# Learning Deep Representations from Clinical Data for Chronic Kidney Disease

Duc Thanh Anh Luong
Department of Computer Science and Engineering
University at Buffalo
Buffalo, New York 14260
Email: ducthanh@buffalo.edu

Varun Chandola
Department of Computer Science and Engineering
University at Buffalo
Buffalo, New York 14260
Email: chandola@buffalo.edu

*Abstract*—We study the behavior of a Time-Aware Long Short-Term Memory Autoencoder, a state-of-the-art method, in the context of learning latent representations from irregularly sampled patient data. We identify a key issue in the way such recurrent neural network models are being currently used and show that the solution of the issue leads to significant improvements in the learnt representations on both synthetic and real datasets. A detailed analysis of the improved methodology for representing patients suffering from Chronic Kidney Disease (CKD) using clinical data is provided. Experimental results show that the proposed T-LSTM model is able to capture the long-term trends in the data, while effectively handling the noise in the signal. Finally, we show that by using the latent representations of the CKD patients obtained from the T-LSTM autoencoder, one can identify unusual patient profiles from the target population.

## I. INTRODUCTION

Chronic Kidney Disease (CKD) is a world-wide public health problem, which was responsible for approximately 956,200 deaths globally in 2013 [1]. In CKD, the primary indicator of disease severity is *estimated glomerular filtration rate* (eGFR) - a clinical marker that measures kidney function. Although CKD is a chronic condition in which kidney function deteriorates over time, the courses of disease progressions among patients are highly heterogeneous. From a clinical perspective, being able to identify sub-groups of patients with similar CKD disease progressions is a first step toward disease subtyping as discovered patient sub-groups may exhibit common biological or physiological patterns that allow us to have better understanding of the underlying mechanism of disease within sub-groups and further develop better treatments for those patients.

From a machine learning perspective, the task of identifying groups of patients with similar disease progressions is an unsupervised machine learning problem in which patient progressions exhibited by their eGFR values over time are clustered. Most recent solutions, in the context of CKD, have relied on fitting a set of functions, e.g., splines, polynomials, etc., to represent the disease trajectories [2–4]. However, such models require pre-specifying the form and number of such functions, which becomes a limiting factor. At the same time, recurrent neural network architectures,

such as *Long Short-Term Memory networks* or LSTMs, have become the workhorses for modeling sequence data. Recently, LSTMs have been applied to clinical data for various analytical tasks [5,6]. To handle the issue of irregular sampling, which is typical in clinical settings, a recent modification called Time-Aware LSTM (T-LSTM) has been proposed [5]. T-LSTMs have been shown to capture short and long-range temporal dependencies in the data and learn latent representations from irregularly sampled observations for patients suffering from Parkinson disease.

While the LSTM-based latent representation learning methodology has been studied in limited disease settings, e.g., Parkinson's disease [5] and Asthma [6], it is unclear how the same architecture would perform in a different disease context, i.e., CKD. In fact, LSTMs offer more than one way of inferring the latent representation for a given input profile. Current applications have used one of those — the activations at the hidden unit (See Section III for more details). Will the same methodology be applicable for CKD? At a more fundamental level, the relationship between the activations within different parts of the network and the properties of the input temporal sequence, remains unstudied.

In this paper we explore the above-mentioned relationship to devise a better strategy for learning latent representation of patient disease progression profiles. Through results on synthetically generated temporal sequences, we discover a better representation and then use the improved representation to understand CKD progression.

*Key contributions*

1) We study the behavior of T-LSTM autoencoder using synthetically generated datasets and show that a latent representation that uses activations at hidden and memory units is better than using the hidden unit activation alone, which is the prevalent strategy. 2) An in-depth analysis of the T-LSTM autoencoder, in the context of CKD, using two real-world datasets, is provided, which includes a data-driven strategy to choose key network parameters. 3) An

application of the improved solution to identify anomalous patient profiles is demonstrated.

The remaining of our paper is structured as follows. In Section II, we provide a brief background on T-LSTM Autoencoder, the primary method that we use in this paper. In Section III, we present our approach to represent the longitudinal profiles in the embedded space by using three different synthetic datasets. In Section IV, we present experimental results with T-LSTM Auto-encoder on CKD datasets. In Section V, we discuss the advantages and drawbacks of this method when applying on CKD as well as other related works. Finally, in Section VI, we give a conclusion of our paper.

## II. BACKGROUND

This section provides necessary background to understand the Time-Aware LSTM autoencoder and its suitability for handling longitudinal clinical data.

### A. Recurrent Neural Networks (RNN)

A typical input sequence for a RNN is of the form $\{\mathbf{x}_t\}_{t=1}^T$ where $T$ is the length of the input sequence and $\mathbf{x}_t$ is a vector of size $N$. In vanilla RNN, at each time step $t$, there is a hidden unit $\mathbf{h}_t \in \mathbb{R}^H$ computed by combining both the input signal at this time step $\mathbf{x}_t$ and the hidden unit of previous time step $\mathbf{h}_{t-1}$:

$$\mathbf{h}_t = \tanh\left(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + b\right)$$

With this network construction, the hidden unit $\mathbf{h}_t$ plays a role of memory that captures the information of all previous inputs $\{\mathbf{x}_i\}_{i=1}^t$.

### B. Long Short-Term Memory (LSTM) Networks

For RNN, learning the network parameters is typically done by applying traditional learning algorithm such as gradient descent with back-propagation through time. However, for applications with very long time lags (many steps between the signal and the output), the learning algorithm may suffer from the problem of exploding or vanishing gradients. This problem is addressed in LSTM architecture [7] by enforcing a "constant error carousel" within special units, called *memory units*. The information stored in these units is controlled by a gated structure.

Typically, at a time step $t$, beside the hidden unit $\mathbf{h}_t$, LSTM has dedicated memory cells $\mathbf{C}_t$. The information in memory cells $\mathbf{C}_t$ is controlled by "forget gate" $\mathbf{f}_t$ and "input gate" $\mathbf{i}_t$. The information from previous memory cells can be moved forward to be stored again in the current memory cells, depending on the values of previous hidden units and current inputs.

### C. Time-Aware LSTM Networks

A key assumption with the LSTM model (and also vanilla RNN) is that temporal sequence is regularly sampled. If the sequence has few missing values, one can extend the dimension of input $\mathbf{x}_t$ to encode whether there is missing values in the input. However, for clinical datasets such as patient's lab test results, the time gaps between consecutive observations can be multiple days or weeks or even months. This makes the approach of encoding missing values within the input $\mathbf{x}_t$ infeasible as there are too many missing values in comparison with actual observations. In addition, the time gaps between observations do carry the signal in themselves. For example, in clinical applications, having fewer lab tests within a certain time span may indicate that a patient has good health and therefore does not require to take lab tests regularly. In contrast, having lab tests more often within a certain period may indicate that a patient has some clinical problems and doctor needs to monitor her physiological markers more frequently. With the lack of proper representation of time lapse between consecutive observations in traditional RNN models, Time-Aware LSTM (T-LSTM) is proposed by Baytas et al. [5] to address this issue. Figure 1 shows the structure of a T-LSTM unit.

In T-LSTM, the previous memory cells $\mathbf{C}_t$ is decomposed into short-term memory $\mathbf{C}_t^S$ and long-term memory $\mathbf{C}_t^T$. Only the short-term memory interacts with the time gap $\Delta_t$ to adjust for the time lapse since previous observation. In particular, the non-increasing function $g()$ is first applied on $\Delta_t$ to penalize for the time passed since previous observation. The quantity $g(\Delta_t)$ now becomes a multiplier to adjust for short-term memory $\mathbf{C}_t^S$. The longer time passed since previous observation, the less amount of short-term memory is kept to use in the next calculation.

Beside the decomposition of previous memory cells and the interaction between short-term memory and time lapse $\Delta_t$, the structure of T-LSTM is exactly the same with LSTM. The detailed formulation of T-LSTM is shown as follows:

$$
\begin{aligned}
\mathbf{C}_{t-1}^S &= \tanh\left(\mathbf{W}_d\mathbf{C}_{t-1} + \mathbf{b}_d\right) && \text{(short-term memory)} \\
\hat{\mathbf{C}}_{t-1}^S &= \mathbf{C}_{t-1}^S * g(\Delta_t) && \text{(discounted short-term memory)} \\
\mathbf{C}_{t-1}^T &= \mathbf{C}_{t-1} - \mathbf{C}_{t-1}^S && \text{(long-term memory)} \\
\mathbf{C}_{t-1}^* &= \mathbf{C}_{t-1}^T + \hat{\mathbf{C}}_{t-1}^S && \text{(adjusted previous memory)} \\
\mathbf{f}_t &= \sigma\left(\mathbf{W}_f\mathbf{x}_t + \mathbf{U}_f\mathbf{h}_{t-1} + \mathbf{b}_f\right) && \text{(forget gate)} \\
\mathbf{i}_t &= \sigma\left(\mathbf{W}_i\mathbf{x}_t + \mathbf{U}_i\mathbf{h}_{t-1} + \mathbf{b}_i\right) && \text{(input gate)} \\
\tilde{\mathbf{C}}_t &= \tanh\left(\mathbf{W}_c\mathbf{x}_t + \mathbf{U}_c\mathbf{h}_{t-1} + \mathbf{b}_c\right) && \text{(candidate memory)} \\
\mathbf{C}_t &= \mathbf{f}_t * \mathbf{C}_{t-1}^* + \mathbf{i}_t * \tilde{\mathbf{C}}_t && \text{(current memory)} \\
\mathbf{o}_t &= \sigma\left(\mathbf{W}_o\mathbf{x}_t + \mathbf{U}_o\mathbf{h}_{t-1} + \mathbf{b}_o\right) && \text{(output gate)} \\
\mathbf{h}_t &= \mathbf{o}_t * \tanh(\mathbf{C}_t) && \text{(current hidden state)}
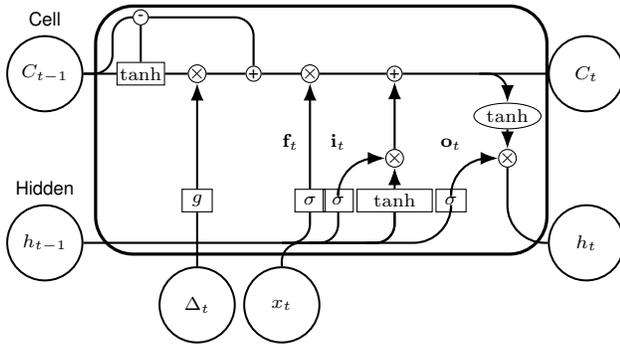\end{aligned}
$$

Fig. 1: Illustration of a T-LSTM unit

## D. *T-LSTM Autoencoder*

With T-LSTM unit as a building block, it is possible to map complex patient longitudinal profile into an embedded space by using an autoencoder structure and analyze disease stratification in this embedded space. In Baytas et al.'s study [5], T-LSTM Autoencoder is represented by two components including one encoder and one decoder. The illustration of T-LSTM autoencoder is shown in Figure 2. In this figure, the structure of encoder and decoder are exactly the same except that the T-LSTM unit in the decoder has another set of learnable parameters $\mathbf{W}_{out}^{<d>}$ and $\mathbf{b}_{out}^{<d>}$ to produce outputs from the hidden unit:

$$\hat{\mathbf{x}}_t = \mathbf{W}_{out}^{<d>}\mathbf{h}_t^{<d>} + \mathbf{b}_{out}^{<d>}$$

In formulation of T-LSTM Autoencoder, we use the superscript $<e>$ and $<d>$ to distinguish the parameters of encoder and decoder respectively. In T-LSTM Autoencoder, the hidden unit and memory unit output from encoder are subsequently passed into the decoder so that it can reconstruct the original longitudinal observations in reverse chronological order. In particular, each T-LSTM unit in the decoder receives two inputs - (1) time gap from previous observation $\Delta_t$ and (2) the current observation $\mathbf{x}_t$; and outputs the prediction of previous observation $\hat{\mathbf{x}}_{t-1}$. For the first time step in the decoder, we define the current observation is a zero vector and time gap from previous observation is zero. The difference between the output prediction $\hat{\mathbf{x}}_t$ and actual observation $\mathbf{x}_t$ is the training signal that gradient-based learning algorithm can use to back-propagate the error through the decoder and encoder and learn to update the model's parameters.

## III. REPRESENTATION OF LONGITUDINAL PROFILES IN T-LSTM AUTO-ENCODER

As illustrated in Figure 2, the hidden unit $\mathbf{h}_T^{<e>}$ and memory unit $\mathbf{C}_T^{<e>}$ at the last time step of the encoder will be fed into the decoder as the initial hidden unit and

memory unit, i.e. $\mathbf{h}_0^{<d>}$ and $\mathbf{C}_0^{<d>}$. Although both hidden unit $\mathbf{h}_T^{<e>}$ and memory unit $\mathbf{C}_T^{<e>}$ are passed from encoder to decoder, in the original paper of Baytas et al. [5], only the hidden unit $\mathbf{h}_T^{<e>}$ is used as the representation of the input longitudinal profile. This raises a question whether the hidden unit alone is sufficient to represent the longitudinal profile and whether the memory unit should be part of the representation. In this section, we investigate different possible representations of a longitudinal profile using T-LSTM Autoencoder and determine the most probable representation.

To simplify our analysis, we only consider longitudinal profile in which each observation has one real value. This is also similar to our application of analyzing CKD that disease progression is monitored via one primary clinical marker - eGFR that measures kidney function. In order to judge the quality of the embedded representations of longitudinal profiles, we will use three synthetic datasets that we know in prior the underlying mechanism of generating data to determine the role of hidden and memory unit in representing longitudinal profile. Each synthetic dataset consists of four clusters and each cluster contains $50$ longitudinal profiles. Each profile belonging to a cluster is sampled by first sampling the $t$ value that represents the time stamp an observation is made and the $x$ value is then sampled by following the underlying line that controls the behavior of the cluster as shown in Table I. The first sample of a profile is always at time $0$ and time gap for subsequent sample is computed by sampling from Poisson distribution with mean value 50, i.e. $\Delta_t \sim Poisson(\lambda = 50)$. This means that the time gap between consecutive observations is approximately $50$ time steps. When time step exceeds $1095$, we stop sampling process for the profile. Figure 3 shows the visualization of our synthetic datasets. As can be observed in both the Figure 3 and the Table I, our synthetic datasets are prepared in a way that each cluster has only one difference from the other within the same dataset. For synthetic dataset 1, each line in the cluster has the same x-intercept and same level of variation and the only difference is the slope. For synthetic dataset 2, all lines within each cluster have the same slope and same level of variations while different clusters have different x-intercepts. For synthetic dataset 3, all clusters are basically horizontal lines with different levels of noises.

Because in all three synthetic datasets, there are only one distinctive difference between clusters (slope for synthetic dataset 1, intercept for synthetic dataset 2 and noise level for synthetic dataset 3), the dimension of hidden unit (and also memory unit) is set to 2 as it is enough to capture the underlying differences between clusters, given that T-LSTM Autoencoder is able to capture this sole difference in the model. The two-dimensional hidden space in hidden unit (and also in memory unit) allows us to visualize
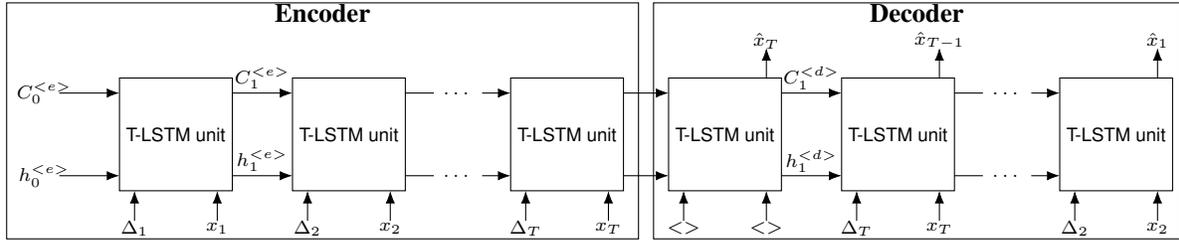
Fig. 2: T-LSTM Autoencoder

| | Synthetic dataset 1 | Synthetic dataset 2 | Synthetic dataset 3 |
|---|---|---|---|
| Cluster 1 | $x = (4/219)t + 60 + \epsilon$ | $x = (4/219)t + 20 + \epsilon$ | $x = 60 + \epsilon_1$ where $\epsilon_1 \sim \mathcal{N}(0, 10)$ |
| Cluster 2 | $x = 60 + \epsilon$ | $x = (4/219)t + 40 + \epsilon$ | $x = 60 + \epsilon_2$ where $\epsilon_2 \sim \mathcal{N}(0, 20)$ |
| Cluster 3 | $x = (-4/219)t + 60 + \epsilon$ | $x = (4/219)t + 60 + \epsilon$ | $x = 60 + \epsilon_2$ where $\epsilon_3 \sim \mathcal{N}(0, 30)$ |
| Cluster 4 | $x = (-8/219)t + 60 + \epsilon$ | $x = (4/219)t + 80 + \epsilon$ | $x = 60 + \epsilon_2$ where $\epsilon_4 \sim \mathcal{N}(0, 40)$ |
| | | where $\epsilon \sim \mathcal{N}(0, 14)$ | |

TABLE I: Summary of three synthetic datasets



(a) Synthetic dataset 1

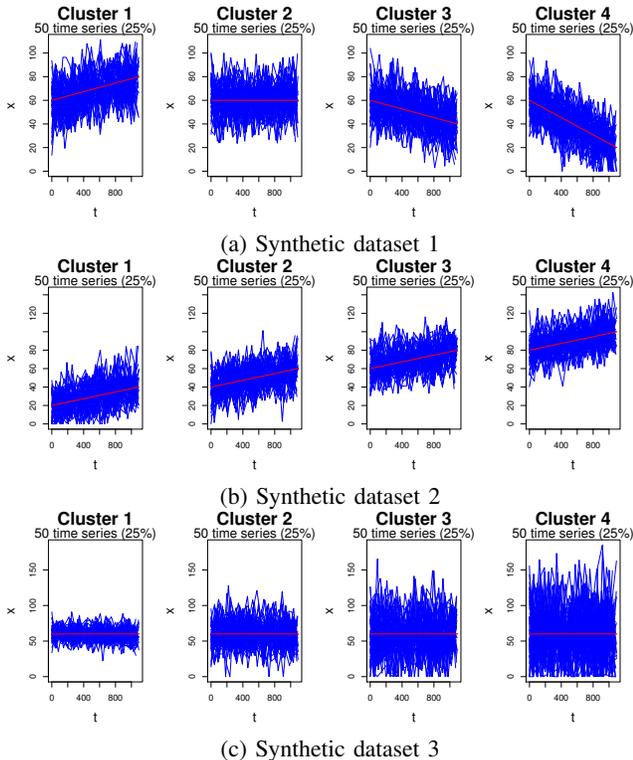(b) Synthetic dataset 2

(c) Synthetic dataset 3

Fig. 3: Visualization of three synthetic datasets

the distribution of longitudinal profiles and inspect the alignment of the representations with their ground-truth clusters. Figure 4 shows visualization of the representations of individual longitudinal profiles in three datasets after being trained for $20,000$ epochs. This visualization provides some interesting insights regarding the role of hidden units and memory units at the last time step in the encoder as

longitudinal profile representations.

As we can see in figure 4a and figure 4d, the distribution of profile representations in synthetic dataset 1 does align with the ground-truth clusters, and this alignment is more visible in memory units than in hidden units. To quantify the alignment of representations in memory units / hidden units with the ground-truth clusters, we compute the silhouette coefficients [8] with ground-truth cluster assignments and distance metric between two longitudinal profiles is the Euclidean distance between their embedded representations. In Silhouette coefficient, the quality of individual data point in a clustering result with respect to a cluster assignment is quantified by the cluster tightness and degree of separation between neighboring clusters. The value of silhouette coefficient ranges from -1 to 1 in which higher value indicates better cluster assignment. The overall quality of a clustering result can be measured as the average of individual Silhouette coefficients. Table II shows the average silhouette coefficient when 3 types of longitudinal representations are used: (1) only hidden unit, (2) only memory unit and (3) both hidden and memory unit. As shown in Table II, for synthetic dataset 1, memory unit is significantly better in distinguishing four different clusters in comparison with hidden unit. When using both hidden unit and memory unit, the clustering quality is still better than hidden unit while only suffers minor reduction in clustering quality in comparison with memory unit.

| Representation | Synthetic dataset 1 | Synthetic dataset 2 |
|---|---|---|
| Hidden unit | 0.0548 | 0.6169 |
| Memory unit | 0.4906 | 0.4769 |
| Hidden and memory unit | 0.4858 | 0.4812 |

TABLE II: Average Silhouette coefficient when using ground-truth cluster membership

For synthetic dataset 2, we can observe in Figure 4b and 4e that there are clear clustering patterns in the visualization of longitudinal representation by memory unit and hidden unit. It is worth to notice that for representation by hidden unit as shown in Figure 4b, the difference between clusters only concentrate in dimension 2 while values in dimension 1 are almost the same for each individual longitudinal profile. On the other hand, when looking at representation by memory unit as shown in Figure 4e, there are clear four clusters spanning in both two dimensions of memory unit. When judging quantitatively the distinguishing power of hidden unit and memory unit, we also use the average silhouette coefficient as similar to synthetic dataset 1. As shown in Table II, for synthetic dataset 2, all three types of representations are well-aligned with the cluster ground-truth as their average silhouette coefficient is very high (the lowest value is $0.4769$ for the case of using memory unit as representation). These average silhouette coefficients show that both hidden and memory unit are well-suited for representing the longitudinal profile, as both of them can capture well the differences in the x-intercept between clusters.

For synthetic dataset 3, we can see in Figure 4c that representations by using hidden unit do not provide any clear clustering patterns. On the other hand, in Figure 4f, most of data points of cluster 1 (cluster with lowest level of noise) concentrates in the center of the space while data points of cluster 4 (cluster with highest level of noise) scatter along the periphery of the space. For cluster 2 and 3, they are blended to some degree between center and the periphery of the space. This is an interesting observation as the level noise tends to play a role in determining the position of the longitudinal profile in latent space of memory unit.

From the examination of the representations of longitudinal profiles in three different synthetic datasets, we can see that memory unit does play an important role in distinguishing different clusters. In a situation as in synthetic dataset 2, using hidden unit as representation does provide better distinguishing power between clusters. Therefore, in subsequent experiments with real clinical datasets, we use both hidden unit and memory unit as the representation for longitudinal profile.

## IV. EXPERIMENTS WITH CHRONIC KIDNEY DISEASE

In this section, we apply T-LSTM Autoencoder to explore the distribution of CKD longitudinal profiles in the latent space and investigate how the embedded representations provide better insights about the progression of CKD. In section IV-A, we briefly describe two CKD datasets that are used in the experiments. In section IV-B, we determine an appropriate dimension of hidden unit (and also memory unit) for T-LSTM Autoencoder to represent the embedded patient profile. In section IV-C, we investigate the reconstruction of original longitudinal profiles output from the decoder of T-LSTM Autoencoder and see how the inputs and outputs are related to each other. In section IV-D, we assess the embedded representations of longitudinal profiles learnt from T-LSTM Autoencoder and identify interesting and unusual longitudinal CKD progressions enabled by these latent representations.

### A. CKD datasets

In our experiments, we use two real-world CKD datasets extracted from two larger clinical datasets: (1) DARTNet dataset [9] and (2) MIMIC-III dataset [10]. From the original datasets, we only consider patients having stage 3 CKD with more than ten eGFR observations and clinical records spanning more than one year to include in our experiments. This preprocessing step is similar to earlier study of Luong and Chandola [3]. After preprocessing, we have $7,142$ patients and $3,082$ patients in CKD cohorts of DARTNet and MIMIC-III dataset respectively. Figure 5 shows the preprocessing steps that we performed to obtain CKD cohorts, which is similar to earlier study of Luong and Chandola [3].

### B. Dimension of hidden units

One of the key hyper-parameter of T-LSTM Autoencoder is the dimension of hidden unit (and also memory unit). This determines the representations of longitudinal profiles after being learned from T-LSTM Autoencoder. In earlier experiments with synthetic datasets, we have decided the dimension of hidden/memory unit to be two as it is enough to capture the variability in our synthetic datasets. However, in the real-world clinical datasets, there are many factors that can affect the progression of a longitudinal clinical profile. Therefore, choosing an appropriate dimension for our embedded representations is not a trivial task. In earlier study of Baytas et al. [5], authors chose hidden dimension to be two for the case of synthetic dataset while they did not provide a clear criterion for choosing the hidden dimension when applying the T-LSTM Autoencoder on the real clinical dataset. In this section, we provide a rigorous process to determine dimension of hidden unit for clinical datasets.

As mentioned earlier, T-LSTM Autoencoder optimizes the reconstruction error. In other words, the reconstructed time series output from the decoder should be as close as possible to the original time series input to the encoder. Therefore, one criterion to estimate the quality of T-LSTM Autoencoder is via the reconstruction error. As T-LSTM Autoencoder optimizes this quantity directly using mini-batch Adam optimizer, the overall reconstruction error should decrease over many epochs during the training phase. However, overfitting may occur as the T-LSTM
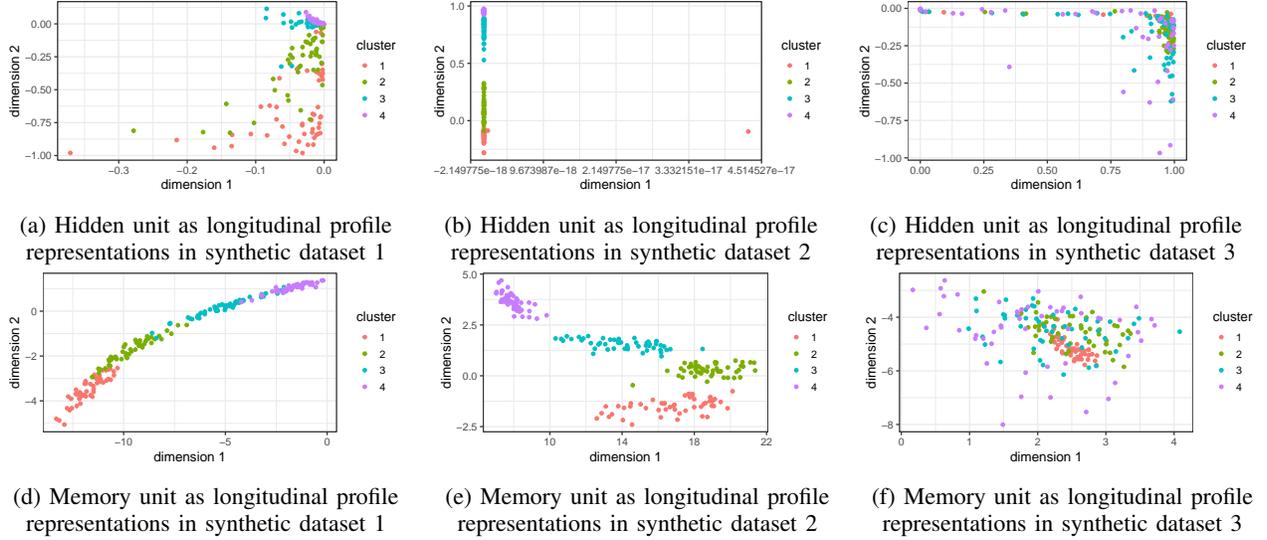
(a) Hidden unit as longitudinal profile representations in synthetic dataset 1

(b) Hidden unit as longitudinal profile representations in synthetic dataset 2

(c) Hidden unit as longitudinal profile representations in synthetic dataset 3

(d) Memory unit as longitudinal profile representations in synthetic dataset 1

(e) Memory unit as longitudinal profile representations in synthetic dataset 2

(f) Memory unit as longitudinal profile representations in synthetic dataset 3

Fig. 4: Visualization of hidden unit and memory unit at the last time step of encoder of T-LSTM Autoencoder after being trained for $20,000$ epochs (Best view in color)
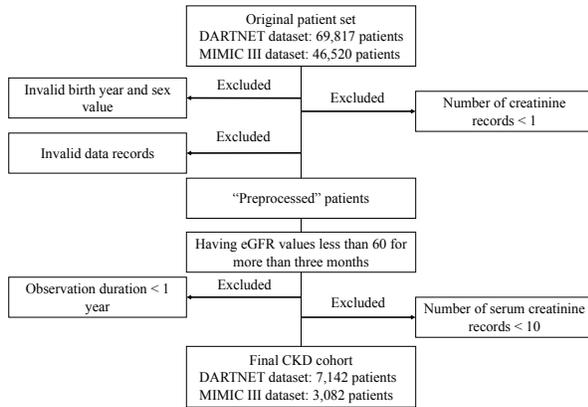


Fig. 5: Preprocessing process to obtain CKD cohorts

Autoencoder may have very small reconstruction error for a training set but fail to provide good reconstruction for a new set of longitudinal profiles that have not been trained before.

For this reason, we use 4-fold cross-validation to determine the appropriate dimension of hidden/memory unit. In particular, the longitudinal profiles in each dataset are split into four parts. T-LSTM model is trained on each three parts (training set) while being validated on the remaining part (validation set). The root-mean-square error (RMSE) between the reconstructed values and original values in the validation set is the measure for quality of T-LSTM Autoencoder. The lower of RMSE in validation set, the better quality of T-LSTM Autoencoder is as it can

generalize well to unseen data. Finally, the overall RMSE is computed by using the RMSE computed for each fold in the cross-validation process, i.e. $\sqrt{\frac{RMSE_1^2 + \cdots + RMSE_4^2}{4}}$.

We perform the above 4-fold cross validation on DARTNet and MIMIC-III datasets and vary the dimension of hidden/memory unit across following values $\{2^2, 2^3, 2^4, \cdots, 2^9\}$. In this experiment, we train T-LSTM Autoencoder with $500$ epochs. Figure 6 shows the overall RMSE of validation set with different values of hidden/memory dimension. For MIMIC-III dataset, we can observe in Figure 6b that $64$ is the optimal value for dimension of hidden/memory unit as the reconstruction error starts increasing after $64$, signaling that T-LSTM Autoencoder starts suffering from overfitting. For DART-Net dataset, the overall RMSE keeps decreasing as we increase the dimension of hidden/memory unit, indicating that overfitting is not a problem in this situation. However, as we can observe in Figure 6a, the overall RMSE decreases rapidly until reaching dimension of value $64$ and becomes saturated for dimension values larger than $64$. Based on the above observations, we choose $64$ as the dimension of hidden/memory unit when applying T-LSTM Autoencoder on two CKD datasets in subsequent experiments. As we discuss earlier in Section III, we will use both hidden and memory unit at the last time step of encoder to represent a longitudinal profile. Therefore, the dimension of our longitudinal profiles derived from T-LSTM Autoencoder is $128$ including $64$ dimensions from hidden unit and another $64$ dimensions from memory unit.
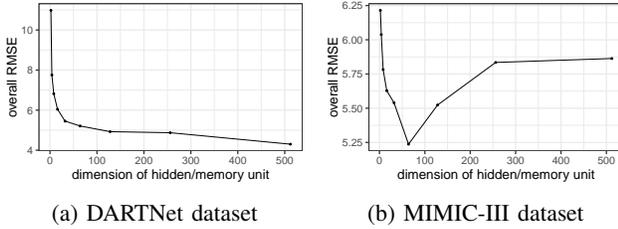
(a) DARTNet dataset      (b) MIMIC-III dataset

Fig. 6: Overall RMSE of validation set with respect to different dimension of hidden/memory unit

### C. Outputs of decoder

As we described earlier, the outputs of decoder in T-LSTM Autoencoder are actually the prediction of the original values in the longitudinal profile in reverse chronological order (see Figure 2 and Section II-D). One of the expectations when applying T-LSTM Autoencoder into CKD is that the disease progression reconstructed from the model should be able to capture the same progression of the original longitudinal profile. In this section, we investigate the reconstructed time series obtained from outputs of the decoder and observe how similar they are to the original time series. Figure 7 shows the outputs of the decoder in comparison with the original longitudinal data for three particular patients in each dataset. In the figure, the reconstructed time series from outputs of decoder can capture the long-term progression in patient CKD profile. In addition, for some particular time stamps in which the original CKD profile experiences sudden jumps in eGFR value, we can observe that the reconstructed time series have smoother transitions from one time step to another time step. In other words, we can view output of decoder in T-LSTM Autoencoder as a denoising summarizer that can capture the trends in CKD progression while simplifying the time series by reducing short-term variance.

In order to verify this observation for the whole dataset, we perform a statistical test to see whether the variance of the reconstructed time series is significantly smaller than the variance original time series. We first compute the variance in original time series and in reconstructed time series for each patient profile. After that, we use the pairs of variances, one for each patient in the dataset to perform a t-test to test whether there is a significant larger amount of variance from the original time series in comparison with the reconstructed one. The p-values of this t-test are both less than $2.2 \times 10^{-16}$ for both DARTNet and MIMIC-III dataset. This shows that the variance of reconstructed time series is significantly less than variance in the original time series, which confirms our observation.

### D. Embedded representation

In this section, we investigate the embedded patient profiles obtained by T-LSTM Autoencoder. As described in Section IV-B, each patient profile is actually represented by a 128 dimensional vector. When analyzing the progression in CKD, it is important to detect patients with interesting and unusual disease progressions. In other words, this task can be viewed as detecting outlier time series from a dataset of many irregularly sampled time series. By applying T-LSTM Autoencoder, each of the longitudinal profile has been transformed into a vector of fixed dimension. Each dimension in this embedded representation is an encoding of complex temporal features of the original longitudinal profile. One way to extract the outlier from this set of longitudinal profiles is to look at each individual dimension and check for the value that is farthest from the mean value of this dimension. In our experiment, by following this approach, we have found a set of few longitudinal profiles that they contain extreme values (farthest from the mean of a particular dimension) in multiple dimensions. Figure 8 shows the set of individual profiles that contain "outliers" in multiple different dimensions. From this figure, we can identify some of the uniqueness in those longitudinal profiles.

First observation we have in Figure 8 is that extreme length of the sequence can cause the embedded representation to have extreme values. In fact, patient $*476$ and $*598$ have two longest sequences in DARTNet dataset with 107 and 102 eGFR readings respectively. Similarly, patient $*33$ and $*18$ are two longest sequences in MIMIC-III dataset with 877 and 655 eGFR readings respectively.[1]

The high fluctuation in consecutive eGFR readings may also be a contributing factor to cause the extreme values. In particular, patient profiles in Figure 8c is highly unusual in the sense that eGFR readings fluctuate in a very wide range between 0 (severe kidney damage) and 100 (healthy kidney function). This can be a result of incorrect reading values by mistakes of input software when entering inputs to the EHR system. We can also observe the fluctuation in Figure 8e of patient $*18$. This patient profile has highly fluctuated observations over a long period of time. However, the range of fluctuation for the observations of this patient is not as extreme as previous patient. Therefore, we suspect that the measurement error instead of input error is the cause of this fluctuated observations.

Another interesting profile is of patient $*14$ in Figure 8f. This patient has a stable eGFR readings over the course of more than 100 days and then unexpectedly experiences a sudden increase in eGFR to more than 100 in the last

---

[1]On average, a patient profile in DARTNet dataset has only 15.8 eGFR observations while a patient in MIMIC-III dataset has only 49.8 observations.
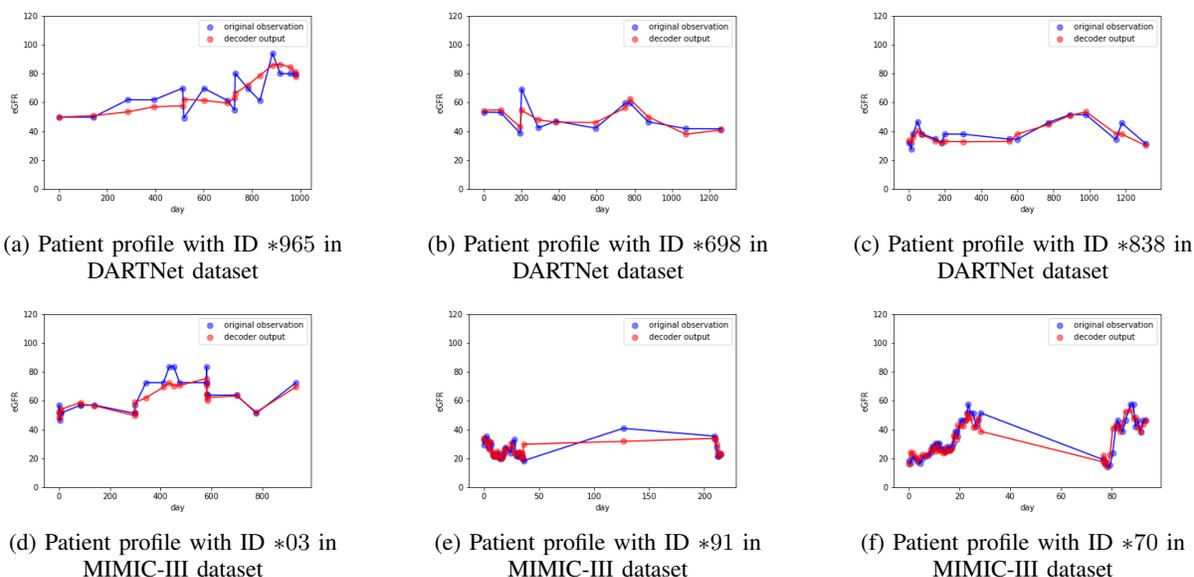
(a) Patient profile with ID ∗965 in DARTNet dataset

(b) Patient profile with ID ∗698 in DARTNet dataset

(c) Patient profile with ID ∗838 in DARTNet dataset

(d) Patient profile with ID ∗03 in MIMIC-III dataset

(e) Patient profile with ID ∗91 in MIMIC-III dataset

(f) Patient profile with ID ∗70 in MIMIC-III dataset

Fig. 7: Reconstructed time series in comparison with original time series. (Best view in color)



(a) Patient ∗476 with extreme values in 15 different dimensions

(b) Patient ∗598 with extreme values in 11 different dimensions

(c) Patient ∗534 with extreme values in 7 different dimensions

(d) Patient ∗33 with extreme values in 16 different dimensions

(e) Patient ∗18 with extreme values in 15 different dimensions

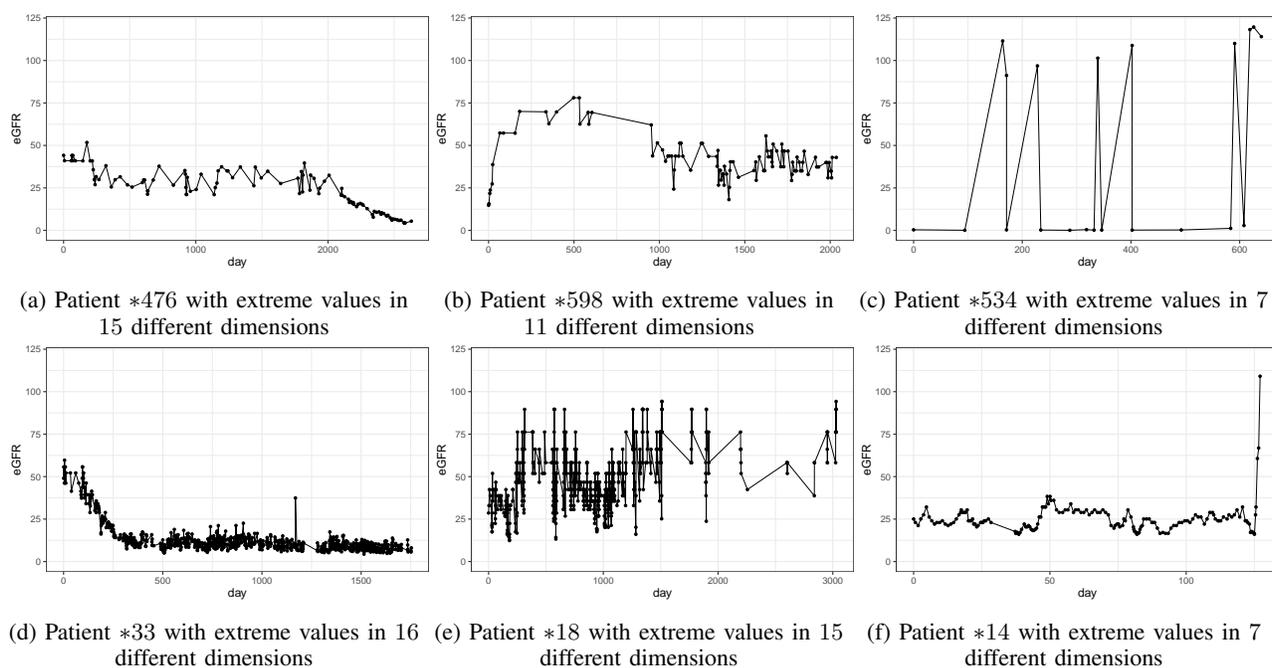(f) Patient ∗14 with extreme values in 7 different dimensions

Fig. 8: Individual CKD profiles which embedded representations have extreme values in multiple dimensions. Patients in Figure (a)-(c) are in DARTNet dataset. Patients in Figure (d)-(f) are in MIMIC-III dataset.

few days of his/her records. This is an interesting case as his/her course of disease is different from our understanding of CKD disease progression. Normally, kidney function in CKD patient deteriorates over the time and as a result, eGFR values will decrease over the time. However, this patient experiences a sudden improvement in eGFR readings, to a level of normal person with healthy kidney function, which is contrary to normal CKD disease progression. Therefore, this patient is a good candidate to further study his/her course of disease progression, which can subsequently help us expand our knowledge about CKD.

## V. DISCUSSION AND RELATED WORKS

One thing to note in T-LSTM Autoencoder is that its training time is highly dependent on the characteristics of training dataset as each step of the optimization requires the training time series to have the same length. For this reason, we have to split the training data into many batches, and each batch contains the set of time series that have the same length. This causes the training process to be dependent on the training dataset. In particular, for a dataset with many sequences having the same length, the training process will be efficient if it can take advantage of parallel computation when training multiple sequences together in each batch. On the other hand, for a dataset with many sequences having very different lengths (which is true for many longitudinal clinical datasets), the training process is reduced to a similar form of stochastic gradient-based optimization as most of the batches contain only one sequence.

In principle, T-LSTM Autoencoder can be used for longitudinal data with multi-dimensional observation. In our analysis on CKD, there is only one dimensional observation being used: eGFR - a main clinical marker of CKD. One may extend the model to have multi-dimensional observation by including other health indicators such as blood pressures, measures of liver functions, measures of body cholesterol level, etc. However, such an extended model may have to deal with missing values in the observation. For example, a patient who has observation of eGFR on a particular day may not have blood pressures and other clinical markers measured on the same day. In a study of Che et al. [11], authors proposed a modification of GRU architecture - another RNN variation that has been showed to perform comparably with LSTM - to handle missing values in clinical observations. However, the focus of Che et al.'s study is supervised learning in which it attempts to predict future events. In our future works, we will further study the capability of T-LSTM Autoencoder to represent missing values in the observation vector.

Med2Vec [12] and Wave2Vec [13] are two earlier studies that attempt to learn the representation of clinical observations by using unsupervised learning. In particular, both Med2Vec and Wave2Vec are inspired by Skip-gram model [14]. In the context of Natural Language Processing, the latent representations of words are learnt in a way that maximizes their co-occurrences with their nearby words. For clinical longitudinal dataset, both Med2Vec and Wave2Vec consider that the sequence of hospital visits for each patient is analogous to the sequence of words in a sentence in which we can learn the embedded representation of a hospital visit or a patient clinical observation by maximizing the co-occurrence of the current observation with few other observations prior to or after the current

one. However, instead of providing the embedded representations for the entire sequence of clinical observations like in the case of T-LSTM Autoencoder, Med2Vec and Wave2Vec only give an embedded representation for each clinical observation in the sequence.

In another study of Schulam and Arora [15], authors attempted to map irregularly sampled clinical time series into latent space by assuming patient's clinical observations follow Gaussian Processes. This study of Schulam and Arora focused on understanding whether there are a small number of degrees of freedom (dimensions of the latent space) that can explain the differences in the progressions of patient. In principle, T-LSTM Autoencoder can also be used to map the whole sequence of clinical observations into a small dimensional latent space, but at the expense of the higher reconstruction error.

## VI. CONCLUSION

In this paper, we have discussed various aspects of T-LSTM Autoencoder as a model to project longitudinal profiles into a latent space and its use on analyzing disease progression in CKD. With synthetic datasets, we demonstrated that both the memory and hidden unit at the last time step of the encoder should be used as the representation of longitudinal profile. In real-world CKD datasets, we presented an approach of using cross-validation to determine the dimension of hidden/memory unit in T-LSTM Autoencoder. In addition, we observed that the outputs of the decoder after being trained on CKD datasets are generally smoother in comparison with the corresponding inputs, while they still capture the trends expressed in the original inputs. Finally, we used the embedded representations of patient eGFR longitudinal profiles learnt from T-LSTM Autoencoder to identify unusual and interesting CKD profiles. The longitudinal profiles obtained by this experiment are candidates for two further analyses: (1) validate whether their inputs are correctly entered or not, and (2) in case the inputs are correctly entered, what is the underlying mechanism that causes the abnormality in the observations and whether we can enhance our understanding of CKD by understanding these unusual profiles.

### REFERENCES

[1] I. Abubakar, T. Tillmann, and A. Banerjee, "Global, regional, and national age-sex specific all-cause and cause-specific mortality for 240 causes of death, 1990-2013: a systematic analysis for the global burden of disease study 2013," *Lancet*, vol. 385, no. 9963, pp. 117–171, 2015.

[2] D. T. A. Luong *et al.*, "Extracting deep phenotypes for chronic kidney disease using electronic health records," *eGEMs*, vol. 5, 2017.

[3] D. T. A. Luong and V. Chandola, "A k-means approach to clustering disease progressions," in *IEEE ICHI*, Aug 2017, pp. 268–274.

[4] P. Singh, V. Chandola, and C. Fox, "Automatic extraction of deep phenotypes for precision medicine in chronic kidney disease," in *ICDH*, 2017, pp. 195–199.

[5] I. M. Baytas, C. Xiao, X. Zhang, F. Wang, A. K. Jain, and J. Zhou, "Patient subtyping via Time-Aware LSTM networks," in *KDD*. ACM, 2017, pp. 65–74.

[6] S. Wu, S. Liu, S. Sohn, S. Moon, C. il Wi, Y. Juhn, and H. Liu, "Modeling asynchronous event sequences with RNNs," *Journal of Biomedical Informatics*, vol. 83, pp. 167 – 177, 2018.

[7] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[8] P. J. Rousseeuw, "Silhouettes: a graphical aid to the interpretation and validation of cluster analysis," *Journal of computational and applied mathematics*, vol. 20, pp. 53–65, 1987.

[9] W. D. Pace, C. Fox, T. White, D. Graham, L. M. Schilling, and R. David, "The DARTNet institute: seeking a sustainable support mechanism for electronic data enabled research networks," *eGEMs*, vol. 2, no. 2, p. 6, 2014.

[10] A. E. Johnson, T. J. Pollard, L. Shen, H. L. Li-wei, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. A. Celi, and R. G. Mark, "MIMIC-III, a freely accessible critical care database," *Scientific data*, vol. 3, 2016.

[11] Z. Che, S. Purushotham, K. Cho, D. Sontag, and Y. Liu, "Recurrent neural networks for multivariate time series with missing values," *Scientific reports*, vol. 8, no. 1, p. 6085, 2018.

[12] E. Choi, M. T. Bahadori, E. Searles, C. Coffey, M. Thompson, J. Bost, J. Tejedor-Sojo, and J. Sun, "Multi-layer representation learning for medical concepts," in *KDD*. ACM, 2016, pp. 1495–1504.

[13] Y. Yuan, G. Xun, Q. Suo, K. Jia, and A. Zhang, "Wave2vec: Deep representation learning for clinical temporal data," *Neurocomputing*, 2018.

[14] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *NIPS*, 2013, pp. 3111–3119.

[15] P. Schulam and R. Arora, "Disease trajectory maps," in *NIPS*, 2016, pp. 4709–4717.