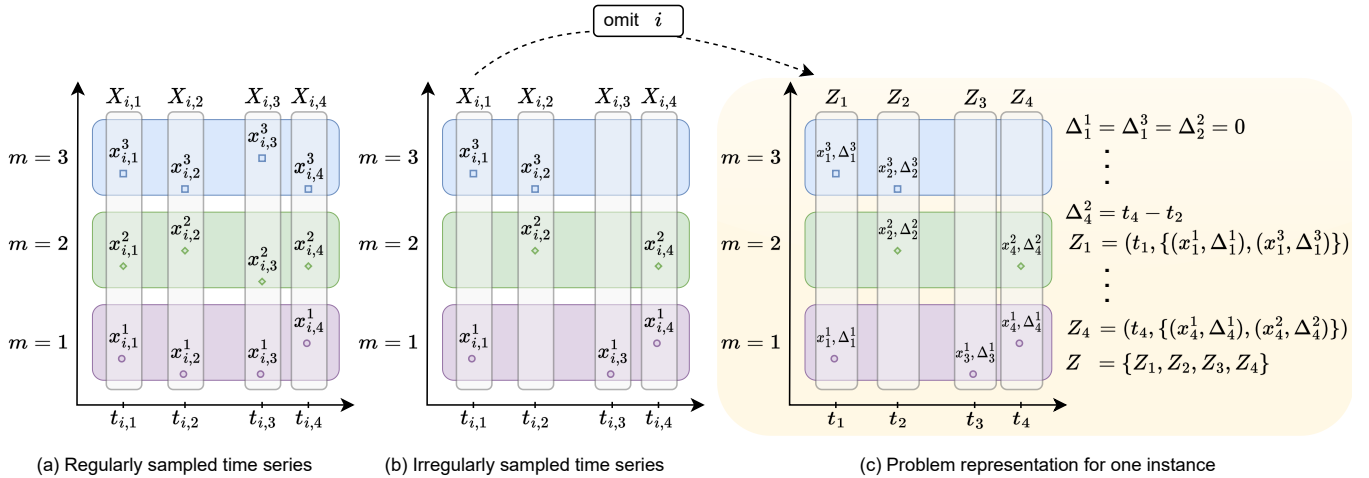


Figure 2: (a) A snapshot of multi-variate regularly sampled time series for i^{th} instance. m represents the index of the sensor. (b) A snapshot of multi-variate irregularly sampled time series (ISTS) for i^{th} instance. (c) Problem representation of the ISTS with respect to one instance by omitting the subscript i . (Best viewed in colour)

observed time series. In most practical situations, these time series are coming from sensor measurement. The model is equipped with a switch layer that enables it to adapt to the input order of measured sensors dynamically. SLAN maintains a global summary state aiding each LSTM with summarised information throughout the observational period. Since the clinical domain data exhibits the characteristics of ISTS with MNAR, we show the efficacy of our model on three such publicly available datasets: MIMIC-III (Johnson et al. 2016), PhysioNet 2012 (Goldberger et al. 2000) and PhysioNet 2019 (Reyna et al. 2019). The main contributions of our work are: (1) SLAN eliminates the need for missing value imputation in ISTS data. (2) A simple switch layer is introduced to guide the activation of appropriate LSTM blocks. (3) We improve the decay function in LSTMs by introducing a sampling rate. (4) We employ an aggregate function to maintain the summary state of the model.

Related Works

Imputation-Based Models Previous works such as (Marlin et al. 2012; Harutyunyan et al. 2017; Rasmussen and Williams 2003) and IP-Nets (Shukla and Marlin 2018) propose to convert ISTS to regularly sampled time series by temporal discretization. They perform imputation using a probabilistic model (Marlin et al. 2012) or interpolation (IP-Nets). The main drawback of these methods is that they require either ad-hoc choices such as width and the aggregation function (Marlin et al. 2012) or a predefined nonlinear form assuming that missingness is at random (IP-Nets), thus inducing bias. GRU-based (Che et al. 2016) methods like GRU-mean and GRU-forward directly impute input by the global mean value and the last measured value, respectively, without any temporal discretization. GRU-simple (Che et al. 2016) and Phased-LSTM (Neil, Pfeiffer, and Liu 2016) go one step ahead and use the missing values indicator as ad-

ditional information along with the imputed input for RNN-based models. These methods lead to a potential distribution shift. Unlike above, other imputation methods use a model-based approach to infer the missing values. GRU-D (Che et al. 2016) imputes the missing input by employing learnable decay on the global mean and the last measured value. Overall, these models perform imputation, assuming an underlying missing mechanism leading to unwanted bias. Albeit they perform well, this does not justify making use of biased data. A detailed discussion on various imputation-based models is presented in *supplementary S1*.

Non-Imputation Models Recent approaches have also explored learning directly from the ISTS data without any form of imputation. Transformer (Vaswani et al. 2017) based models have been widely used for modelling time series data, including ISTS. These approaches mainly replace the positional encoding with an encoding of time and model sequences using self-attention and concatenate it with the input representation. The main drawback of these methods is the permutation-invariant nature of self-attention, which can be problematic in capturing dependencies within each time series. Another work, SeFT (Horn et al. 2020) proposed to learn from irregularly-sampled data via set-based data representation, treating time series as an unordered set of measurements. However, the order-invariant nature of set representation fails to capture the irregularity information, which is order-variant and increases with time. Raindrop (Zhang et al. 2021) proposes a graph neural network to learn a sensor dependency graph. Raindrop leverages inter-sensor dependency to train latent embeddings. However, Raindrop does not exploit the irregularity information of the sensors. This is one of the main motivations of our model, as we not only exploit the irregularity explicitly but also preserve the temporal information and maintain a global summary state throughout the observational period. More information on

the non-imputation models is given in *supplementary S1*.

Problem Formulation

Regularly Sampled Time Series Consider a dataset represented by $D = \{X, Y\}$, where X is a set of instances given by $X = \{X_1, \dots, X_n\}$, Y is the set of label given by $Y = \{y_1, \dots, y_n\}$ and n is the total number of instances. X_i is a time series for i^{th} instance given by $X_i = \{X_{i,1}, \dots, X_{i,l_i}\}$ where l_i is the number of time steps i^{th} instance was measured. $X_{i,j}$ is the set of measured values of all sensors at time $t_{i,j}$, given by $X_{i,j} = \{x_{i,j}^1, \dots, x_{i,j}^s\}$. Here, $x_{i,j}^m$ represents the measured value of sensor m for i^{th} instance at time $t_{i,j}$ and s is the total number of sensors/features. We represent all sensors by their indices and the set of indices of sensors is given by $\mathbb{M} = \{1, \dots, s\}$. We present a snapshot of multi-variate regularly sampled time series data of i^{th} instance in Fig. 2a considering $s = 3$ and $l_i = 4$. Note that $t_{i,2} - t_{i,1}$ need not be equal to $t_{i,3} - t_{i,2}$.

Irregularly Sampled Time Series (ISTS) ISTS follows the definition of regularly sampled time series, except not all sensors will be measured at each time step leading to an irregular sampling of each sensor. This is mathematically given as $X_{i,j} \subseteq \{x_{i,j}^1, \dots, x_{i,j}^s\}$. A snapshot of ISTS for i^{th} instance is shown in Fig. 2b where $X_{i,1} = \{x_{i,1}^1, x_{i,1}^3\}$, $X_{i,2} = \{x_{i,2}^2, x_{i,2}^3\}$, $X_{i,3} = \{x_{i,3}^1\}$ and $X_{i,4} = \{x_{i,4}^1, x_{i,4}^2\}$.

Problem Representation For simplicity, we omit the subscript i representing an instance and consider only one instance to discuss the problem and the working of the proposed model. In that sense, each measured value is given by x_j^m (instead of $x_{i,j}^m$) and it is received at time t_j (instead of $t_{i,j}$). Let us denote this instance by $Z = X_i$ and the set of values of measured sensors at time t_j by Z_j . Based on this, at each time step t_j , we represent the measured sensors as $Z_j = (t_j, \bigcup_{m \in \mathbb{A}_j} \{(x_j^m, \Delta_j^m)\})$ where \mathbb{A}_j is the set of sensors measured at time t_j , given by $\mathbb{A}_j \subseteq \mathbb{M}$. Δ_j^m denotes the time delay between two successive values measured by sensor m , i.e., $\Delta_j^m = t_j - t_k$, where $k = \max(1, \dots, j-1)$ such that $m \in \mathbb{A}_k$. Thus the whole input data is given by $Z = \bigcup_{j=1}^l \{Z_j\}$ where l is the number of time steps. Following the previous paragraph, we visually present the data representation and the corresponding equations in Fig. 2c.

SLAN: Switch LSTM Aggregation Network

SLAN is an adaptive LSTM-based model that dynamically changes its architecture depending on the measured sensors at any time point by utilizing a switch layer. The architecture of SLAN is presented in Fig. 3a. It consists of a pack of LSTMs such that there is a one-on-one connection between a sensor and an LSTM block. This is facilitated by the switch layer (see the yellow-colored box in Fig. 3a). Each sensor is connected to its corresponding LSTM block by a switch. A switch goes "on" if its corresponding sensor is measured; otherwise, it stays off. The "on" switch results in activating its corresponding LSTM block, thus, eliminating the need for any imputation. The LSTM block outputs a long-term memory (LTM) and a short-term memory (STM) (Fig. 3b).

Algorithm 1: Switch LSTM Aggregation Network

Require: Model M with s LSTM block, switch layer and aggregation function as shown in Fig. 3a
for i in instances **do**
 Initialize h_0^m and c_0 for each m
 while time t_j **do**
 Receive measured sensors input Z_j (see Fig. 2c)
 Create switch layer \mathbb{S}_j from measured sensors \mathbb{A}_j
 Activate LSTM blocks corresponding \mathbb{S}_j
 for L^m in active LSTM blocks **do**
 Calculate h_j^m, c_j^m by eq. 3
 end for
 Calculate c_j using aggregation function by eq. 2
 end while
 Decay final hidden states ($h_{t_j^*}^m$) of each sensor (eq. 5)
 Concat all decayed final hidden states ($\tilde{h}_{t_j^*}^m$) and final summary state (c_i) to get concat layer C
 Predict \hat{y} using eq. 4
 Update M based on the loss
end for

The LTM of each activated LSTM block is aggregated to produce a global summary state and passed on to all the LSTM blocks for the next time step as input. This aids the LSTM blocks with summarised information. The previous short-term memory (STM) of each activated LSTM block is decayed based on the time delay and decay function (Fig. 3c) and passed as an input for the next measured time point. This acts as a local summary for each sensor. Finally, at the last time point, the STM from each LSTM block is decayed and concatenated with the aggregated LTM as seen in the concat layer in Fig. 3a. The concat layer is then fully connected to a 2-node output layer for binary classification. The unrolled SLAN architecture for the data example given in Fig. 2c is presented in Fig. 3e and discussed in detail in *supplementary S2*. The pseudo-code of the SLAN architecture is given in Algorithm 1.

Architecture

SLAN consists of s LSTM blocks $\{L^1, \dots, L^s\}$ where L^m is associated with sensor m . We define the switch layer (\mathbb{S}_j) as the set of switches kept "on" based on the measured sensors at time t_j . Since there is a one-on-one correspondence between a switch and its corresponding measured sensor, we borrow the representation of \mathbb{S}_j as the indices of the sensors measured at time t_j from the *Problem Formulation* section, thus $\mathbb{S}_j = \mathbb{A}_j$. Each active LSTM block (L^m) at time t_j takes the sensor value (x_j^m), STM (h_{j-1}^m), LTM (c_{j-1}^m) and time delay (Δ_j^m) as inputs and outputs h_j^m and c_j^m , given by

$$(h_j^m, c_j^m) = L^m(x_j^m, h_{j-1}^m, c_{j-1}^m, \Delta_j^m) \quad \forall m \in \mathbb{S}_j \quad (1)$$

where $L^m \forall m \in \mathbb{A}_j$ are active based on \mathbb{S}_j at time t_j . An aggregate function is employed on the LTM of the active LSTM blocks to get a summary state (c_j) at t_j . Any function that can group multiple values to give a single summary value can be used as an aggregation function and is represented by $agg()$. Some examples of aggregation functions

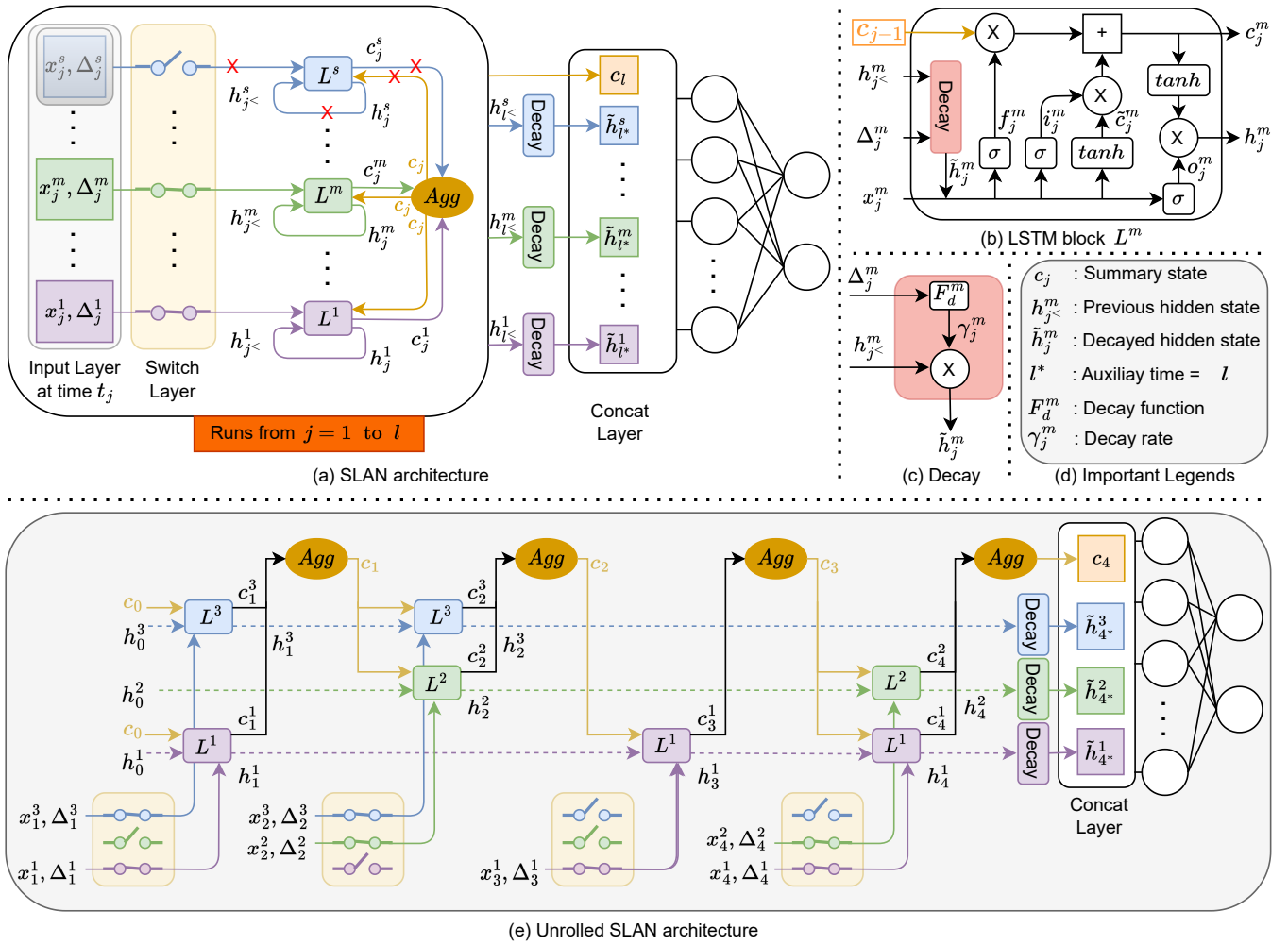


Figure 3: (a) Snapshot of SLAN Architecture. Here $x_j^m, \Delta_j^m, h_j^m, c_j^m$ denotes the input, time delay, short-term memory, and long-term memory at time t_j of m^{th} sensor and LSTM block. The closed circuit in the switch layer means that a particular switch is "on", otherwise it is "off". The X sign in red implies that there is no input or output to the corresponding LSTM block. Agg denotes the aggregate function. (b) The inner working of the LSTM block is given here. (c) The decay of the previous hidden state is shown here. (d) Important notations are denoted in the legend. (e) Unrolled SLAN architecture based on the example from Figure 2c is presented here. (Best viewed in colour)

are mean, max and attention. The summary state is given as

$$c_j = \text{agg}\left(\bigcup_{m \in \mathbb{S}_j} \{c_j^m\}\right) \quad (2)$$

The c_j is used as an input for the next time step for every active LSTM block. An active LSTM at t_j might not be active at t_{j-1} . Thus, the STM input to L^m at t_j is represented by $h_{j^<}^m$ (instead of h_{j-1}^m) where $j^< = \max(1, \dots, j-1)$ such that $m \in \mathbb{S}_{j^<}$. Thus, eq. 1 can be updated as:

$$(h_j^m, c_j^m) = L^m(x_j^m, h_{j^<}^m, c_{j-1}, \Delta_j^m) \quad \forall m \in \mathbb{S}_j \quad (3)$$

Finally, the hidden states (or STM) of all the LSTM block and the summary state is concatenated to give a final output. The hidden states are decayed before concatenating to ensure that the contribution of each hidden state is accounted for based on the last time it was measured. The decayed

hidden state at t_j is represented by \tilde{h}_j^m and is given by the eq. 5. For simplicity, we introduce an auxiliary time t_{l^*} which is equal to t_l . All the hidden states are decayed (represented by $\tilde{h}_{l^*}^m$) based on $\Delta_{l^*}^m$ where $\Delta_{l^*}^m = t_l - t_k$ and $k = \max(1, \dots, l)$ such that $m \in \mathbb{A}_k$. Thus the concat layer is given by $C = \{c_l, \tilde{h}_{l^*}^1, \dots, \tilde{h}_{l^*}^s\}$. A fully-connected network is employed to get a final prediction from C as follows

$$\hat{y} = F(C) \quad (4)$$

Hidden State Decay Since the measurement of each sensor is irregular, we employ a time-decay function on the hidden states inspired by (Che et al. 2016). The time-decay function can be any monotonic non-increasing function F_d^m (e.g. e^{-z}), ensuring that the previous local summary is minimized if the sensor is measured after a long time interval

Table 1: Description of the MIMIC-III, Physionet 2012 and Physionet 2019 (early sepsis prediction) datasets. #Instances denotes the number of patient records in the datasets, #Sensors denotes the number of features/sensors in each instance, #Observations denotes the average number of observations recorded in each instance, i.e., the number of time steps, Measured ratio denotes the ratio of the total number of data measured by the total number of data measured in a fully available dataset, #Num-Imputation is the number of imputation or missing values and Imbalance denotes the percentage of instances with a minority class label.

Dataset	#Instances	#Sensors	#Observations(avg.)	Measured ratio (%)	#Num-Imputation	Imbalance (%)
MIMIC-III	21110	17	77.7	32.5	7.5×10^6	13.22
Physionet 2012	11988	37	74.9	11.6	28.0×10^6	14.24
Physionet 2019 (ESP)	40333	34	38.5	19.8	42.3×10^6	1.79

i.e., a large value of Δ_j^m . Since each sensor is different, the decaying function should differ for each. Thus, a trainable time-decay function is employed. Moreover, each sensor has a different sampling rate, defined as the number of diagnoses of each sensor per hour. Thus, we introduce a sampling rate in the decay function to inculcate the different decay intervals of sensors. The sampling rate of sensor m is represented by r^m , and the decay of the hidden state is given by

$$\begin{aligned} \gamma_j^m &= \exp\{-\max(0, a_\gamma r^m \Delta_j^m + b_\gamma)\} \\ \tilde{h}_j^m &= \gamma_j^m \times h_j^m \end{aligned} \quad (5)$$

where γ is the decay ratio and a, b are learnable decay parameters. For the concat layer, the decay is achieved by substituting j with l^* in eq. 5 where $h_{l^*}^m = h_k^m$ and $k = \max(1, \dots, l)$ such that $m \in \mathbb{S}_k$.

Working of an LSTM Block The gates of L^m at time t_j are denoted by forget gate (f_j^m), input gate (i_j^m) and output gate (o_j^m). Based on the decayed hidden states (\tilde{h}_j^m) given by eq. 5 and summary state (c_{j-1}) given by eq. 2, the working of L^m is given in eq. 6. We indicate the changes in the SLAN LSTM block compared to a vanilla LSTM (Hochreiter and Schmidhuber 1997) by underlining in the eq. 6.

$$\begin{aligned} f_j^m &= \sigma(W_f^m x_j^m + V_f^m \tilde{h}_j^m + b_f^m) \\ i_j^m &= \sigma(W_i^m x_j^m + V_i^m \tilde{h}_j^m + b_i^m) \\ o_j^m &= \sigma(W_o^m x_j^m + V_o^m \tilde{h}_j^m + b_o^m) \\ \tilde{c}_j^m &= \tanh(W_c^m x_j^m + V_c^m \tilde{h}_j^m + b_c^m) \\ c_j^m &= f_j^m c_{j-1} + i_j^m \tilde{c}_j^m \\ h_j^m &= o_j^m \tanh(c_j^m) \end{aligned} \quad (6)$$

Discussion

SLAN vs GRU-D When data is missing, GRU-D performs input imputation and hidden state decay. The input is imputed as $x_t^d = m_t^d x_t^d + (1 - m_t^d) \gamma_{x_t^d} x_{t'}^d + (1 - m_t^d) (1 - \gamma_{x_t^d}) \tilde{x}^d$ where m_t^d is the masking value, γ is the decay factor, $x_{t'}^d$ is the last observation of the d^{th} variable ($t' < t$) and \tilde{x}^d is the empirical mean of the d^{th} variable. The hidden state is decayed as $\gamma_j = \exp\{-\max(0, a_\gamma \Delta_j + b_\gamma)\}$. When the data is not missing, GRU-D just performs the hidden state decay to capture richer knowledge from missingness. Thus, *GRU-D performs imputation only when the data*

is missing. Whereas, in SLAN, when data is missing, the switch of the LSTM corresponding to that data value is 'off'. Hence, *we don't perform any form of imputation*. When data is not missing, SLAN performs the hidden state decay (eq. 5) to capture information from time delay (Δ_j^m). The decay in SLAN also differs from GRU-D as SLAN introduces a sampling rate (r^m) in the decay function (eq. 5) to inculcate the different delay intervals of sensors. Hence, SLAN is different from GRU-D and not an imputation model.

Switch Layer in SLAN vs Observation Mask in Transformer The switch layer dynamically changes the architecture of SLAN (see *supplementary S2*) to adapt to missing values by explicitly informing the model which LSTM blocks will be active. This is different from the observation mask in Transformers. The architecture of the Transformer is fixed, where the observation masks are concatenated with the input value and passed as input to the model. The Transformer implicitly learns the meaning of observational masks via training.

Experiments

Datasets We consider MIMIC-III, Physionet 2012 and Physionet 2019 datasets to showcase the efficacy of SLAN. These datasets are widely used for the ISTS study. The description of the datasets is given in Table 1. We prepare the datasets by following SeFT (Horn et al. 2020). For MIMIC-III and Physionet 2012, the mortality prediction task is considered and for Physionet 2019, early sepsis prediction (ESP) is done. Therefore, in ESP, after each observation time, sepsis is predicted. We refer to MIMIC-III, Physionet 2012 and Physionet 2019 (ESP) as M-3, P-12 and P-19, respectively. The datasets are skewed with 13.22%, 14.24% and 1.79% positive labels for M-3, P-12 and P-19, respectively, making them challenging datasets. The number of missingness is 7.5×10^6 , 28×10^6 and 42.3×10^6 for M-3, P-12 and P-19 dataset leading to high irregularity. A detailed description of the dataset is presented in *supplementary S3*.

Baselines We consider both non-imputation and imputation baselines. Among imputation, GRU-D (Che et al. 2016), IP-Nets (Shukla and Marlin 2018), GRU-Simple (Che et al. 2016) and Phased-LSTM (Neil, Pfeiffer, and Liu 2016) are considered. The non-imputation baselines are Transformer (Vaswani et al. 2017), SeFT (Horn et al. 2020) and Raindrop (Zhang et al. 2021). We do not run most of the baseline models since we follow the exact data preprocessing steps,

Table 2: Comparison of various methods on M-3, P-12 and P-19 datasets. *Imp* and *No-Imp* represent imputation models and non-imputation models, respectively. The **best** and 2nd best performance is represented by **bold** and UNDERLINE, respectively. The metric is reported as the mean \pm standard deviation of three runs with different seeds.

Model		MIMIC-III		Physionet 2012		Physionet 2019 (ESP)			
		AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	B _{Acc}	U _{norm}
Imp	GRU-D	52.0\pm0.8	85.7\pm0.2	53.7 \pm 0.9	86.3\pm0.3	5.3 \pm 0.4	67.4 \pm 1.2	57.4 \pm 0.2	12.6 \pm 1.1
	IP-Nets	48.3 \pm 0.4	83.2 \pm 0.5	51.0 \pm 0.6	86.0 \pm 0.2	5.1 \pm 0.8	74.2 \pm 1.2	63.8 \pm 0.9	-11.9 \pm 4.0
	GRU-SIMPLE	43.6 \pm 0.4	82.8 \pm 0.0	42.2 \pm 0.6	80.8 \pm 1.1	6.1 \pm 0.7	78.1 \pm 1.5	<u>71.0\pm1.4</u>	26.9 \pm 4.1
	Phased-LSTM	37.1 \pm 0.5	80.3 \pm 0.4	38.7 \pm 1.5	79.0 \pm 1.0	5.5 \pm 0.9	75.4 \pm 1.3	67.5 \pm 1.7	20.2 \pm 3.2
No-Imp	TRANSFORMER	42.6 \pm 1.0	82.1 \pm 0.3	52.8 \pm 2.2	86.3\pm0.8	3.6 \pm 0.9	65.8 \pm 3.7	53.6 \pm 1.7	-43.9 \pm 10.0
	SeFT	46.3 \pm 0.5	83.9 \pm 0.4	52.4 \pm 1.1	85.1 \pm 0.4	4.8 \pm 0.2	76.8 \pm 0.9	70.9 \pm 0.8	25.6 \pm 1.9
	RAINDROP	34.8 \pm 1.4	79.3 \pm 0.9	48.8 \pm 3.1	84.3 \pm 1.1	<u>7.6\pm0.2</u>	<u>78.1\pm0.4</u>	69.3 \pm 0.8	48.4\pm0.1
	SLAN (Ours)	45.7 \pm 0.9	<u>84.9\pm0.2</u>	54.9\pm0.4	86.2 \pm 0.2	9.7\pm0.6	80.5\pm2.0	71.8\pm2.9	<u>48.1\pm0.3</u>

splits and metrics implementation of SeFT. We report their performance directly from the SeFT paper in Table 2. We only run RAINDROP ourselves. See *supplementary S4* for its implementation details.

Comparison Metrics The datasets are imbalanced, thus, we use the area under the receiver operating characteristic (AUROC) and the area under the precision-recall curve (AUPRC) as comparison metric. Moreover, for P-19, we also consider the balanced accuracy and utility metric as reported in SeFT. See *supplementary S5* for more details.

Implementation Details We consider the train-val-test split of all datasets provided in SeFT. To handle the imbalance (Table 1), we resort to a weighted-oversampling strategy for M-3 and P-12 and weighted loss for P-19. Weighted oversampling involves preparing the training batch by sampling the data based on the class weights given by the inverse frequency of the class. For weighted loss, the class weights are 0.5094 and 26.9973 for negative and positive labels of P-19. The models are trained for 30 epochs with an early stopping of 5 on AUPRC to avoid overfitting. SLAN uses cross-entropy loss, AdamW optimizer, data standardization, and mean aggregate function. The dropout is 0.3 and the short-term and long-term memory size is 32. The learning rate is 0.0005, 0.0001 and 0.0005 for M-3, P-12 and P-19 datasets, respectively. The learning rate is adaptive with decay by a factor of 0.5 after each epoch without improvement. The batch size is 32, 16 and 32 for M-3, P-12 and P-19 datasets, respectively. Since P-12 and P-19 dataset has 6 and 4 static features, respectively, the embedding of these features is concatenated in the final concat layer before applying a fully connected layer for prediction. The size of the embedding is kept equal to the size of the global summary state. We ran experiments on an Intel Xeon E5-2630 v3 (8 cores/CPU) machine. More details on implementation are given in *supplementary S4*.

Results The results of M-3, P-12 and P-19 datasets are presented in Table 2. For the P-12 dataset, SLAN performs best across all the methods in terms of AUPRC and second best in terms of AUROC. In terms of AUROC, SLAN performs only marginally poor (86.2) in comparison to the best

performance (86.3). In M-3, SLAN gives the second-best result in AUROC and 4th best in AUPRC. SLAN outperforms all the other models by significant margins for P-19 datasets in AUPRC, AUROC and balanced accuracy. Overall, SLAN performs competitively on all the datasets. Moreover, the performance of SLAN improves compared to the best imputation model as the number of imputation values increases from the M-3 (7.5×10^6) to the P-19 dataset (42.3×10^6) (see Fig. 1). Thus, SLAN justifies the need for a better non-imputation model.

Ablation Studies

We perform several ablation studies to investigate the efficacy of SLAN architecture. All the following experiments of SLAN in M-3 and P-12 datasets are performed by following the implementation details given in the previous section, unless otherwise stated.

Different Aggregation Methods We compare the performance of SLAN for mean, max and simple attention (Bahdanau, Cho, and Bengio 2014) as the aggregation function to calculate the global summary state. In attention, the normalized weightage of the LTM of each active LSTM block is determined using a single-layer feed-forward neural network followed by the weighted average of LTMs to output the global summary state. In max, the element-wise max is performed over the candidate summary states. The comparison is reported in the upper half of Table 3. The performance of max is lower than both attention and mean because max may downplay the contribution of many LTMs by highlighting just one, thus becoming sensitive to outliers. Among mean and attention, attention performs the best in P-12, whereas mean gives better results in M-3.

De-mistifying Concat Layer SLAN’s concat layer consists of a global summary state and the local summary state of each sensor. We remove the global summary state from the concat layer to check the informativeness of the local summary state (see *Only L.S.* in Table 3). When compared with the default setting of SLAN, *i.e.* *G.S. + L.S.*, *Only L.S.* is slightly poorer (max by $\sim 0.9\%$). This is due to the fact

Table 3: Comparison of the performance of SLAN for different aggregation functions and variants of concat layer. Att stands for attention. G.S. stands for global summary state and L.S. stands for local summary state. G.S. + L.S. is the default setting of SLAN.

	MIMIC-III		Physionet 2012	
	AUPRC	AUROC	AUPRC	AUROC
	<i>Aggregation Function</i>			
Max	43.5±0.5	84.1±0.6	54.8±0.3	85.9±0.3
Att	44.8±1.2	84.4±0.6	55.3±0.4	86.2±0.4
Mean	45.7±0.9	84.9±0.2	54.9±0.4	86.2±0.2
	<i>Concat</i>			
Only G.S.	40.4±1.3	82.8±0.6	48.1±1.3	83.1±0.4
Only L.S.	45.3±1.2	84.9±0.2	54.5±0.8	86.0±0.2
G.S.+L.S.	45.7±0.9	84.9±0.2	54.9±0.4	86.2±0.2

that the local summary state contains individual sensor information and is also aided by the global summary state at each time step. SLAN with only global summary states in the concat layer (*Only G.S.*) performs 12.4% and 3.6% poorer than *G.S. + L.S.* in terms of AUPRC and AUROC on P-12. *Only G.S.* performs 11.6% and 2.5% poorer than *G.S. + L.S.* in AUPRC and AUROC, respectively for M-3. Indeed *Only G.S.* performs significantly poorer than *G.S. + L.S.*, still it contains sufficient summarised information to outperform baseline models like Phased-LSTM in both datasets, GRU-Simple in P-12 and Latent-ODE in M-3 (see Table 2).

Data Scalability In the practical setting, it is important for any model to have data scalability meaning the performance of the model on test data should improve as the amount of training data increases. We consider the first 25%, 50%, 75% and 100% training data for both P-12 and M-3 and train our model on them. The average and the 95% confidence interval of 3 different runs of SLAN on the test data are shown in Fig. 4. The performance of SLAN steadily increases with the increasing amount of training data on both datasets. The percentage improvement of AUPRC for M-3, when trained on 50% data compared to 25% data is 6.2%, 75% data compared to 50% data is 5.2% and 100% data compared to 75% data is 2.9%. Transitioning from 25% to 50% data, we double the number of instances thus the percentage improvement is highest. Whereas when trained on 100% data compared to 75% data, we add only 1/3rd data thus the percentage improvement is lowest. The same trend is followed in the AUROC of M-3, AUPRC and AUROC of P-12. The exact value of both metrics is given in *supplementary S6*.

Limitations (1) *Time Requirement*: SLAN processes all the activated LSTM at t_j , then calculates the global summary state, which is used for time t_{j+1} . Thus, SLAN works in a *serialized* manner. Moreover, the current implementation of SLAN processes each activated LSTM block one by one at time t_j . This results in an algorithmic complexity of

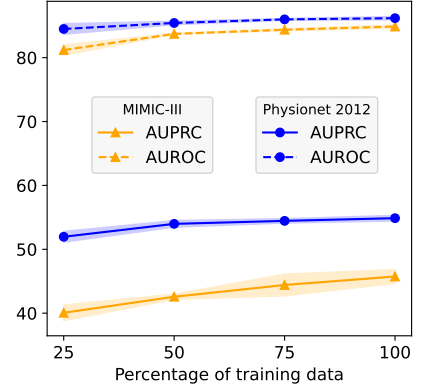


Figure 4: Performance of SLAN on different percentages of training datasets. The average with 95% confidence interval of 3 runs with random seed is reported here.

$O(npq)$, where n is #Instances, p is #Sensors and q is #Observations (see Table 1 for definition). SLAN takes 2211.1 and 3749.4 seconds per epoch compared to 7.62 and 20.1 seconds per epoch of SeFT on P-12 and M-3, respectively. High time requirement is a limitation of our model with respect to other faster models. We plan to eradicate this problem in our future work.

(2) *Scalability to #Sensors*: In the current implementation of SLAN, the time complexity is linearly dependent on the number of sensors. Therefore, SLAN may not be scalable to applications with many sensors. However, it should be noted that the datasets used in this paper are standard for ISTS applications and can be considered a good representation of practical applications. SLAN handles the 17 features of the M-3 dataset as well as the 37 features of the P-12 dataset seamlessly. Thus SLAN is sufficiently scalable to handle practical scenarios.

Conclusion

We propose a Switch LSTM Aggregate Network (SLAN) to handle missing not at random multivariate irregularly-sampled time series data without any imputation. The competitive performance of SLAN and various ablation studies empirically demonstrate the effectiveness of our proposed model on ISTS data. We also establish the sub-optimal performance of the imputation model compared to SLAN as the number of imputations increases. Moreover, the SLAN framework can be extended for modelling multi-modality data, like adding clinical notes to sensor measurement data (Zhang et al. 2023). SLAN can also be leveraged for streaming data modelling in an online setting with time-variant dimensions (Agarwal et al. 2023a,b). It would also be interesting to see the application of SLAN in other fields. We plan to explore the above fields in the future.

References

- Agarwal, R.; Agarwal, K.; Horsch, A.; and Prasad, D. K. 2023a. Auxiliary network: scalable and agile online learning for dynamic system with inconsistently available inputs. In *Neural Information Processing: 29th International Conference, ICONIP 2022, Virtual Event, November 22–26, 2022, Proceedings, Part I*, 549–561. Springer.
- Agarwal, R.; Gupta, D.; Horsch, A.; and Prasad, D. K. 2023b. Aux-Drop: Handling Haphazard Inputs in Online Learning Using Auxiliary Dropouts. *Transactions on Machine Learning Research*.
- Bahdanau, D.; Cho, K.; and Bengio, Y. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Brodersen, K. H.; Ong, C. S.; Stephan, K. E.; and Buhmann, J. M. 2010. The balanced accuracy and its posterior distribution. In *2010 20th international conference on pattern recognition*, 3121–3124. IEEE.
- Chang, Z.; Liu, S.; Qiu, R.; Song, S.; Cai, Z.; and Tu, G. 2023. Time-aware neural ordinary differential equations for incomplete time series modeling. *The Journal of Supercomputing*, 1–29.
- Che, Z.; Purushotham, S.; Cho, K.; Sontag, D. A.; and Liu, Y. 2016. Recurrent Neural Networks for Multivariate Time Series with Missing Values. *Scientific Reports*, 8.
- Goldberger, A. L.; Amaral, L. A.; Glass, L.; Hausdorff, J. M.; Ivanov, P. C.; Mark, R. G.; Mietus, J. E.; Moody, G. B.; Peng, C.-K.; and Stanley, H. E. 2000. PhysioBank, PhysioToolkit, and PhysioNet: components of a new research resource for complex physiologic signals. *circulation*, 101(23): e215–e220.
- Harutyunyan, H.; Khachatrian, H.; Kale, D. C.; and Galstyan, A. G. 2017. Multitask learning and benchmarking with clinical time series data. *Scientific Data*, 6.
- Hochreiter, S.; and Schmidhuber, J. 1997. Long Short-Term Memory. *Neural Computation*, 9: 1735–1780.
- Horn, M.; Moor, M.; Bock, C.; Rieck, B.; and Borgwardt, K. 2020. Set functions for time series. In *International Conference on Machine Learning*, 4353–4363. PMLR.
- Ipsen, N. B.; Mattei, P.-A.; and Frellsen, J. 2020. not-MIWAE: Deep Generative Modelling with Missing not at Random Data. In *International Conference on Learning Representations*.
- Johnson, A. E.; Pollard, T. J.; Shen, L.; Lehman, L.-w. H.; Feng, M.; Ghassemi, M.; Moody, B.; Szolovits, P.; Anthony Celi, L.; and Mark, R. G. 2016. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3(1): 1–9.
- Lim, D. K.; Rashid, N. U.; Oliva, J. B.; and Ibrahim, J. G. 2021. Handling non-ignorably missing features in electronic health records data using importance-weighted autoencoders. *arXiv e-prints*, arXiv–2101.
- Ma, C.; and Zhang, C. 2021. Identifiable generative models for missing not at random data imputation. *Advances in Neural Information Processing Systems*, 34: 27645–27658.
- Marlin, B. M.; Kale, D. C.; Khemani, R. G.; and Wetzel, R. C. 2012. Unsupervised pattern discovery in electronic health care data using probabilistic clustering models. In *Proceedings of the 2nd ACM SIGHIT international health informatics symposium*, 389–398.
- Mudelsee, M. 2002. TAUEST: a computer program for estimating persistence in unevenly spaced weather/climate time series. *Computers & Geosciences*, 28: 69–72.
- Neil, D.; Pfeiffer, M.; and Liu, S.-C. 2016. Phased LSTM: Accelerating Recurrent Network Training for Long or Event-based Sequences. In *NIPS*.
- Rasmussen, C. E.; and Williams, C. K. I. 2003. Gaussian Processes for Machine Learning. In *Adaptive computation and machine learning*.
- Ravuri, S. V.; Lenc, K.; Willson, M.; Kangin, D.; Lam, R. R.; Mirowski, P. W.; Fitzsimons, M.; Athanassiadou, M.; Kashem, S.; Madge, S.; Prudden, R.; Mandhane, A.; Clark, A.; Brock, A.; Simonyan, K.; Hadsell, R.; Robinson, N. H.; Clancy, E.; Arribas, A.; and Mohamed, S. 2021. Skilful precipitation nowcasting using deep generative models of radar. *Nature*, 597: 672 – 677.
- Reyna, M. A.; Josef, C.; Seyedi, S.; Jeter, R.; Shashikumar, S. P.; Westover, M. B.; Sharma, A.; Nemati, S.; and Clifford, G. D. 2019. Early prediction of sepsis from clinical data: the PhysioNet/Computing in Cardiology Challenge 2019. In *2019 Computing in Cardiology (CinC)*, Page–1. IEEE.
- Rubanava, Y.; Chen, T. Q.; and Duvenaud, D. K. 2019. Latent Ordinary Differential Equations for Irregularly-Sampled Time Series. In *Neural Information Processing Systems*.
- Rubin, D. B. 1976. Inference and missing data. *Biometrika*, 63(3): 581–592.
- Shrive, F. M.; Stuart, H.; Quan, H.; and Ghali, W. A. 2006. Dealing with missing data in a multi-question depression scale: a comparison of imputation methods. *BMC Medical Research Methodology*, 6: 57 – 57.
- Shukla, S. N.; and Marlin, B. 2018. Interpolation-Prediction Networks for Irregularly Sampled Time Series. In *International Conference on Learning Representations*.
- Vaswani, A.; Shazeer, N. M.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and Polosukhin, I. 2017. Attention is All you Need. In *NIPS*.
- Wu, Y.; Hernández-Lobato, J. M.; and Ghahramani, Z. 2013. Dynamic Covariance Models for Multivariate Financial Time Series. In *International Conference on Machine Learning*.
- Yadav, P.; Steinbach, M.; Kumar, V.; and Simon, G. 2018. Mining electronic health records (EHRs) A survey. *ACM Computing Surveys (CSUR)*, 50(6): 1–40.
- Zeng, F.; and Gao, W. 2022. Early Rumor Detection Using Neural Hawkes Process with a New Benchmark Dataset. In *North American Chapter of the Association for Computational Linguistics*.
- Zhang, X.; Li, S.; Chen, Z.; Yan, X.; and Petzold, L. R. 2023. Improving medical predictions by irregular multi-modal electronic health records modeling. In *International Conference on Machine Learning*, 41300–41313. PMLR.

Zhang, X.; Zeman, M.; Tsiligkaridis, T.; and Zitnik, M. 2021. Graph-Guided Network for Irregularly Sampled Multivariate Time Series. In *International Conference on Learning Representations*.

Supplementary

S1 - Baseline Models

Imputation-based models

GRU-D (Che et al. 2016) proposed GRU-D which exploits missingness by considering two main missingness representation methods: masking and timestamps, to devise effective solutions to characterize the missing patterns. The proposed model aims to use the masking information and temporal pattern in the missingness via the two trainable decay terms. The decay is calculated as

$$\gamma_t = \exp\{-\max(0, W_\gamma \delta_t + b_\gamma)\} \quad (7)$$

where γ is the decay parameter at time t , W and b are model parameters to learn the decay. GRU-D decays the hidden states as

$$h_{t-1} = \gamma_{h_t} \odot h_{t-1} \quad (8)$$

where h_{t-1} is the hidden state from time $t - 1$ and γ_{h_t} is decay value of hidden state at time t .

GRU-D further imputes the input missing value whenever the input data is missing. The imputation is done by the following equation

$$x_t^d = m_t^d x_t^d + (1 - m_t^d) \gamma_{x_t^d} x_{t'}^d + (1 - m_t^d)(1 - \gamma_{x_t^d}) \tilde{x}^d \quad (9)$$

Here, m_t^d represents the masking value which is 1 if the sensor is measured otherwise 0, $\gamma_{x_t^d}$ is the decay factor, $x_{t'}^d$ is the last observation of the d^{th} variable ($t' < t$) and \tilde{x}^d is the empirical mean of the d^{th} variable. Thus, the missing input feature is imputed whenever not measured.

IP-Nets (Shukla and Marlin 2018) proposed Interpolation-Prediction Networks which consist of an interpolation network followed by a prediction network.

In the interpolation network, IP-Nets convert irregularly sampled time series to regularly sampled time series (RSTS). It uses the information from each time series to interpolate values of all the other time series. IP-Nets considers a set of reference time points $r = [r_1, \dots, r_T]$. All the reference time points are evenly spaced within its interval. For each sensor of an instance, IP-Nets output three interpolants (cross-channels, transient component and intensity) corresponding to each reference point and a sensor. Thus, the interpolation network takes i^{th} ISTS instance (X_i) as input and outputs i^{th} RSTS interpolated output (\hat{X}_i) where the dimension of \hat{X}_i is $(3s) \times T$. Here, s is the number of sensors/features, T is the number of reference time points and 3 represents the number of interpolants corresponding to each time point for each sensor.

Finally, in the prediction network, \hat{X}_i is used as an input to produce the final prediction as $\hat{y}_i = g_\theta(\hat{X}_i)$.

GRU-Simple (Che et al. 2016) considers GRU-Simple architecture as a baseline for GRU-D. Here for each missing value, GRU-Simple concatenates the imputed inputs, missing indicator and the time delay (time duration for which the input was missing) and passed it as an input to the GRU. This is given by

$$x_t^n = [x_t^n; m_t^n; \delta_t^n] \quad (10)$$

where x_t^n is imputed via the global mean of each variable or using forward imputation.

Latent ODE is an imputation-based model (Chang et al. 2023). (Rubanova, Chen, and Duvenaud 2019) introduced Latent-ODE which does not impute data as a preprocessing step but rather models the probability of observation times using Poisson processes to estimate missing observations. It has been shown to perform well in the interpolation and extrapolation tasks.

Phased-LSTM (Neil, Pfeiffer, and Liu 2016) proposed an extension of LSTM with an additional time gate that operates with parameterized rhythmic oscillation and contributes to the updation of memory cells to deal with the irregular pattern of sampling. Phased-LSTM cannot handle unaligned measurements as subject to its design which is for asynchronous, event-triggered updates in continuous time. Thus, in order to deal with the problem of missing values resulting from the ISTS data, imputation is required. Any imputation technique, such as, mean, forward-filling, etc., can be used. SeFT uses imputation similar to GRU-Simple (See Appendix of (Horn et al. 2020)).

Non-Imputation models

Transformer The architecture of the Transformer (Vaswani et al. 2017) is fixed, where the observation masks are concatenated with the input value and passed as input to the model. The transformer learns implicitly via training the meaning of positional encoding (observational masks) and which sensors are not measured. Moreover, the architecture of SLAN is completely different from Transformer and can better capture the temporal dependencies within each time series. This can be further empirically seen by the performance of SLAN compared to the transformer in MIMIC-III, Physionet 2012 and Physionet 2019 datasets, where SLAN outperforms the transformer significantly in most of the metrics.

SeFT (Horn et al. 2020) proposed set functions for time series (SeFT) to handle ISTS. SeFT considers sequential ISTS data as a set of observations, losing the order-variant nature of the data. Each observed input of a sensor is considered a data point and passed to a multi-layer perceptron to give an output. All the outputs are aggregated and a neural network is applied to the aggregate for a final prediction. The order-invariant nature of set representation fails to capture the irregularity information which is order-variant and increases with time. Since ISTS data is sequential, considering it as a set-based representation is detrimental to the performance of SeFT. Moreover, ISTS data exhibits a temporal nature, therefore, the informativeness associated with missingness grows with time. For e.g., the result of a diagnosis of the patient on admission day has an influence on the medical attention they need the next day. Therefore, we advocate that to better understand the missingness information found in ISTS data, the temporal order is important and hence sequential processing is important too. SeFT, although passes the time information as input, still does not capture the temporal order nor exploits the sequential processing.

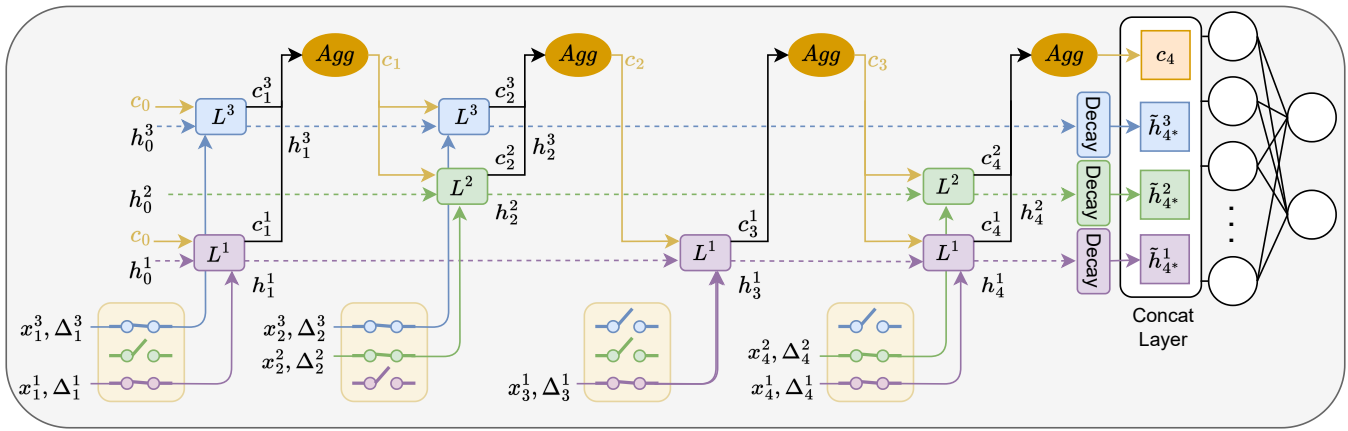


Figure 5: Unrolled SLAN architecture based on the example in Figure 2c of the main manuscript. (Best viewed in color)

Raindrop The authors (Zhang et al. 2021) propose a graph neural network-based method to jointly learn dependency graphs and fixed dimensional embedding for each sensor for downstream tasks. Each node in a graph corresponds to a sensor. At each observed input, Raindrop updates the corresponding nodes embeddings. Moreover, it considers that there is a relation between different nodes. Thus Raindrop uses message aggregation to estimate the observation embedding for each related node. This is achieved using sensor dependency graphs.

S2 - Unrolled SLAN: Showcasing the Dynamic Adaptability of SLAN

An unrolled architecture (see Figure 3e in the main manuscript) based on the snapshot of an instance (see Figure 2c in the main manuscript) is given in the main paper. We present the detailed workflow of the unrolled SLAN architecture here. For ease of reading, we also provide here the same unrolled SLAN architecture in Figure 5 (from left to right). The progression of SLAN at each time step is discussed next.

- At time t_1 , we receive input $Z_1 = (t_1, \{(x_1^1, \Delta_1^1), (x_1^3, \Delta_1^3)\})$. Based on the measured sensors, the switch layer is $\mathbb{S}_1 = \{1, 3\}$, indicating switch 1 and switch 3 are "on". Thus, the associated LSTM blocks L^1 and L^3 are activated. The hidden states (h_0^1, h_0^2, h_0^3) and summary state (c_0) is initialized randomly. The time delay is $\Delta_1^1 = \Delta_1^3 = 0$. Using eq. 3 (see main manuscript), we get $(h_1^1, c_1^1) = L^1(x_1^1, h_0^1, c_0, \Delta_1^1)$ and $(h_1^3, c_1^3) = L^3(x_1^3, h_0^3, c_0, \Delta_1^3)$ where the inner working of L^m is given by eq. 6 (see the main manuscript). The LTM (c_1^1, c_1^3) is aggregated to give the next summary state c_1 .
- At time step t_2 , $Z_2 = (t_2, \{(x_2^2, \Delta_2^2), (x_2^3, \Delta_2^3)\})$, thus $\mathbb{S}_2 = \{2, 3\}$. Corresponding LSTM L^2 and L^3 are kept active. The previous hidden state of L^2 and L^3 are h_0^2 and h_1^3 respectively. The time delay is $\Delta_2^2 = \Delta_2^3 = t_2 - t_1$. We calculate $(h_2^2, c_2^2) = L^2(x_2^2, h_0^2, c_1, \Delta_2^2)$ and $(h_2^3, c_2^3) = L^3(x_2^3, h_1^3, c_1, \Delta_2^3)$ using eq. 3 (see the main

manuscript). Based on this, the summary state c_2 is given by $agg(c_2^2, c_2^3)$.

- At time t_3 , the input is $Z_3 = (t_3, \{(x_3^1, \Delta_3^1)\})$, time delay is $\Delta_3^1 = t_3 - t_1$ and switch layer is $\mathbb{S}_3 = 1$. Thus we compute $(h_3^1, c_3^1) = L^1(x_3^1, h_0^1, c_2, \Delta_3^1)$ and $c_3 = agg(c_3^1)$.
- Next at time t_4 , we get $(h_4^1, c_4^1) = L^1(x_4^1, h_3^1, c_3, \Delta_4^1)$ and $(h_4^2, c_4^2) = L^2(x_4^2, h_2^2, c_3, \Delta_4^2)$ where $\Delta_4^1 = t_4 - t_3$ and $\Delta_4^2 = t_4 - t_2$.

Finally, we decay the hidden states to concatenate. To achieve this, an auxiliary time t_{4^*} is introduced such that $t_{4^*} = t_4$ and $\Delta_{4^*}^1 = 0$, $\Delta_{4^*}^2 = 0$ and $\Delta_{4^*}^3 = t_4 - t_2$. The hidden states are decayed using eq. 5 (see the main manuscript) to get $\tilde{h}_{4^*}^1, \tilde{h}_{4^*}^2$ and $\tilde{h}_{4^*}^3$. The decayed hidden states and the summary state c_4 where $c_4 = agg(c_4^1, c_4^2)$ are concatenated to give $C = \{c_4, \tilde{h}_{4^*}^3, \tilde{h}_{4^*}^2, \tilde{h}_{4^*}^1\}$. A fully connected layer is employed to give the final prediction as $\hat{y} = F(C)$ (see eq. 4 in the main manuscript). This demonstrates the simplicity of SLAN in dynamically adapting the irregularly sampled sensor measurements without any need for imputation.

S3 - Datasets

MIMIC-III MIMIC (Johnson et al. 2016) is a dataset of stays of patients in the critical care unit at a large tertiary care hospital. It has 21142 stays of unique patients (instances) with a median length of stay of 2.1 days. A total of 17 physiological measurements, like vital signs, medications, etc., are recorded for each patient. Following SeFT (Horn et al. 2020), we remove 32 instances. The discarded instances contained dramatically different recording frequencies compared to the rest of the dataset. Thus, the total number of instances is 21110. The statistics of the datasets can be found in Table 1 in the main manuscript. We train our model for the in-hospital mortality prediction tasks.

Physionet 2012 Physionet Mortality Prediction Challenge 2012 (Goldberger et al. 2000) is a dataset of 12000 patient records (instances) containing measurements taken during the first 48 hours of the ICU stays. Each instance is associated with 37 time series variables (sensors) like blood

Table 4: Hyperparameters of all the models. This table is adapted from SeFT (Horn et al. 2020). * All the experiments of Raindrop are done by us.

Model	MIMIC-III	Physionet 2012	Physionet 2019
GRU-D	lr: 0.00016, bs: 32, n_units: 256, dropout: 0.0, recurrent_dropout: 0.2	lr: 0.00138, bs: 512, n_units: 128, dropout: 0.1, recurrent_dropout: 0.1	lr: 0.0069, bs: 128, n_units: 512, dropout: 0.3, recurrent_dropout: 0.3
IP-Nets	lr: 0.00062, bs: 16, n_units: 256, dropout: 0.2, recurrent_dropout: 0.1, imputation_stepsize: 1.0, reconstr_fraction: 0.2	lr: 0.00035, bs: 32, n_units: 32, dropout: 0.4, recurrent_dropout: 0.3, imputation_stepsize: 1.0, reconstr_fraction: 0.75	lr: 0.0008, bs: 16, n_units: 32, dropout: 0.3, recurrent_dropout: 0.4, imputation_stepsize: 1.0, reconstr_fraction: 0.5
GRU-Simple	lr: 0.00011, bs: 32, n_units: 512, dropout: 0.3, recurrent_dropout: 0.4	lr: 0.00022, bs: 256, n_units: 256, dropout: 0.0, recurrent_dropout: 0.0	lr: 0.00024, bs: 64, n_units: 1024, dropout: 0.3, recurrent_dropout: 0.3
Phased-LSTM	lr: 0.00576, bs: 32, n_units: 1024, use_peepholes: False, leak: 0.01, period_init_max: 1000.0	lr: 0.00262, bs: 256, n_units: 128, use_peepholes: True, leak: 0.01, period_init_max: 1000.0	lr: 0.00069, bs: 32, n_units: 512, use_peepholes: False, leak: 0.001, period_init_max: 100.0
Transformer	lr: 0.00204, bs: 256, warmup_steps: 1000, n_dims: 512, n_heads: 8, n_layers: 2, dropout: 0.4, attn_dropout: 0.0, aggregation_fn: mean, max_timescale: 100.0	lr: 0.00567, bs: 256, warmup_steps: 1000, n_dims: 512, n_heads: 2, n_layers: 1, dropout: 0.3, attn_dropout: 0.3, aggregation_fn: max, max_timescale: 1000.0	lr: 0.00027, bs: 128, warmup_steps: 1000, n_dims: 128, n_heads: 2, n_layers: 4, dropout: 0.1, attn_dropout: 0.4, aggregation_fn: mean, max_timescale: 100.0
SeFT	lr: 0.00245, bs: 512, n_phi_layers: 3, phi_width: 64, phi_dropout: 0.1, n_psi_layers: 2, psi_width: 64, psi_latent_width: 128, dot_prod_dim: 128, n_heads: 4, attn_dropout: 0.1, latent_width: 256, n_rho_layers: 2, rho_width: 512, rho_dropout: 0.1, max_timescale: 1000.0, n_positional_dims: 8	lr: 0.00081, bs: 512, n_phi_layers: 4, phi_width: 128, phi_dropout: 0.2, n_psi_layers: 2, psi_width: 64, psi_latent_width: 128, dot_prod_dim: 128, n_heads: 4, attn_dropout: 0.5, latent_width: 32, n_rho_layers: 2, rho_width: 512, rho_dropout: 0.0, max_timescale: 100.0, n_positional_dims: 4	lr: 0.00011, bs: 64, n_phi_layers: 4, phi_width: 32, phi_dropout: 0.0, n_psi_layers: 2, psi_width: 64, psi_latent_width: 128, dot_prod_dim: 128, n_heads: 4, attn_dropout: 0.1, latent_width: 512, n_rho_layers: 3, rho_width: 128, rho_dropout: 0.0, max_timescale: 10.0, n_positional_dims: 16
Raindrop*	lr: 0.001, bs: 128, dropout: 0.3, obs_embedding_size: 4, n_layers: 2, regularization_scale: 0.02, lr_decay: 0.1, time_representation_size: 16, linear_mapping_size: 20, auxiliary_attributes_dimension: 17, edge_pruning: 50%	lr: 0.0005, bs: 128, dropout: 0.2, obs_embedding_size: 4, n_layers: 2, regularization_scale: 0.02, lr_decay: 0.1, time_representation_size: 16, linear_mapping_size: 20, auxiliary_attributes_dimension: 37, edge_pruning: 50%	lr: , bs: 128, dropout: , obs_embedding_size: 4, n_layers: 2, regularization_scale: 0.02, lr_decay: 0.1, time_representation_size: 16, linear_mapping_size: 20, auxiliary_attributes_dimension: 34, edge_pruning: 50%
SLAN	lr: 0.0005, bs: 32, dropout: 0.3, early_stopping: 5, short_term_memory_size: 32, long_term_memory_size: 32, lr_decay_rate: 0.5, aggregation_function: mean, data_standardization: True	lr: 0.0001, bs: 16, dropout: 0.3, early_stopping: 5, short_term_memory_size: 32, long_term_memory_size: 32, lr_decay_rate: 0.5, static_feat_embeddings_size: 32, aggregation_function: mean, data_standardization: True	lr: 0.0005, bs: 32, dropout: 0.3, early_stopping: 5, short_term_memory_size: 32, long_term_memory_size: 32, lr_decay_rate: 0.5, static_feat_embeddings_size: 32, aggregation_function: mean, data_standardization: True, class_weights: (0.5094, 26.9973)

Table 5: Performance of state-of-the-art RAINDROP (Zhang et al. 2021) on Physionet 2012 and MIMIC-III datasets. DP and LR represent *dropout* and *learning rate* hyperparameter setting respectively. All the metrics are reported as the mean \pm standard deviation of 3 runs with different seeds.

DP	LR	MIMIC-III		Physionet 2012	
		AUPRC	AUROC	AUPRC	AUROC
0.2	5e-2	14.0 \pm 3.1	52.5 \pm 7.6	18.7 \pm 3.1	60.2 \pm 5.6
	1e-2	29.9 \pm 3.5	74.4 \pm 3.4	23.1 \pm 0.1	67.6 \pm 0.2
	5e-3	33.8 \pm 0.5	78.2 \pm 0.3	30.5 \pm 9.3	72.3 \pm 6.1
	1e-3	34.6 \pm 1.4	79.2 \pm 0.9	47.9 \pm 1.3	84.0 \pm 0.6
	5e-4	33.6 \pm 0.8	78.9 \pm 0.4	48.8\pm3.1	84.3\pm1.1
	1e-4	34.7 \pm 0.7	79.0 \pm 0.5	47.3 \pm 2.3	83.3 \pm 0.9
0.3	5e-2	11.8 \pm 1.6	49.2 \pm 5.1	16.8 \pm 3.7	55.4 \pm 7.7
	1e-2	30.4 \pm 2.3	76.0 \pm 2.3	22.7 \pm 0.2	67.4 \pm 0.5
	5e-3	35.0\pm1.9	79.2 \pm 0.4	28.4 \pm 7.3	71.3 \pm 4.7
	1e-3	34.8 \pm 1.4	79.3\pm 0.9	46.9 \pm 2.2	83.3 \pm 1.5
	5e-4	33.7 \pm 1.0	79.0 \pm 0.5	48.3 \pm 1.6	84.0 \pm 0.5
	1e-4	34.1 \pm 0.6	79.0 \pm 0.5	48.1 \pm 2.3	83.6 \pm 0.6
SLAN(Ours)		45.7\pm0.9	84.9\pm0.2	54.9\pm0.4	86.2\pm0.2

pressure, lactate, respiration rate, etc and six static descriptor features (i.e., RecordID, Age, Gender, height, ICUType and Weight). We follow the SeFT (Horn et al. 2020) paper and remove 12 instances that do not contain any time series information at all. The weight feature is considered a time series since it is measured multiple times in the observation period. The final dataset has 11988 instances with 37 features. We train our model on the in-hospital mortality task, which is a binary classification task to predict if the patient dies before being discharged by using the data of the first 48 hours of the ICU admission. The statistics of the datasets can be found in Table 1 in the main manuscript.

Physionet 2019 Physionet 2019 Sepsis Early Prediction Challenge (Reyna et al. 2019) launched a challenge for early prediction of sepsis from clinical data. Physionet 2019 has 40333 instances from three different hospitals in the US. It has 34 time series variables (like Heart rate, Platelets, Cholesterol, Temperature *etc*) and four static variables (age, gender, height and ICU type). The data contains a binary label corresponding to each observed time for each instance (or patient), indicating the onset of sepsis within the next 6h to 12h. We consider the train-val-test split provided in SeFT (Horn et al. 2020). Since, at each observation time, sepsis is predicted, the number of labels for Physionet 2019 (ESP) data is 1552093. Considering the sepsis prediction task, the dataset is highly imbalanced, with 1.852% positive labels.

S4 - Implementation Details

We report the hyperparameters of all the models in Table 4. This table is adapted from the SeFT (Horn et al. 2020) paper.

GRU-D, IP-Nets, GRU-Simple, Phased-LSTM, Transformer, SeFT Since we follow the exact data preprocessing steps and splits of SeFT, we do not run these baseline models ourselves. We directly take the performance reported in the SeFT paper and report it in Table 2 of the main

manuscript. The hyperparameters of all these models can be found in Table A.5. of the SeFT paper. For ease of reading, we provide the same hyperparameters reported in SeFT in Table 4 of this supplementary.

Raindrop We use the Raindrop code provided in the official RAINDROP GitHub repository¹ to conduct all experiments on Raindrop model (Zhang et al. 2021). All the default settings of Raindrop are considered for training: observation embedding size (R_u) is set as 4, the number of heads in the transformer encoder is chosen as 2, λ_r is set as 0.02 to adjust the regularization scale, the batch size is fixed at 128, the binary cross-entropy loss function is used and the scheduler with a learning decay factor of 0.1 and patience equals to 1 is employed. Furthermore, we follow the data imbalance handling strategy provided by RAINDROP where sampling is done from the collection of one times the majority class and three times the minority class samples (Zhang et al. 2021). To determine the best performance of Raindrop, We perform the hyperparameter search for learning rate and dropout. Two different dropouts for all the datasets are considered: 0.2 and 0.3. Six different learning rate for Physionet 2012 and MIMIC-III datasets is considered: 0.05, 0.01, 0.005, 0.001, 0.0005 and 0.0001. Physionet 2019 (ESP) requires high amount of GPU and run-time. So, we ran Raindrop on P-19 with the default setting of dropout = 0.2 and learning rate = 0.0005. All the models are run 3 times with different seeds for 30 epochs. The performance of Raindrop on these different hyperparameters is shown in Table 5.

SLAN A detail of different hyperparameters used for the final performance is reported in Table 4. To determine the best performance of SLAN, we varied the values of different hyperparameters. We opted for finding the best hyperparameter one at a time. First, we fixed the learning rate as

¹<https://github.com/mims-harvard/Raindrop/>

Table 6: Performance of SLAN when trained on different percentages of training data. The AUPRC and AUROC are reported for Physionet 2012 and MIMIC-III datasets. The metric is reported as the mean \pm standard deviation of three runs with different seeds.

Model	Physionet 2012		MIMIC-III	
	AUPRC	AUROC	AUPRC	AUROC
25%	52.0 \pm 0.7	84.5 \pm 0.7	40.1 \pm 1.1	81.2 \pm 0.7
50%	54.0 \pm 0.5	85.4 \pm 0.2	42.6 \pm 0.4	83.7 \pm 0.1
75%	54.4 \pm 0.3	86.0 \pm 0.1	44.4 \pm 1.5	84.4 \pm 0.1
100%	54.9 \pm 0.4	86.2 \pm 0.2	45.7 \pm 0.9	84.9 \pm 0.2

0.0005 and varied the batch size to 16, 32, 64, 125, 256 and 512 to find the best batch size. Then, we varied the size of hidden units (short-term memory and long-term memory) to 1, 2, 4, 8, 16, 32 and 64. Finally, based on the best batch size and hidden unit size, we varied the learning rate to .0001, .0005, .001 and .005. All the other hyperparameter values were fixed heuristically. Indeed, the above procedure of finding the best hyperparameter is not optimum. But, since the model takes a large training time, we opted for the above process.

S5 - Comparison Metrics

AUROC AUROC stands for the area under the receiver operating characteristics and it informs the discriminative ability of the model between positive labels and negative labels. For binary classification, based on different thresholds, different true positive rates (TPR) and false positive rate (FPR) is achieved. This gives an ROC curve and the area under this curve is AUROC.

AUPRC The area under the precision-recall curve (AUPRC) is similar to AUROC but instead of the TPR as the y-axis, precision is used and instead of FPR as the x-axis, recall is used. It is mainly used for imbalanced data where the focus is on correctly classifying positive labels.

Balanced Accuracy The balanced accuracy (Brodersen et al. 2010) metric deals with imbalanced datasets and is defined as the average of recall obtained on each class.

Utility Metrics The models are also compared on utility metric introduced by (Reyna et al. 2019). For each instance (or patient) i , based on the prediction at time t_j , a utility function $U(i, t_j)$ is used to calculate the reward for the model. All the reward is summed over all the instances (n) across all the time points (l_i), i.e., $U_{total} = \sum_{i=1}^n \sum_{j=1}^{l_i} U(i, t_j)$. The score is then normalized as $U_{norm} = \frac{U_{total}}{\sum_{i=1}^n l_i}$.

S6 - Data Scalability

In the practical setting, it is important for any model to have data scalability meaning the performance of the model on test data should improve as the amount of training data increases. We consider the first 25%, 50%, 75% and 100% training data for both Physionet 2012 and MIMIC-III and train our model on them. The exact values of AUPRC and

AUROC are shown in Table 6 and the results are discussed in the main manuscript (see paragraph *Data Scalability* in the section *Ablation Studies*). The average \pm standard deviation of 3 runs is reported here.