

Reconstructing Missing EHRs Using Time-Aware Within- and Cross-Visit Information for Septic Shock Early Prediction

1st Ge Gao

Department of Computer Science
North Carolina State University
Raleigh, USA
ggao5@ncsu.edu

2nd Farzaneh Khoshnevisan

Intuit Inc.
San Diego, USA
farzaneh_khoshnevisan@intuit.com

3rd Min Chi

Department of Computer Science
North Carolina State University
Raleigh, USA
mchi@ncsu.edu

Abstract—Real-world Electronic Health Records (EHRs) are often plagued by a high rate of missing data. In our EHRs, for example, the missing rates can be as high as 90% for some features, with an average missing rate of around 70% across all features. We propose a Time-Aware Dual-Cross-Visit missing value imputation method, named *TA-DualCV*, which spontaneously leverages multivariate dependencies across features and longitudinal dependencies both *within- and cross-visit* to maximize the information extracted from limited observable records in EHRs. Specifically, *TA-DualCV* captures the latent structure of missing patterns across measurements of different features and it also considers the time continuity and capture the latent temporal missing patterns based on both time-steps and irregular time-intervals. *TA-DualCV* is evaluated using three large real-world EHRs on two types of tasks: an *unsupervised imputation task* by varying mask rates up to 90% and a *supervised 24-hour early prediction of septic shock* using Long Short-Term Memory (LSTM). Our results show that *TA-DualCV* performs significantly better than all of the existing state-of-the-art imputation baselines, such as DETROIT and TAME, on both types of tasks.

Index Terms—Electronic Health Records(EHRs), EHRs Imputation, Septic Shock Early Prediction

I. INTRODUCTION

Sepsis, defined as life-threatening organ dysfunction in response to infection, is the leading cause of mortality and the most expensive condition associated with in-hospital stay, accounting for more than \$24 billion in annual costs in the United States [1]. In particular, *Septic shock*, the most advanced complication of sepsis due to severe abnormalities of circulation and/or cellular metabolism [2], reaches a mortality rate as high as 50% [3] and the annual incidence keeps rising [4]. It is estimated that as many as 80% of sepsis deaths could be prevented with early diagnosis and intervention; indeed prior studies have demonstrated that *early diagnosis* and treatment of septic shock can significantly decrease patients' mortality and shorten their length of stay [5]–[7]. Formerly, multiple complex patient health scoring systems have been defined and employed for early diagnosis and early intervention of sepsis, such as SOFA score [8] and MEDS [9]. Despite these efforts, early diagnosis of septic shock is still a

challenging problem, because of subtle but fast progression of sepsis/septic shock at early stages with lack of information.

On the other hand, the rapid growth in volume and diversity of Electronic Health Records (EHRs), makes it possible to apply many machine learning and data mining methods for disease diagnosis. EHRs are collections of multivariate time series clinical events recorded during patients' visits. *Each visit* consists of a sequence of events pertaining to the changing health status of the patient during the visit. EHRs typically acquire measurements irregularly. For example, vital signs are measured every 8 hours while lab values are measured only every 24 hours. Hence there may not be available readings for lab results when a new event is created for vital signs. Due to this integration of irregular data, EHRs usually contain high missing rates. For example, in this work, EHRs collected from general medical system, Christiana Care Health System and Mayo Clinic, have 73% and 68% missing rates, which are fairly typical of most real-world EHRs. This missingness in EHRs can severely limit the application of most data mining techniques that require complete data as inputs.

Our work further differentiates between two types of missing rates: *within-visit missing rate* and *cross-visit missing rate*. Essentially, the former refers to the average value of the missing portion of the data within one visit, while the latter refers to the average value of the missing portion of the data across all visits, and most existing missing rates in literature are cross-visit rates. Table II in section IV shows, in EHRs some features have a higher average within-visit missing rate while others have a higher cross-visit missing rate. Our goal here is to maximize the information extracted from limited observable records in EHRs by spontaneously leverages various dependencies both *within- and cross-visit*.

We propose a *Time-Aware Dual-Cross-Visit missing value imputation* approach, named *TA-DualCV*, to impute EHRs with high missing rates. Specifically, it takes advantages of all data points within EHRs by modeling *longitudinal and multivariate* dependencies both *within-visit* and *cross-visit*, to resolve the high missing rates in EHRs. *Multivariate dependencies* are captured among measurements across features,

and *longitudinal dependencies* consider the time continuity of measurements. Longitudinal dependencies can be captured from two perspectives: *temporal perspective* across time-steps and *time-aware perspective* across time-intervals (i.e., the elapsed time from a starting point till when each event occurs). Each type of dependencies can be modeled either *within-visit* or *cross-visit*. The core part of TA-DualCV is the *dual-cross-visit imputation* (DualCV), with which missingness can be modeled jointly at features and time steps across-visits using chained equations. In particular, it leverages cross-visit information from both feature-perspective and temporal-perspective. Additionally, we design a time-aware imputation (TA) within-visit, which captures longitudinal dependencies across the time intervals. DualCV is combined with TA visit-by-visit to ensure that dependencies are modeled across visits as well as within visits.

Various imputation methods have been proposed previously to handle missingness in EHRs. Based on the purpose of tasks, these prior approaches can be divided into those evaluated on *unsupervised learning tasks* and those designed towards *supervised learning tasks including disease diagnosis*. In the former, different approaches can be categorized into non-Neutral Network (NN)- such as 3D-MICE [10] and NN-based approaches such as DeEp impuToR Of mIssing Temporal data (DETROIT) [11], Time-aware Multi-modal Auto-Encoder (TAME) [12]. Our proposed TA-DualCV is a non-NN based approach using chained equations; as a result, it does not require specific characteristics or assumptions to be followed in EHRs, such as assumptions about the underlying density distributions. For the supervised learning tasks, missing indicator (MI) [13] show great success in disease progression modeling [14]–[17]. However, MI cannot be directly used for unsupervised learning tasks.

Our proposed TA-DualCV can be applied to both unsupervised and supervised learning tasks. Consequently, its effectiveness is assessed against the existing state-of-the-art imputation baselines that can be applied to both types of tasks. For *unsupervised learning tasks*, we assess imputation performance in terms of normalized root-mean-square error (nRMSE), by masking the observations with two commonly used strategies: masking with different rates, and masking one measurement per feature per visit. For *supervised learning tasks*, we utilize the imputed data incorporating Long Short-Term Memory (LSTM) network to predict septic shock 24 hours before its onset. LSTM has been shown to achieve the state-of-the-art results in many real-world applications including disease diagnosis [18] through deep hierarchical feature construction. Moreover, it can capture long-range dependencies in time series data in an effective manner [19].

Furthermore, the robustness of TA-DualCV is evaluated on EHRs from three different medical systems. Most existing methods are evaluated by using EHRs from single medical system, such as MIMIC-III (Medical Information Mart for Intensive Care) [20]. As EHR characteristics across different medical systems differ dramatically [21], it is yet unclear whether existing methods will hold up across different medical

systems. To summarize, our work has at least two main contributions:

- 1) To handle high missing EHRs, TA-DualCV integrates both cross-visit and within-visit dependencies, by exploiting dependencies among features, time-interval, and time-steps. As far as we know, TA-DualCV is the first *non-NN-based framework* designed for both unsupervised and supervised learning tasks in EHRs.
- 2) The generalizability and robustness of TA-DualCV are evaluated on both unsupervised imputation tasks and a supervised task, 24 hours early prediction of septic shock, using EHRs from three different medical systems. The results show that TA-DualCV outperforms state-of-the-art both non-NN- and NN-based baselines in EHRs with high missing rates.

II. RELATED WORK

A. Septic Shock Prediction

Various machine learning approaches, such as Temporal Belief Memory (TBM) [19] and Time-aware subGroup Basis Approach with Forecasted events (TGBA-F) [22], have been applied to missing EHRs to predict septic shock. For example, TBM [19] captures latent missing patterns based on irregular time intervals of EHRs and imputes missing values within a time window during recurrent neural network (RNN) based classifier training for septic shock prediction. TGBA-F [22] utilizes matrix decomposition for data imputation by considering forecasting future events, irregular time intervals, and patients subgroups, which is concurrent with RNN classifier training for septic shock prediction. Though both approaches can impute missing values, these need to combine a classifier and cannot be standalone for missing EHRs imputation. As our main focus is constructing a standalone imputing method that can be used for both unsupervised and supervised learning tasks, we don't compare TA-DualCV to these approaches.

B. Missing EHRs Imputation

Existing approaches of missing EHRs imputation in general can be divided into *those evaluated on unsupervised learning tasks* and *those designed towards supervised learning tasks*.

1) *Approaches evaluated on unsupervised learning tasks:* Table I summarizes the characteristics of some recently proposed imputation approaches. These approaches can be categorized into *non-Neutral Network (NN)-* and *NN-based approaches*. For example, MICE [23] captures multivariate dependencies using chained equations, and 3D-MICE [10] combines MICE and Gaussian process to integrate multivariate dependencies cross-visit and time-aware dependencies within-visit. Both MICE and 3D-MICE are non-NN-based and evaluated on MIMIC-III with a moderate missing rate (i.e., 17%) across their selected 13 features in [10]. Benefiting from chained equations, both approaches do not require specific characteristics or assumptions to be followed in EHRs, which motivates us to use chained equations for EHRs imputation.

More recently, NN-based approaches have been proposed and achieved better performance than existing non-NN-based

TABLE I
CHARACTERISTICS OF EXISTING APPROACHES VS. TA-DUALCV

Type	Approach	Native Missing Rate	Tensor Shape	Sliding Window	Prefill
Non-NN-Based	MICE (2010)	17%	No	No	No
	3D-MICE (2018)	17%	No	No	No
NN-Based	BRITS (2018)	78%	Yes	No	Yes
	DETROIT (2019)	7%	No	Yes	Yes
	GP-VAE (2020)	78%	Yes	No	No
	CATSI (2020)	7%	No	No	Yes
	TAME (2020)	54%	Yes	No	No
Non-NN-Based	TA-DualCV (Ours)	73%	No	No	No

- Native missing rate: it refers to the cross-visit native missing rate of the selected features in the original work. If the approach is evaluated on more than one EHRs, we report the highest one.

- For NN-based approaches, BRITS, CATSI, and TAME are bidirectional NN-based approaches and TAME achieves the best performance. Also, DETROIT is recognized as a strong baseline in [12]. Thus we use TAME and DETROIT as NN-based baselines in this study. GP-VAE requires tensor-shape input, thus cannot produce results under our EHRs with varied numbers of events across visits.

approaches. For example, DETROIT [11] prefills missing values with means and impute missing values based on neural networks by leveraging multivariate and time-aware dependencies from observed values within a 5-length sliding window, and it is evaluated on a subset of MIMIC-III with a small missing rate (i.e., 7%) across 13 lab analyte features. GP-VAE [24] takes tensor-shape data (i.e., a fixed number of events across visits) as model inputs and uses deep variational autoencoders to map the missing data into a latent space without missingness, where it models the low-dimensional dynamics with a Gaussian process. It is evaluated on PhysioNet data [25] with 78% missing rate across 35 features. BRITS [26] and CATSI [27] prefill missing values and adopt bidirectional recurrent neural networks (RNNs) to model data. BRITS is evaluated on PhysioNet data, and CATSI is evaluated on the same subset of MIMIC-III as DETROIT. BRITS further evaluates their work on in-hospital death classification, however, diseases diagnosis would be more desirable as important supervised learning tasks in healthcare that can guide clinicians construct treatments and intervention. Without prefilling, TAME [12] utilizes bidirectional RNNs and within-visit multi-modal embedding that takes data including demographics, diagnosis, medication, features, and time-intervals as inputs, and it outputs tensor-shape data without missingness. It is evaluated on MIMIC-III data with 54% average missing rate across 29 features.

Despite NN-based approaches’ superior performance in imputation, they have at least one of the three major limitations: i) *Tensor-shape input/output*. They require tensor-shape data as input or output, while in clinical practice visits under different health conditions can have a varied number of events. Truncating data into tensor shape can exclude substantial information. ii) *Sliding window*. They impute each missing value using observations within a fixed number of time steps around the missing value, which may not work on highly sparse data. iii) *Prefill*. They need to prefill the missing values to provide fixed-size inputs for NNs, while prefilling can introduce bias which limits model performance. Unlike these approaches, our framework is non-NN-based and does not require tensor-shape input/output, sliding window, or prefilling. BRITS and GP-VAE can handle high missing EHRs. However, they require

tensor-shape inputs. As real-world EHRs often contain varied numbers of events across visits, we don’t consider BRITS and GP-VAE as baselines in this work. Also, as BRITS, CATSI, and TAME are all bidirectional-RNNs-based among which TAME has achieved the best performance, and DETROIT has been recognized as a strong baseline in [11], we use TAME and DETROIT as baselines in this study.

2) *Approaches designed towards supervised learning tasks*: Missing indicator (MI) [13] shows great success in disease progression modeling [14]–[16]. However, it only indicates whether a value is observed or missing, thus cannot be directly used for imputation. In this work, we investigate incorporating imputed data with MI for septic shock prediction. Recently, TBM [19] and TGBA-F [22] impute the data concurrent with their classifier training. However, as we have discussed earlier, they cannot be standalone for missing EHRs imputation.

III. TA-DUALCV

Figure 1 shows the architecture of TA-DualCV. Specifically, a dual-cross-visit (DualCV) imputation method imputes the missing data using features and temporal dependencies, respectively, then the results from two streams are fused together. The results from DualCV are augmented with GPs which consider the within-visit time-aware dependencies, to constitute the final imputation results. We are releasing the source code of our implementation freely to the research community and it can be accessed at <https://github.com/fay067/TA-DualCV>.

A. Problem Formulation

1) *EHRs Imputation*: EHRs can be represented as $\mathbf{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_N\}$, where N is the total number of hospital visits. Each visit \mathbf{X}_i consists of a sequence of events: $\mathbf{X}_i = \{\mathbf{X}_{i,1}, \dots, \mathbf{X}_{i,T_i}\}$ where $\mathbf{X}_{i,j} \in \mathbb{R}^D$ represents an event recording measurements across D features at time step j in patient i ’s record. T_i is the number of events (i.e., time-steps) in visit i that can vary across different visits. Each visit i is also associated with a vector $\mathcal{T}_i = (\tau_{i,1}, \dots, \tau_{i,T_i}) \in \mathbb{R}^{T_i}$ recording the time intervals (e.g., minutes) from a starting point (e.g., patient arrival time) till when each event occurs.

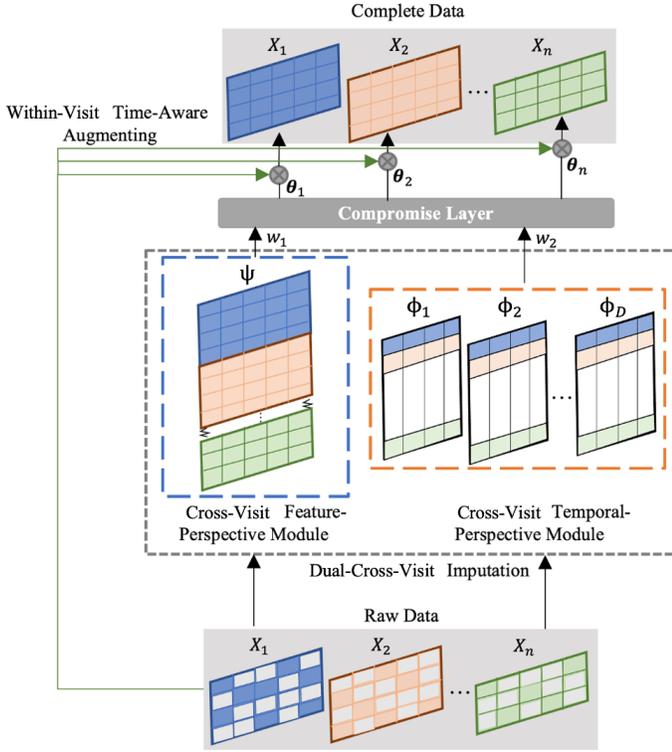


Fig. 1. Schematic of TA-DualCV. It consists of two parts: 1) Dual-cross-visit (DualCV) imputation that consists of cross-visits feature-perspective (CFP) and cross-visits temporal-perspective (CTP) modules. We get complete data from a compromise layer which combines two separate results from both modules; 2) Within-visit time-aware imputation (TA). The result from DualCV is augmented visit-by-visit with result from TA.

Moreover, the events of \mathbf{X}_i are in ascending order of the time intervals. In practice, \mathbf{X}_i may contain unknown *native missing values* due to the irregular data collection process. We apply additional masks to the observed data which introduce *masked missing values* with known groundtruth. Our goal is to devise an imputation model f which generates complete multivariate time series $\hat{\mathbf{X}}: \hat{\mathbf{X}} = f(\mathbf{X}, \mathcal{T})$.

B. DualCV Imputation

DualCV is designed to capture multivariate and temporal dependencies cross-visit, using chained equations. Chained equations are advantageous here because they do not require specific characteristics or assumptions to be followed, such as assumptions about the underlying density distributions. Specifically, we design two chained equation based modules, a CFP module to capture multivariate dependencies cross-visit, as well as a CTP module to capture longitudinal dependencies from temporal perspective cross-visit. We combine the results from the two modules to obtain imputation results, where we call dual-cross-sectional.

1) *Cross-Visits Feature-Perspective Module (CFP)*: To apply CFP, we first transform the data \mathbf{X} into a single matrix $\Psi = (\mathbf{X}_{1,\cdot}, \dots, \mathbf{X}_{N,\cdot})^\top$ by concatenating N visits on all features, where we then denote $\Psi_{j,\cdot}$ as the measurements on j -th event and $\Psi_{\cdot,k}$ represents measurements on k -th feature

(e.g., white blood cell count), as shown in the Figure 1 (CFP). Each visit can contain different numbers of events and time-intervals between any pair of consecutive events can be different. $\Psi = (\Psi_{\cdot,1}, \dots, \Psi_{\cdot,D})$ represents the observed part and $\Psi^{imp} = (\Psi_{\cdot,1}^{imp}, \dots, \Psi_{\cdot,D}^{imp})$ represents the missing part of the visits to be imputed. The predictors of $\Psi_{\cdot,k}$ is denoted as $\Psi_{\cdot,-k} = (\Psi_{\cdot,1}, \dots, \Psi_{\cdot,k-1}, \Psi_{\cdot,k+1}, \dots, \Psi_{\cdot,D})$ showing the collection of measurements on the $(D-1)$ features except feature k . $\Psi_{\cdot,-k}$ can be incomplete and the correlation between any pair of $\Psi_{\cdot,k}$ and $\Psi_{\cdot,-k}$ can be complex (e.g. nonlinear). Therefore, the hypothetically complete data $\hat{\Psi} = (\hat{\Psi}_{\cdot,1}, \dots, \hat{\Psi}_{\cdot,D})$ can be sampled from the D -variate multivariate distribution $P(\Psi_{\cdot,1}, \dots, \Psi_{\cdot,D} | \gamma)$ where γ are parameters [23]. Then the chain equations are designed to get the posterior distribution of γ by sampling iteratively from the conditional distribution

$$P(\Psi_{\cdot,1} | \Psi_{\cdot,-1}, \gamma_1) \dots P(\Psi_{\cdot,D} | \Psi_{\cdot,-D}, \gamma_D). \quad (1)$$

Specifically, the t -th iteration of chained equations is a Gibbs sampler which draws:

$$\begin{aligned} \gamma_1^{imp(t)} &\sim P(\gamma_1 | \Psi_{\cdot,1}, \Psi_{\cdot,-1}^{t-1}), \\ \Psi_{\cdot,1}^{imp(t)} &\sim P(\Psi_{\cdot,1} | \Psi_{\cdot,-1}, \Psi_{\cdot,-1}^{t-1}, \gamma_1^{imp(t)}), \\ &\vdots \\ \gamma_D^{imp(t)} &\sim P(\gamma_D | \Psi_{\cdot,D}, \Psi_{\cdot,-D}^{t-1}), \\ \Psi_{\cdot,D}^{imp(t)} &\sim P(\Psi_{\cdot,D} | \Psi_{\cdot,-D}, \Psi_{\cdot,-D}^{t-1}, \gamma_D^{imp(t)}). \end{aligned} \quad (2)$$

$\hat{\Psi}_{\cdot,k}^t = (\Psi_{\cdot,k}, \Psi_{\cdot,k}^{imp(t)})$ is defined as the complete data of k -th feature across visits at iteration t . Thus, $\hat{\Psi}^t = (\hat{\Psi}_{\cdot,1}^t, \dots, \hat{\Psi}_{\cdot,D}^t)$ is the complete data after t -th iteration.

2) *Cross-Visits Temporal-Perspective Module (CTP)*: To apply CTP, we then re-formulated the data \mathbf{X} into D matrices $\Phi = \{\Phi_1, \dots, \Phi_D\}$ by transposing each visit \mathbf{X}_i and concatenating N visits over the time-steps (e.g., $1, \dots, T_i$) as shown in Figure 1 (CTP). We fix the number of events for each visit to T^{med} , which denotes the median value of $\{T_i | i \in [1, N]\}$. Thus, $\Phi_j \in \mathbb{R}^{N \times T^{med}}$ denotes the matrix recording measurements on j -th feature, where $\Phi_{j(i,\cdot)}$ represents measurements of i -th visit on j -th feature and $\Phi_{j(\cdot,m)}$ represents measurements of the m -th time-step on j -th feature. If the number of time-steps within a record is less than T^{med} , the visit is padded to T^{med} by filling with placeholder values representing missingness (e.g., NaN's). If the number of time-steps is greater than T^{med} , the length of the visit is truncated to T^{med} from the front. The reason we use median values is that it is not sensitive to outliers. Also, as the computational complexity of chained-equation based algorithms is directly correlated with the number of variables [23], by using median it can balance between efficiency and the number of variables.

Similar to the first step described in Section III-B1, we apply chained equations on each Φ_j separately. Specifically, measurements on the m -th time-step of j -th feature is denoted by $\Phi_{j(\cdot,m)}$ while its complementary part is denoted as $\Phi_{j(\cdot,m)}^{imp}$ which is missing and subjected to imputation. The complete

data (i.e., after imputation) on j -th feature is denoted as $\hat{\Phi}_j = (\hat{\Phi}_{j(\cdot,1)}, \dots, \hat{\Phi}_{j(\cdot, T^{med})})$ which can be sampled from the T^{med} -variate multivariate distribution $P(\Phi_{j(\cdot,1)}, \dots, \Phi_{j(\cdot, T^{med})} | \beta)$. The conditional distributions used to obtain the posterior distribution of Φ_j , parametrized by β , are

$$\begin{aligned} P(\Phi_{j(\cdot,1)} | \Phi_{j(\cdot,-1)}, \beta_1), \dots, \\ P(\Phi_{j(\cdot, T^{med})} | \Phi_{j(\cdot, -T^{med})}, \beta_{T^{med}}), \end{aligned} \quad (3)$$

where $\Phi_{j(\cdot, -m)} = (\Phi_{j(\cdot,1)}, \dots, \Phi_{j(\cdot, m-1)}, \Phi_{j(\cdot, m+1)}, \dots, \Phi_{j(\cdot, T^{med})})$. Then the Gibbs sampler is used to sample β and Φ_j iteratively following the chained equations, at the s -th iteration, as

$$\begin{aligned} \beta_1^{imp(s)} &\sim P(\beta_1 | \Phi_{j(\cdot,1)}, \Phi_{j(\cdot,-1)}^{s-1}), \\ \Phi_{j(\cdot,1)}^{imp(s)} &\sim P(\Phi_{j(\cdot,1)} | \Phi_{j(\cdot,1)}, \Phi_{j(\cdot,-1)}^{s-1}, \beta_1^{imp(s)}), \\ &\vdots \\ \beta_{T^{med}}^{imp(s)} &\sim P(\beta_{T^{med}} | \Phi_{j(\cdot, T^{med})}, \Phi_{j(\cdot, -T^{med})}^{s-1}), \\ \Phi_{j(\cdot, T^{med})}^{imp(s)} &\sim P(\Phi_{j(\cdot, T^{med})} | \Phi_{j(\cdot, T^{med})}, \Phi_{j(\cdot, -T^{med})}^{s-1}, \beta_{T^{med}}^{imp(s)}). \end{aligned} \quad (4)$$

$\hat{\Phi}_{j(\cdot, m)}^s = (\Phi_{j(\cdot, m)}, \Phi_{j(\cdot, m)}^{imp(s)})$ is the complete data on j -th feature m -th time-step at iteration s . Then the complete (i.e., imputed) data on j -th feature, after the s -th iteration, is obtained as $\hat{\Phi}_j^s = (\hat{\Phi}_{j(\cdot,1)}^s, \dots, \hat{\Phi}_{j(\cdot, T^{med})}^s)$.

3) *Compromising Layer*: We now have two complete results, $\hat{\Psi}$ and $\hat{\Phi}$, obtained using CFP and CTP which considers the dependencies captured from the feature and temporal space respectively. Then $\hat{\Psi}$ and $\hat{\Phi}$ are fused together to constitute the complete data $\hat{\mathbf{X}}^M$. Specifically, the value on i -th visit, j -th time-step, and k -th feature of $\hat{\mathbf{X}}^M$, $\hat{X}_{i,j,k}^M$, is obtained as following

$$\begin{aligned} \hat{X}_{i,j,k}^M &= [w_1 \hat{\Psi}_{\sum_{x=1}^{i-1} T_{x,k}} + w_2 \hat{\Phi}_{k(i,j)}^s] \mathbb{1}\{T_i \leq T^{med}\} \\ &\quad + \hat{\Psi}_{\sum_{x=1}^{i-1} T_{x,k}} \mathbb{1}\{T_i > T^{med}\}, \end{aligned} \quad (5)$$

where $\mathbb{1}$ is indicator function, w_1 and w_2 are hyperparameters s.t. $w_1 + w_2 = 1$.

C. Within-Visit Time-Aware Augmenting Mechanism

Besides DualCV, which captures the feature and temporal dependencies *across* visits, we employ GP to perform time-aware imputations on each visit to capture patient-specific correlations *within* each visit. The difference between CTP and within-visit time-aware imputation is that CTP focuses on ‘‘outer’’ dependencies among time-steps of events across samples within population, while GP captures the ‘‘inner’’ dependencies among time-intervals of events within each visit.

Specifically, we apply single-task GP imputation over each feature on visit \mathbf{X}_i individually. The observed data on the l -th feature of visit \mathbf{X}_i is denoted as $\mathbf{X}_{i,\cdot,l}$, with $\mathbf{X}_{i,\cdot,l}^{imp}$ representing its complementary missing part which is subject to be imputed. For l -th feature, we split the time-interval vector $\mathcal{T}_i^{(l)}$ into $\mathcal{T}_i^{obs(l)}$ and $\mathcal{T}_i^{imp(l)}$, where $\mathcal{T}_i^{obs(l)}$ and $\mathcal{T}_i^{imp(l)}$ denote the time-intervals corresponding to $\mathbf{X}_{i,\cdot,l}$ and $\mathbf{X}_{i,\cdot,l}^{imp}$ respectively. The target output $\mathbf{X}_{i,\cdot,l}^{imp}$ is modeled as:

$$\mathbf{X}_{i,\cdot,l}^{imp} = \mu + z(\mathcal{T}_i^{imp(l)}) \quad (6)$$

where μ is overall mean and $z(\mathcal{T}_i^{imp(l)})$ is a GP that we assume to reach a local optima on its covariance structure:

$$Cov(z(\mathcal{T}_i^{imp(l)}), z(\mathcal{T}_i^{imp(m)})) = \sigma_z^2 R_{lm} \quad (7)$$

In general, $\mathbf{X}_{i,\cdot,l}^{imp} = (\mathbf{X}_{i,\cdot,1}^{imp}, \dots, \mathbf{X}_{i,\cdot,D}^{imp})^\top$ has a multivariate normal distribution. We specifically follow the Gaussian correlation function defined by Ranjan et al. [28]:

$$R_{lm} = \prod_{k=1}^{t_i} \exp\{-\alpha_k | \mathcal{T}_{i,k}^{imp(l)} - \mathcal{T}_{i,k}^{imp(m)}|^2\} \text{ for all } l, m \quad (8)$$

where $\alpha = (\alpha_1, \dots, \alpha_{t_i})$ is a vector of hyperparameters.

According to the Equations 6, 7, and 8, we can use $(\mathcal{T}_i^{obs(l)}, \mathbf{X}_{i,\cdot,l}^{obs})$ to fit GP and estimate the α following the negative profile log-likelihood [29]

$$\begin{aligned} -2\log(L_\alpha) \propto \\ \log(|R|) + n\log[(\mathbf{X}_{i,\cdot,l}^{obs} - \mathbf{1}_n \hat{\mu}(\alpha))^\top R^{-1} (\mathbf{X}_{i,\cdot,l}^{obs} - \mathbf{1}_n \hat{\mu}(\alpha))], \end{aligned} \quad (9)$$

where $|R|$ denotes the determinant of the Gaussian correlation function R defined in (8). Then we can get the predicted $\mathbf{X}_{i,\cdot,l}^{imp}$ by evaluating the GP model with input $\mathcal{T}_i^{imp(l)}$. Finally, we can obtain the complete (i.e., imputed) data of visit i as $\hat{\mathbf{X}}_i^G = (\mathbf{X}_i, \mathbf{X}_i^{imp})$ which constitutes the imputed dataset $\hat{\mathbf{X}}^G = (\hat{\mathbf{X}}_1^G, \dots, \hat{\mathbf{X}}_N^G)$.

Finally, we augment the results from DualCV with the outputs from GP through weighted averaging $\hat{\mathbf{X}}_i = \theta_{i,1} \cdot \hat{\mathbf{X}}_i^M + \theta_{i,2} \cdot \hat{\mathbf{X}}_i^G$, where $\theta_{i,1} = \frac{\sigma_{G_i}}{\sigma_{M_i} + \sigma_{G_i}}$, $\theta_{i,2} = \frac{\sigma_{M_i}}{\sigma_{M_i} + \sigma_{G_i}}$, σ_{M_i} and σ_{G_i} are the standard deviations of the imputations output from DualCV and within-visit time-aware imputations given visit i , respectively. Note the weight is designed to penalize the results (either from DualCV or GP) that are associated with the larger standard deviations.

IV. IMPUTATION WITH HIGH MISSING RATES

In this section, we first compare TA-DualCV with baselines in unsupervised learning tasks, imputation, by masking the observations with two commonly-used strategies separately:

1) Following [30], we vary random mask rates on 13 general laboratory analyse features from 60% up to 90% of observed values using widely-used MIMIC-III so as to evaluate the performance of imputation methods when the missing rate is high and close to the general clinical data.

2) Following [10], we randomly mask one measurement per feature per visit using 9 sepsis-related features in EHRs from three healthcare systems, CCHS, Mayo, and MIMIC-III, to study the performance of imputation and to obtain complete data for supervised learning tasks: septic shock early prediction.

Features 22 features, including 13 lab analyte and 9 sepsis-related features, are investigated in this study. The 13 lab analyte features are commonly used in existing works [10], [12], [31], [32]: Chloride, Potassium, Bicarb, Sodium, Hematocrit, Hemoglobin, MCV, Platelet, WBC, RDW, BUN, Creatinine, and Glucose. We use short forms CI, K, HCO3, Na, Hct, Hb, MCV, PLT, WBC, RDW, BUN, Cr, GLC, respectively,

TABLE II
MISSING RATES OF OVERALL 22 FEATURES IN CCHS, MAYO, AND MIMIC-III.

Feature	Cross-Visit			Within-Visit		
	MIMIC	CCHS	Mayo	MIMIC	CCHS	Mayo
CI	26%			26%		
K	21%			20%		
HCO3	27%			27%		
Na	24%			24%		
Hct	18%			19%		
Hb	30%			30%		
MCV	31%			31%		
PLT	28%			28%		
WBC	30%			30%		
RDW	31%			30%		
BUN	26%			25%		
Cr	25%			25%		
GLC	30%			30%		
Temp	97%	81%	91%	95%	82%	88%
RespRate	45%	63%	49%	61%	65%	54%
HeartRate	44%	61%	54%	61%	64%	55%
FiO2	84%	85%	94%	96%	84%	91%
PulseOx	45%	64%	34%	62%	68%	40%
OFlow	94%	78%	97%	94%	87%	92%
DBP	33%	72%	63%	25%	82%	64%
SBP	33%	72%	63%	25%	72%	64%
MAP	33%	77%	63%	25%	77%	68%
Median	31%	72%	63%	29%	77%	64%
Mean	39%	73%	68%	40%	76%	68%

- The median/mean refers to the median/average missing rate of the *selected features*, which can be different from the median/average missing rate across all features.

to represent these features for simplicity. The 9 sepsis-related features [22] are from four categories: 1) *Vital Signs*: Temperature, RespiratoryRate, HeartRate; 2) *Respiratory System*: FiO2, PulseOx; 3) *Cardiovascular System*: DiastolicBP, SystolicBP, MAP; 4) *Others*: OxygenFlow.

A. Datasets

- **MIMIC-III** [20]¹ contains admissions of adult patients (i.e., age > 16) who are admitted to intensive care units (ICU) at a tertiary referral hospital between 2001 and 2012, corresponding to 53,423 visits containing ~11 million events.
- **Christiana Care Health System (CCHS)** contains visit records of adult patients (i.e., age > 18) admitted to the Christiana Care hospital from July 2013 to December 2015, corresponding to 210,289 visits containing ~9 million events.
- **Mayo** contains visit records of adult patients (i.e., age > 18) admitted to Mayo Clinic Hospital from July 2013 to December 2015, corresponding to 121,019 visits containing ~52 million events.

1) *Data Preprocessing for Imputation with Varied Masking Rates*: With the selected 13 lab analyte features, to provide sufficient data for evaluation with high masking rates and to meet the base implementing requirements of baselines such as those that need a sliding window, we apply exclusion criteria for data sampling in MIMIC-III: i) exclude visits that have less than 11 events; ii) exclude events if they contain more than half

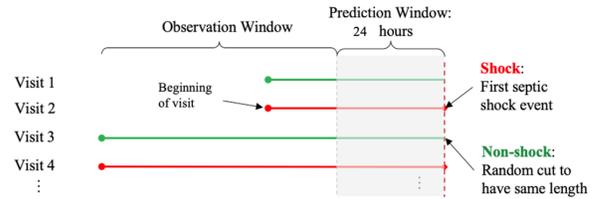


Fig. 2. Illustration of septic shock early prediction setting.

missing measurements. Thus, 19,693 visits containing 658,690 events are obtained from MIMIC-III.

2) *Ground Truth Labeling & Data Preprocessing for Septic Shock*: Supervised learning models depend heavily on the accurate label of the training set. However, acquiring the true label (i.e., septic shock and non-shock) can be challenging. Although diagnosis codes, such as International Classification of Diseases, Ninth Revision (ICD-9), are widely used for clinical labeling, solely relying on ICD-9 can be problematic as it has been proven to have limited reliability due to the fact that its coding practice is used mainly for administrative and billing purpose. Indeed, it has been widely argued that ICD-9 cannot be used for establishing reliable gold standards for various clinical conditions [33], [34]. More importantly, ICD-9 cannot tell when septic shock occurs at event level, which is essential for our task. On the basis of the Third International Consensus Definitions for Sepsis and Septic Shock [35], our domain experts identified septic shock as any of the following conditions are met:

- Persistent hypertension as shown through two consecutive readings (≤ 30 minutes apart).
 - Systolic Blood Pressure (SBP) < 90 mmHg
 - Mean Arterial Pressure (MAP) < 65 mmHg
 - Decrease in SBP ≥ 40 mmHg with an 8-hour period
- Any vasopressor administration.

When combing both ICD-9 and the domain experts' rules, we identify: i) 4,918 visits (2,459 shock positive visits and 2,459 negative visits) containing ~1 million events from MIMIC-III. ii) 144,119 visits (2,963 shock positive visits and 141,156 negative visits) containing 795,314 events from CCHS. iii) 84,897 visits (3,499 shock positive visits and 81,398 negative visits) containing 7,612,360 events from Mayo. We are given all of a patient's EHRs until 24 hours before the septic shock onset (shock group) or end of the sequences (non-shock group), and supervised learning task is to use the complete data to predict whether or not the patient will develop septic shock exactly 24 hours later. To conduct this task, we *right aligned* all the shock sequences by septic shock onset and all non-shock by the end of their sequences and include all the EHRs until 24 hours before the end of sequences (see Fig. 2). Our *prediction window* here is 24-hour window before the onset of septic shock or end of sequence. To simulate early prediction in practice, we exclude events in prediction window, accordingly the visits that contain no events in observation window are excluded for both imputation and septic shock

¹Dataset is available at <http://mimic.physionet.org>.

TABLE III
IMPUTATION RESULTS (nRMSE) WITH MASKING RATES 60%-90% ON 13 FEATURES OF MIMIC-III.

Model	CI	K	HCO3	Na	Hct	Hb	MCV	PLT	WBC	RDW	BUN	Cr	GLC	Average
Mask rate: 90%														
MeanFill	1.296	1.208	1.143	1.120	4.227	2.729	2.337	8.754	5.065	4.721	3.368	3.398	4.488	3.373
ECF	0.270	0.256	0.272	0.267	0.275	0.278	0.295	0.294	0.279	0.304	0.283	0.280	0.260	0.278
MICE	0.283	0.269	0.287	0.280	0.282	0.284	0.310	0.310	0.295	0.321	0.304	0.318	0.271	0.293
3D-MICE	0.259	0.253	0.262	0.259	0.262	0.263	0.270	0.274	0.270	0.265	0.267	0.261	0.256	0.263
DETROIT	4.457	0.289	0.311	2.855	0.421	0.371	0.981	0.578	2.121	0.303	1.720	2.014	1.585	1.385
TAME	0.381	0.388	0.318	0.323	0.387	0.380	0.734	0.436	0.465	0.441	0.482	1.290	0.365	0.491
TA-DualCV ^{-C}	0.284	0.276	0.285	0.284	0.288	0.290	0.296	0.300	0.296	0.289	0.292	0.284	0.279	0.288
TA-DualCV ^{-I}	0.254	0.264	0.255	0.257	0.258	0.261	0.266	0.267	0.277	0.257	0.257	0.260	0.267	0.261
TA-DualCV	0.221**	0.225**	0.224**	0.222**	0.229**	0.233**	0.232**	0.236**	0.232**	0.237**	0.229**	0.223**	0.218**	0.228**
Mask rate: 80%														
MeanFill	1.201	1.105	1.065	1.072	4.325	2.486	2.174	7.354	4.189	4.394	3.032	3.211	3.878	3.037
ECF	0.267	0.254	0.268	0.265	0.273	0.276	0.292	0.289	0.276	0.299	0.279	0.276	0.258	0.275
MICE	0.287	0.270	0.290	0.283	0.286	0.289	0.312	0.313	0.297	0.322	0.307	0.324	0.273	0.296
3D-MICE	0.257	0.250	0.259	0.256	0.259	0.261	0.267	0.273	0.267	0.264	0.264	0.258	0.253	0.260
DETROIT	0.272	0.257	0.256	0.257	0.256	0.251	0.407	0.287	0.321	0.270	0.314	1.544	0.243	0.380
TAME	0.317	0.294	0.314	0.317	0.389	0.402	0.513	0.527	0.602	0.624	0.343	0.484	0.389	0.424
TA-DualCV ^{-C}	0.277	0.272	0.280	0.279	0.283	0.284	0.290	0.292	0.290	0.280	0.286	0.279	0.274	0.282
TA-DualCV ^{-I}	0.246	0.251	0.250	0.247	0.242	0.242	0.268	0.262	0.262	0.271	0.256	0.259	0.257	0.255
TA-DualCV	0.222**	0.219**	0.224**	0.222**	0.228**	0.231**	0.235**	0.236**	0.233**	0.234**	0.225**	0.226**	0.221**	0.227**
Mask rate: 70%														
MeanFill	0.847	0.787	0.789	0.764	2.611	1.752	1.726	4.394	3.225	3.247	2.170	2.674	2.159	2.088
ECF	0.255	0.249	0.258	0.255	0.264	0.266	0.279	0.277	0.264	0.285	0.265	0.264	0.254	0.264
MICE	0.271	0.264	0.279	0.269	0.264	0.265	0.304	0.302	0.287	0.314	0.293	0.310	0.268	0.284
3D-MICE	0.243	0.242	0.248	0.245	0.244	0.243	0.251	0.256	0.254	0.246	0.250	0.245	0.246	0.247
DETROIT	0.239	0.244	0.301	0.266	0.227	0.220**	0.258	0.361	0.449	0.238**	0.790	0.849	0.263	0.362
TAME	0.274	0.320	0.281	0.265	0.237	0.235	0.497	0.326	0.370	0.410	0.330	0.509	0.444	0.346
TA-DualCV ^{-C}	0.255	0.249	0.258	0.255	0.264	0.266	0.279	0.277	0.264	0.285	0.265	0.264	0.254	0.264
TA-DualCV ^{-I}	0.231	0.235	0.240	0.234	0.223	0.225	0.247	0.248	0.256	0.246	0.238	0.242	0.247	0.239
TA-DualCV	0.220**	0.221**	0.221**	0.229**	0.221**	0.222	0.235**	0.231**	0.229**	0.239	0.225**	0.224**	0.214**	0.225**
Mask rate: 60%														
MeanFill	0.657	0.589	0.624	0.591	1.801	1.279	1.544	2.958	1.782	2.538	1.685	2.370	1.601	1.540
ECF	0.247	0.246	0.250	0.248	0.258	0.260	0.272	0.266	0.257	0.275	0.253	0.255	0.252	0.257
MICE	0.260	0.261	0.271	0.259	0.247	0.247	0.301	0.298	0.283	0.310	0.285	0.296	0.266	0.276
3D-MICE	0.233	0.234	0.239	0.235	0.229	0.228	0.240	0.244	0.245	0.232	0.238	0.236	0.238	0.236
DETROIT	0.154**	0.241	0.168**	0.154**	0.149**	0.196**	0.164**	0.481	0.458	0.154**	0.216**	0.411	0.170**	0.240
TAME	0.227	0.271	0.257	0.245	0.211	0.212	0.355	0.321	0.320	0.352	0.270	0.423	0.321	0.291
TA-DualCV ^{-C}	0.271	0.275	0.271	0.274	0.281	0.283	0.284	0.281	0.285	0.269	0.271	0.275	0.282	0.277
TA-DualCV ^{-I}	0.218	0.237	0.230	0.223	0.205	0.207	0.260	0.245	0.247	0.256	0.234	0.240	0.249	0.235
TA-DualCV	0.219	0.223**	0.220	0.224	0.218	0.218	0.229	0.230**	0.223	0.222	0.225	0.223**	0.209	0.222**

- For each block, the best approach is in **bold**; The best approach across ALL is labeled with ******.

prediction tasks. As a result, the final datasets have: i) 772 visits (386 shock positive visits and 386 negative visits) containing 37,246 events from MIMIC-III. Notably, MIMIC-III is ICU-specific and sepsis onset can present itself early, thus the number of remaining visits in MIMIC-III is relatively small after applying a 24-hour cut before the onset. ii) 2,842 visits (1,347 shock positive visits and 1,495 negative visits) containing 376,887 events from CCHS. iii) 6,836 visits (3,457 shock positive visits and 3,379 negative visits) containing 1,547,893 events from Mayo. Note that the labels are not used for imputation tasks, but only for septic shock prediction tasks in Section V.

3) *Missing Data Analysis*: Since different features are measured at different events, plenty of missing entries exist in EHRs. For instance, vital signs are generally measured every 8 hours, while lab values are measured every 24 hours. Hence there may not be available readings for lab results when a new event is created for vital signs. Table II shows both cross-visit and within-visit missing rates of 22 features selected from the three datasets in our study. We also calculate median and mean value of cross-visit and within-visit missing rates across features. The missingness across visits within one medical

system and across different medical systems can diff a lot.

B. Baselines

We compare TA-DualCV with six imputation approaches for comparison, which are originally designed towards either supervised or unsupervised learning tasks.²

Approaches towards supervised learning tasks:

- **MeanFill** imputes missing values by mean values cross-visit on each feature.
- **Expert Carry Forward (ECF)** [19] fills the missing values by the last value within a fixed-length window (i.e., 24 hours for lab analytes and 8 hours for vital signs). Then it fills the remaining with the mean value of the corresponding features. ECF is a strong method for both imputation and disease diagnosis in [19].

Approaches towards unsupervised learning tasks:

Non-NN-Based:

- **MICE** [23] imputes missing values by capturing multi-variate dependencies cross-visit.

²Note that Missing indicator (MI) is not included as a baseline in this task, as it only indicates whether a measurement is missing without imputation.

TABLE IV
IMPUTATION RESULTS (nRMSE) WITH CLASSIC MASKING RATES ON 9 SEPSIS-RELATED FEATURES OF CCHS, MAYO, AND MIMIC-III.

Model	Temp	RespRate	HeartRate	FiO2	PulseOx	OFlow	DBP	SBP	MAP	Average
CCHS										
MeanFill	1.929	0.624	1.704	2.267	0.655	12.927	1.092	1.404	0.867	2.825
ECF	1.872	0.564	1.498	2.267	0.655	12.927	1.092	1.404	0.867	2.785
MICE	1.224	0.644	0.744	0.391	0.572	0.557	2.564	1.644	1.268	1.043
3D-MICE	1.219	0.637	0.714	0.311	0.549	0.473	2.485	1.483	1.179	0.984
TAME	0.808	0.631	0.986	0.980	0.880	1.484	0.790	1.401	1.801	1.085
TA-DualCV ^{-C}	0.355	0.218**	0.487	0.343	0.374	0.236**	0.303	0.281	0.222	0.313
TA-DualCV ^{-I}	0.397	0.337	0.522	0.399	0.471	0.516	0.380	0.416	0.409	0.427
TA-DualCV	0.314**	0.222	0.334**	0.303**	0.346**	0.295	0.272**	0.267**	0.211**	0.285**
Mayo										
MeanFill	0.710	0.680	1.998	45.391	0.684	5.217	0.880	1.310	1.662	7.109
ECF	0.640	0.467	0.949	45.391	0.684	5.217	0.880	1.310	1.662	6.942
MICE	0.812	0.460	0.716	24.445	0.645	0.435	0.652	0.775	0.780	3.617
3D-MICE	0.820	0.452	0.716	24.445	0.546	0.427	0.549	0.778	0.779	3.592
TAME	0.812	1.526	1.101	31.405	0.955	1.329	0.587	0.730	0.743	4.806
TA-DualCV ^{-C}	0.677	0.397	0.542	0.528	0.469	0.246**	0.490	0.692	0.702	0.527
TA-DualCV ^{-I}	0.541	0.372	0.543	0.530	0.554	0.660	0.543	0.517	0.506	0.530
TA-DualCV	0.363**	0.276**	0.321**	0.329**	0.326**	0.310	0.316**	0.295**	0.209**	0.305**
MIMIC-III										
MeanFill	1.094	0.780	4.273	0.847	0.946	1.286	1.400	1.345	1.122	1.455
ECF	0.375	0.587	3.171	0.847	0.946	1.286	1.400	1.345	1.122	1.231
MICE	0.656	0.647	0.618	0.739	0.769	0.792	1.050	0.878	1.312	0.829
3D-MICE	0.607	0.642	0.684	0.712	0.759	0.729	1.076	0.895	1.305	0.823
TAME	0.287	0.291	0.253	0.747	0.246	0.451	0.140**	0.179	0.089**	0.297
TA-DualCV ^{-C}	0.228**	0.380	0.400	0.431	0.291	0.319	0.213	0.209	0.288	0.293
TA-DualCV ^{-I}	0.475	0.497	0.336	0.477	0.378	0.406	0.262	0.283	0.351	0.384
TA-DualCV	0.240	0.278**	0.226**	0.333**	0.223**	0.306**	0.189	0.174**	0.228	0.244**

- For each block, the best approach is in **bold**; The best model across ALL is labeled with **.
- DETROIT cannot produce results so we don't include it in the table.

- **3D-MICE** [10] imputes missing values by integrating multivariate dependencies cross-visit and time-interval dependencies within-visit.

NN-Based:

- **DETROIT** [11] prefills missing values with means and impute missing values based on neural networks by leveraging multivariate and time-aware dependencies from observed values within a 5-length window.
- **TAME** [12] imputes missing values based on bidirectional RNNs and within-visit multi-modal embedding that takes data including demographics, diagnosis, medication, features, and time-intervals as inputs.

To evaluate the performance of our proposed modules, we implement another two variant versions of TA-DualCV.

- **TA-DualCV^{-C}** removes DualCV but imputes based on within-visit time-aware information.
- **TA-DualCV^{-I}** removes within-visit time-aware augmentation that is used to augmenting results from DualCV.

C. Experimental Setup

To obtain results of baselines, for ECF, we implement the code strictly following the instructions in the authors' paper. For MICE, we use R package *mice* with 10 iterations. For 3D-MICE, DETROIT and TAME, we use the source code provided by the authors with their default parameter settings

as suggested in [30]. We implement our proposed work with R and Python. For parameters settings, we use 10 iteration with multiple chain equations. In compromising layer, we use $w_1 = w_2 = 0.5$ for all experiments, to provide a general view of TA-DualCV's performance taking balanced efforts from both cross-visit and within-visit imputation in this study.

1) *Evaluation Metric:* Following [10], [12], normalized root-mean-square error (nRMSE) is used to evaluate the imputation performance of the methods. We evaluate the imputation performance of all methods on masked values by calculating a popular matrix, normalized root-mean-square error (nRMSE). nRMSE assigns large residuals with a disproportionately large effect, thus more sensitive to outliers [36]. nRMSE is widely adopted in clinical data imputation tasks that contain different scales of values [10], [12]. As normalization of RMSE has different calculations in different approaches. In this study, we used the common choice of range normalization as in [10].

D. Results

Table III compares the various approaches' average nRMSE of 13 general lab analytes features across different masking rates on MIMIC-III. The best approaches for each block is indicated in **bold**; the best approaches across all blocks is indicated by **. In approaches towards supervised learning tasks, ECF outperforms MeanFill at all masking rates and

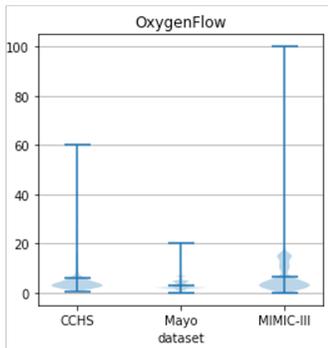


Fig. 3. Value distribution of OxygenFlow on CCHS, Mayo, and MIMIC-III. The blue-shaded area represents the density of the data. Blue-vertical lines represent the maximum, mean, and minimum values of the data.

performance holds steady. For approaches towards unsupervised learning tasks, in non-NN-based approaches, 3D-MICE outperforms MICE at all masking rates. In NN-based approaches, TAME can outperform DETROIT when masking rate is as high as 90%, while DETROIT outperforms TAME when masking rate is 60%. The nRMSE of DETROIT is increased faster when masking rate increased. A possible reason is that DETROIT needs to prefill missing values, which can introduce more bias when masking rate is high. 3D-MICE achieves the best performance among four unsupervised learning imputation methods. To our surprise, DETROIT and TAME, especially DETROIT, perform worse than 3D-MICE. There is a possible reason that the datasets in our experiments are highly sparse across visits. DETROIT predicts missing values according to neighboring measurements, so it may be difficult to obtain reliable results when the data is highly sparse and even fails to provide results when the number of events is small (e.g., 1). TAME can take in multi-modal information including demographics, diagnosis, medication, and lab analytes, and it has achieved superior performance in its original settings using multi-modal inputs containing 29 features with 54% average missing rate, as described in [12]. However, its performance can suffer when highly sparse 9 features are provided, as information is too limited to construct NN model. Finally, Table III shows that TA-DualCV achieves the best performance in terms of average nRMSE across all masking rates on this task.

On the CCHS, Mayo, and MIMIC-III datasets, Table IV provides the overall average nRMSE across nine sepsis-related features for the compared methods. Across all three medical systems, these approaches performances varied considerably. One explanation for this could be that the features are distributed differently in the three medical systems (e.g., OxygenFlow). ECF scores better than MeanFill across the three medical systems. 3D-MICE again performs well among four unsupervised learning imputation methods on CCHS and Mayo, while TAME performs well on MIMIC-III. Note that we were unable to make DETROIT converge to produce results. Moreover, 3D-MICE outperforms ECF across all three medical systems, and all methods perform better on CCHS

and MIMIC-III than Mayo. Finally, Table IV shows that TA-DualCV outperforms all baselines as it not only achieves the least overall nRMSE among all three medical systems but its performance is relatively stable.

In summary, Table III and Table IV show that TA-DualCV is capable of handling a diverse set of feature imputation tasks in both general and ICU-specific real-world medical systems when missing rate is high and achieving state-of-the-art performance.

V. SEPTIC SHOCK EARLY PREDICTION

To examine how our proposed approach performs on a supervised task, we utilize the complete data with nine sepsis-related features generated by implementing different imputation methods, to build an LSTM network to predict an extremely challenging condition in practice, septic shock. We predict septic shock 24 hours before its onset, as suggested by domain experts to enable early treatments that can prevent 80% of deaths [7], [37], [38]. Furthermore, we combine each imputation method with missing indicator (MI), for prediction and exploit *to what extent the missing pattern information can affect the prediction results*. The missing indicator $MI_{i,j,k}$ of a missing value $X_{i,j,k}$ is represented as

$$MI_{i,j,k} = \mathbb{1}\{X_{i,j,k} \text{ is missing}\}.$$

A. Experimental Setup

In this experiment, we use the three datasets that containing events within their observation window as described in Section IV-A and the same baselines as the ones from Section IV-B except DETROIT, as it does not produce complete data as described in Section IV-D. For parameters settings, in LSTM, we use a general setting with 1 hidden layer with 128 hidden neurons, 0.005 initial learning rate. We adapt Adam optimizer [39] with 64 batches and 25 epochs.

1) *Evaluation Matrix*: Metrics of accuracy (Acc), recall, precision (Prec), F-score (F1) and area under the ROC curve (AUC) were employed for evaluating our models. Accuracy is the proportion of patients whose labels are correctly identified. Recall indicates what proportion of patients that actually have septic shock can be correctly diagnosed by the model as septic shock. Precision tells what proportion of patients who are diagnosed as septic shock actually have septic shock. F1 is the harmonic mean of precision and recall that sets their trade-off. AUC calculates the tradeoff between recall and specificity. Experiment results are averaged from 5-fold cross-validation.

B. Results

Table V presents the 24-hours early prediction results of septic shock without and with missing indicators. Either without or with missing indicators, there is no clear winner between the two approaches towards supervised learning tasks. Among the three approaches towards unsupervised learning tasks, in general, MICE and 3D-MICE perform better than TAME. Interestingly, TAME performs worse than other baselines without missing indicators, probably because it outputs tensor-shape data and some useful information regarding septic

TABLE V
SEPTIC SHOCK 24-HOUR EARLY PREDICTION RESULTS (MEAN±STD) USING 9 SEPSIS-RELATED FEATURES ON CCHS, MAYO, AND MIMIC-III.

Method	Raw					Combined with Missing Indicators				
	Accuracy	Precision	Recall	F1	AUC	Accuracy	Precision	Recall	F1	AUC
CCHS										
MeanFill	.632±.025	.637±.032	.611±.036	.623±.036	.679±.021	.682±.019	.673±.031	.710±.058	.689±.058	.740±.026
ECF	.631±.023	.636±.024	.602±.079	.617±.079	.676±.023	.673±.020	.668±.018	.687±.044	.677±.044	.729±.023
MICE	.672±.020	.670±.026	.680±.048[‡]	.673±.048	.729±.033	.708±.043	.713±.059	.692±.050	.703±.050	.770±.038[‡]
3D-MICE	.666±.024	.673±.015	.643±.058	.656±.058	.726±.034	.704±.049	.707±.056	.692±.086	.697±.086	.766±.054
TAME	.482±.010	.291±.238	.600±.490	.392±.490	.513±.039	.643±.032	.633±.048	.714±.074[‡]	.666±.074	.707±.016
TA-DualCV ^{-C}	.679±.024	.688±.032	.652±.055	.668±.055	.726±.026	.714±.017 [‡]	.714±.018 [‡]	.711±.027	.712±.027 [‡]	.769±.025
TA-DualCV ^{-I}	.690±.025**	.691±.024 [‡]	.689±.051**	.689±.051**	.730±.036 [‡]	.701±.033	.701±.041	.701±.039	.700±.039	.765±.044
TA-DualCV	.687±.037 [‡]	.700±.033**	.662±.034	.674±.034 [‡]	.741±.020**	.721±.035**	.716±.039**	.729±.044**	.721±.055**	.777±.031**
Mayo										
MeanFill	.695±.020	.699±.033	.690±.025	.693±.025	.783±.016	.722±.020	.711±.039	.752±.055	.729±.055	.805±.014
ECF	.701±.014	.699±.023	.705±.026	.702±.026	.770±.014	.733±.021	.730±.033	.739±.029	.734±.029	.816±.022
MICE	.711±.014	.715±.015	.701±.044	.707±.044	.786±.016	.727±.009	.726±.027	.735±.032	.729±.032	.814±.011
3D-MICE	.697±.016	.706±.030	.673±.036	.689±.036	.778±.007	.735±.009	.732±.023	.743±.018	.737±.018	.819±.009
TAME	.515±.034	.523±.029	.576±.268	.510±.268	.549±.028	.730±.017	.710±.033	.780±.041**	.742±.041	.787±.010
TA-DualCV ^{-C}	.713±.021	.717±.035 [‡]	.703±.043	.709±.043	.793±.018	.735±.016	.735±.028	.739±.050	.735±.050	.819±.013
TA-DualCV ^{-I}	.722±.014**	.716±.007	.733±.052**	.724±.052**	.796±.013**	.746±.006 [‡]	.740±.028 [‡]	.762±.047[‡]	.749±.038 [‡]	.834±.005**
TA-DualCV	.717±.017 [‡]	.732±.022**	.724±.026 [‡]	.719±.026 [‡]	.795±.018 [‡]	.752±.009**	.760±.023**	.751±.043	.750±.043**	.830±.010 [‡]
MIMIC-III										
MeanFill	.816±.037	.836±.039	.789±.084	.808±.084	.889±.028	.848±.025	.839±.036	.863±.047	.850±.047	.901±.010
ECF	.834±.041	.848±.052	.813±.061	.829±.061	.899±.023	.845±.034	.845±.022	.846±.088	.842±.088	.907±.018[‡]
MICE	.837±.037	.850±.059[‡]	.826±.046	.835±.046	.887±.026	.846±.020	.855±.028[‡]	.836±.061	.844±.061	.904±.011
3D-MICE	.825±.015	.831±.053	.823±.045	.824±.045	.890±.017	.846±.022	.847±.024	.847±.050	.846±.050	.906±.014
TAME	.510±.025	.616±.196	.493±.349	.436±.349	.581±.058	.782±.144	.739±.128	.957±.042	.825±.042	.845±.142
TA-DualCV ^{-C}	.825±.021	.846±.036	.796±.036	.819±.036	.892±.021	.845±.032	.851±.034	.839±.072	.842±.072	.905±.015
TA-DualCV ^{-I}	.838±.035 [‡]	.844±.030	.829±.058 [‡]	.836±.058 [‡]	.900±.025 [‡]	.854±.021 [‡]	.847±.037	.863±.038	.854±.038 [‡]	.903±.015
TA-DualCV	.856±.024**	.852±.028**	.864±.034**	.858±.034**	.905±.018**	.863±.032**	.856±.026**	.880±.045[‡]	.864±.049**	.909±.011**

- For each block, the best model is in **bold**; The best and the second-best models across ALL are labeled with ** and ‡, respectively.

shock prediction is discarded. Across Table V, TA-DualCV is the best approach across all approaches and MI settings. When combined with MI, the performance of each imputation approach is improved compared to the original corresponding approach without MI. This suggests that the missing pattern is indeed an important characteristic of EHRs for prediction. For example, when a patient is in a severe condition, events are likely to be recorded more frequently than when a patient is in a relatively “healthier” condition. Also, TA-DualCV generally benefits less from MI compared to other methods. For example, on MIMIC-III, the AUC of TA-DualCV only increases by 0.4% after combined with MI, while the AUC of TAME increases by 45%. The reason might be TA-DualCV has already captured various dependencies within data, which could contain missing pattern information revealed by MI. Moreover, imputation methods benefit more from MI on CCHS and Mayo, but not much on MIMIC-III probably because the former two contain much more native missing values compared to MIMIC-III.

VI. DISCUSSIONS & CONCLUSION

In this work, we present *TA-DualCV*, a non-Neural Network-based imputation framework towards both unsupervised and supervised learning tasks in EHRs. TA-DualCV integrates both cross-visit and within-visit dependencies, by exploiting dependencies among features, time-interval, and time-steps. The robustness of TA-DualCV is evaluated on two tasks: *unsupervised imputation task* and *supervised 24-hour septic shock early prediction task* using EHRs from three different medical systems. In both tasks, TA-DualCV is compared to state-of-the-art baselines. Among the different

baselines, 3D-MICE consistently outperformed other baselines on the unsupervised learning tasks, while on the supervised task of early septic shock prediction, there was no clear winner. TA-DualCV, on the other hand, shows its robustness to handle high missing rates of EHRs across medical systems by achieving the best performance on both imputation and septic shock early prediction tasks on all evaluation metrics.

It is also important to note that we cannot produce results in our experiments with some popular imputation approaches, which are data-driven to a specific type of missingness or data characteristic (e.g., ICU lab analytes data with a small native missing rate). Mayo, for example, can have events ranging from 1 to 5,000 across different visits. Consequently, approaches including Mix-MI [40] and GP-VAE [24] that require tensor-shape inputs cannot be utilized. For approaches relying on sliding windows, including DETROIT [11], a lower bound on the number of events is needed to compute a sliding window, and a specific density of neighboring measurements must be observed to predict missing values. They are not able to produce certain results in our experiments. Additionally, our experiments on CCHS and Mayo have fewer training data but higher average native missing rates (up to 73%) than those used the existing literature. In practice, as discussed in [31], deep learning relies on high-quality representations of the output of substantive data, whereas their imputation performance dwindles with limited training data.

ACKNOWLEDGEMENTS

This research was supported by the NSF Grants: Generalizing Data-Driven Technologies to Improve Individualized STEM Instruction by Intelligent Tutors (2013502), Integrated Data-driven Technologies for Individualized Instruction in

STEM Learning Environments (1726550), CAREER: Improving Adaptive Decision Making in Interactive Learning Environments (1651909), and S.E.P.S.I.S.: Sepsis Early Prediction Support Implementation System (1522107).

We would also like to thank the anonymous reviewers, Xi Yang (North Carolina State University), and Qitong Gao (Duke University) for insightful comments that leads to improved paper presentations.

REFERENCES

- [1] V. Liu, G. J. Escobar, J. D. Greene, J. Soule, A. Whippy, D. C. Angus, and T. J. Iwashyna, "Hospital deaths in patients with sepsis from 2 independent cohorts," *Jama*, vol. 312, no. 1, pp. 90–92, 2014.
- [2] R. C. Bone, R. A. Balk, F. B. Cerra, R. P. Dellinger, A. M. Fein, W. A. Knaus, R. M. Schein, and W. J. Sibbald, "Definitions for sepsis and organ failure and guidelines for the use of innovative therapies in sepsis," *Chest*, vol. 101, no. 6, pp. 1644–1655, 1992.
- [3] G. Martin, D. Mannino *et al.*, "The epidemiology of sepsis in the united states from 1979 through 2000," *New England Journal of Medicine*, 2003.
- [4] R. Dellinger, M. Levy *et al.*, "Surviving sepsis campaign: international guidelines for management of severe sepsis and septic shock: 2008," *Intensive care medicine*, 2008.
- [5] A. Kumar, D. Roberts, K. Wood, B. Light, J. Parrillo, S. Sharma, R. Suppes, D. Feinstein, S. Zanotti, L. Taiberg, D. Gurka, A. Kumar, and M. Cheang, "Duration of hypotension before initiation of effective antimicrobial therapy is the critical determinant of survival in human septic shock," *CCM*, vol. 34, no. 6, June 2006.
- [6] V. Coba, M. Whitmill, R. Mooney, H. M. Horst, M.-M. Brandt, B. Digiovine, M. Mlynarek, B. McLellan, G. Boleski, J. Yang *et al.*, "Resuscitation bundle compliance in severe sepsis and septic shock: improves survival. is better late than never," *Journal of intensive care medicine*, vol. 26, no. 5, pp. 304–313, 2011.
- [7] K. E. Henry, D. N. Hager, P. J. Pronovost, and S. Saria, "A targeted real-time early warning score (trewscore) for septic shock," *Science translational medicine*, vol. 7, no. 299, pp. 299ra122–299ra122, 2015.
- [8] J.-L. Vincent, R. Moreno, J. Takala, S. Willatts, A. De Mendonça, H. Bruining, C. Reinhart, P. Suter, and L. Thijs, "The sofa (sepsis-related organ failure assessment) score to describe organ dysfunction/failure," *Intensive care medicine*, vol. 22, no. 7, pp. 707–710, 1996.
- [9] N. I. Shapiro, R. E. Wolfe, R. B. Moore, E. Smith, E. Burdick, and D. W. Bates, "Mortality in emergency department sepsis (meds) score: a prospectively derived and validated clinical prediction rule," *Critical care medicine*, vol. 31, no. 3, pp. 670–675, 2003.
- [10] Y. Luo, P. Szolovits, A. S. Dighe, and J. M. Baron, "3d-mice: integration of cross-sectional and longitudinal imputation for multi-analyte longitudinal clinical data," *Journal of the American Medical Informatics Association*, vol. 25, no. 6, pp. 645–653, 2018.
- [11] C. Yan, C. Gao, X. Zhang, Y. Chen, and B. Malin, "Deep imputation of temporal data," in *2019 IEEE International Conference on Healthcare Informatics (ICHI)*. IEEE, 2019, pp. 1–3.
- [12] C. Yin, R. Liu, D. Zhang, and P. Zhang, "Identifying sepsis subphenotypes via time-aware multi-modal auto-encoder," in *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, 2020, pp. 862–872.
- [13] D. Rubin and R. Little, *Statistical Analysis with Missing Data*. Hoboken, NJ: John Wiley & Sons, 1987.
- [14] J. C. Ho, C. H. Lee, and J. Ghosh, "Septic shock prediction for patients with missing data," *ACM transactions on management information systems (TMIS)*, vol. 5, no. 1, pp. 1–15, 2014.
- [15] J. Fagerström, M. Bång, D. Wilhelms, and M. S. Chew, "Lisep lstm: a machine learning algorithm for early detection of septic shock," *Scientific reports*, vol. 9, no. 1, pp. 1–8, 2019.
- [16] Z. Lipton, D. Kale, and R. Wetzal, "Directly modeling missing data in sequences with RNNs: Improved classification of clinical time series," *JMLR*, vol. 56, 2016.
- [17] Q. Gao, D. Wang, J. D. Amason, S. Yuan, C. Tao, R. Henao, M. Hadzi-hmetovic, L. Carin, and M. Pajic, "Gradient importance learning for incomplete observations," *International Conference on Learning and Representations (ICLR)*, 2022.
- [18] F. Khoshnevisan, J. Ivy, M. Capan, R. Arnold, J. Huddleston, and M. Chi, "Recent temporal pattern mining for septic shock early prediction," in *2018 IEEE International Conference on Healthcare Informatics (ICHI)*. IEEE, 2018, pp. 229–240.
- [19] Y.-J. Kim and M. Chi, "Temporal belief memory: Imputing missing data during rnn training," in *In Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI-2018)*, 2018.
- [20] 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, no. 1, pp. 1–9, 2016.
- [21] F. Khoshnevisan and M. Chi, "An adversarial domain separation framework for septic shock early prediction across ehr systems," in *2020 IEEE International Conference on Big Data (Big Data)*. IEEE, 2020, pp. 64–73.
- [22] X. Yang, Y. Zhang, and M. Chi, "Time-aware subgroup matrix decomposition: Imputing missing data using forecasting events," in *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 2018, pp. 1524–1533.
- [23] S. v. Buuren and K. Groothuis-Oudshoorn, "mice: Multivariate imputation by chained equations in r," *Journal of statistical software*, pp. 1–68, 2010.
- [24] V. Fortuin, D. Baranchuk, G. Rätsch, and S. Mandt, "Gp-vae: Deep probabilistic time series imputation," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 1651–1661.
- [25] I. Silva, G. Moody, D. J. Scott, L. A. Celi, and R. G. Mark, "Predicting in-hospital mortality of icu patients: The physionet/computing in cardiology challenge 2012," in *2012 Computing in Cardiology*. IEEE, 2012, pp. 245–248.
- [26] W. Cao, D. Wang, J. Li, H. Zhou, L. Li, and Y. Li, "Brits: Bidirectional recurrent imputation for time series," *arXiv preprint arXiv:1805.10572*, 2018.
- [27] K. Yin, L. Feng, and W. K. Cheung, "Context-aware time series imputation for multi-analyte clinical data," *Journal of Healthcare Informatics Research*, vol. 4, no. 4, pp. 411–426, 2020.
- [28] P. Ranjan, R. Haynes, and R. Karsten, "A computationally stable approach to gaussian process interpolation of deterministic computer simulation data," *Technometrics*, vol. 53, no. 4, pp. 366–378, 2011.
- [29] B. MacDonald, P. Ranjan, H. Chipman *et al.*, "Gpfit: An r package for fitting a gaussian process model to deterministic simulator outputs," *Journal of Statistical Software*, vol. 64, no. i12, 2015.
- [30] X. Miao, Y. Wu, J. Wang, Y. Gao, X. Mao, and J. Yin, "Generative semi-supervised learning for multivariate time series imputation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 10, 2021, pp. 8983–8991.
- [31] X. Xu, X. Liu, Y. Kang, X. Xu, J. Wang, Y. Sun, Q. Chen, X. Jia, X. Ma, X. Meng *et al.*, "A multi-directional approach for missing value estimation in multivariate time series clinical data," *Journal of Healthcare Informatics Research*, vol. 4, pp. 365–382, 2020.
- [32] X. Zhang, C. Yan, C. Gao, B. A. Malin, and Y. Chen, "Predicting missing values in medical data via xgboost regression," *Journal of Healthcare Informatics Research*, vol. 4, no. 4, pp. 383–394, 2020.
- [33] J. Ho, C. Lee, and J. Ghosh, "Septic shock prediction for patients with missing data," *Management Information Systems*, vol. 5, no. 1, April 2014.
- [34] Y. Zhang, C. Lin, M. Chi, J. Ivy, M. Capan, and J. M. Huddleston, "Lstm for septic shock: Adding unreliable labels to reliable predictions," in *2017 IEEE International Conference on Big Data (Big Data)*. IEEE, 2017, pp. 1233–1242.
- [35] M. Singer, C. S. Deutschman, C. W. Seymour, M. Shankar-Hari, D. Annane, M. Bauer, R. Bellomo, G. R. Bernard, J.-D. Chiche, C. M. Coopersmith *et al.*, "The third international consensus definitions for sepsis and septic shock (sepsis-3)," *Jama*, vol. 315, no. 8, pp. 801–810, 2016.
- [36] R. G. Pontius, O. Thonteh, and H. Chen, "Components of information for multiple resolution comparison between maps that share a real variable," *Environmental and Ecological Statistics*, vol. 15, no. 2, pp. 111–142, 2008.
- [37] H. Sohn, K. Park, and M. Chi, "Mulan: Multilevel language-based representation learning for disease progression modeling," in *IEEE International Conference on Big Data, Big Data 2020, Atlanta, GA, USA, December 10-13, 2020*. IEEE, 2020, pp. 1246–1255. [Online]. Available: <https://doi.org/10.1109/BigData50022.2020.9377829>

- [38] Y. Zhang, X. Yang, J. Ivy, and M. Chi, "Time-aware adversarial networks for adapting disease progression modeling," in *2019 IEEE International Conference on Healthcare Informatics (ICHI)*. IEEE, 2019, pp. 1–11.
- [39] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [40] Y. Xue, D. Klabjan, and Y. Luo, "Mixture-based multiple imputation model for clinical data with a temporal dimension," in *2019 IEEE International Conference on Big Data (Big Data)*. IEEE, 2019, pp. 245–252.