

# GENERATING SYNTHETIC MIXED-TYPE LONGITUDINAL ELECTRONIC HEALTH RECORDS FOR ARTIFICIAL INTELLIGENT APPLICATIONS

Jin Li<sup>†‡</sup>, Benjamin J. Cairns<sup>‡</sup>, Jingsong Li<sup>†\*</sup>, Tingting Zhu<sup>‡\*</sup>

<sup>†</sup>Zhejiang University, <sup>‡</sup>University of Oxford

{jin.li, tingting.zhu}@eng.ox.ac.uk, ben.cairns@ndph.ox.ac.uk, ljs@zju.edu.cn

## ABSTRACT

The recent availability of electronic health records (EHRs) have provided enormous opportunities to develop artificial intelligence (AI) algorithms. However, patient privacy has become a major concern that limits data sharing across hospital settings and subsequently hinders the advances in AI. *Synthetic data*, which benefits from the development and proliferation of generative models, has served as a promising substitute for real patient EHR data. However, the current generative models are limited as they only generate *single type* of clinical data for a synthetic patient, i.e., either continuous-valued or discrete-valued. To mimic the nature of clinical decision-making which encompasses various data types/sources, in this study, we propose a generative adversarial network (GAN) entitled EHR-M-GAN which simultaneously synthesizes *mixed-type* timeseries EHR data. EHR-M-GAN is capable of capturing the multidimensional, heterogeneous, and correlated temporal dynamics in patient trajectories. We have validated EHR-M-GAN on three publicly-available intensive care unit databases with records from a total of 141,488 unique patients, and performed privacy risk evaluation of the proposed model. EHR-M-GAN has demonstrated its superiority over state-of-the-art benchmarks for synthesizing clinical timeseries with high fidelity, while addressing the limitations regarding data types and dimensionality in the current generative models. Notably, prediction models for outcomes of intensive care performed significantly better when training data was augmented with the addition of EHR-M-GAN-generated timeseries. EHR-M-GAN may have use in developing AI algorithms in resource-limited settings, lowering the barrier for data acquisition while preserving patient privacy.

## 1 INTRODUCTION

The past decade has witnessed ground-breaking advancements been made in computational health, owing to the explosion of medical data, such as electronic health records (EHRs) Artzi et al. (2020); Raket et al. (2020); Menger et al. (2019). The secondary uses of EHRs give rise to research in a wide range of varieties, especially machine learning (ML)-based digital health solutions for improving the delivery of care Wilkinson et al. (2020); Watson et al. (2019); Futoma et al. (2020); Esteva et al. (2021); Rajkomar et al. (2019). However, in practice, the benefits of data-driven research are limited to healthcare organizations (HCOs) who possess the data Wirth et al. (2021); Dinov (2016). Due to concerns about patient privacy, HCO stakeholders are reluctant to share patient data Miotto et al. (2018); Kim et al. (2021); Simon et al. (2019). Access to clinical data is often restricted, or can be prohibitively expensive to obtain, meaning that ML in biomedical research lags behind other areas in AI.

To accelerate the progress of developing AI methods in medicine, one promising alternative is for the data holder to create *synthetic* yet realistic data Jordon et al. (2018); Frid-Adar et al. (2018). By avoiding “one-to-one” mapping to the genuine data compared with data anonymization, synthetic data offers a solution to circumvent the issue of privacy, while the correlations in the original data

\*Corresponding authors.

---

distributions are preserved for downstream AI applications. There have been successes in the literature using synthetic data to improve AI models where otherwise not possible due to limited availability of resources Jordon et al. (2020); Chen et al. (2021); Tucker et al. (2020); El Emam et al. (2021). For example, large-scale data sharing programs have been demanded for advancing studies related to COVID-19, such as in National COVID Cohort Collaborative (N3C) *N3C Synthetic Data Workstream* (n.d.), and Clinical Practice Research Datalink (CPRD) database in the UK *Synthetic data at CPRD, howpublished = <https://www.cprd.com/content/synthetic-data>* (n.d.).

Recent advances in generative adversarial networks (GANs) Goodfellow et al. (2014) and their variants offer efficacious means to generate EHRs for a wide range of clinical applications Kearney et al. (2020); Yang et al. (2018); Marouf et al. (2020). Over the past years, EHR synthesizers have evolved from generating static patient information to producing longitudinal EHR timeseries Esteban et al. (2017); Lee et al. (2020); Zhang et al. (2021). As longitudinal EHRs contain patient trajectories for describing the underlying health condition, synthesizing such EHR timeseries, therefore, enables new clinical applications related to the status of disease progression Zhang et al. (2022), such as dynamic forecasting of risks, predicting the onset of diseases, and survival analysis based on the time-to-event data. However, existing studies focus on synthesizing the longitudinal EHRs of a single data type Esteban et al. (2017); Yoon et al. (2019); Lee et al. (2020), whereas the clinical decision-making in real practice includes a variety of information sources in the form of *mixed-type* timeseries. For example, patient physiological signals and laboratory test results are collected in the EHR as *continuous-valued* timeseries, while the medication and diagnostic information are recorded as *discretized-valued* data as binary indicators or categorical ICD codes. Information provided in these mixed-type longitudinal EHRs offer opportunities for more precise and complex clinical analysis. Furthermore, the predictive power and robustness of the ML models can be boosted by utilizing longitudinal EHR timeseries with various types/sources.

Existing GANs are limited in simulating mixed-type EHRs due to two reasons. Firstly, it is intrinsically difficult to model the underlying joint distribution of mixed data type timeseries using a single unified framework. Since GANs require the network architectures of the generator and discriminator to be fully differentiable Hjelm et al. (2017), its success is typically limited to generating real-valued, continuous data while facing obstacles for directly generating sequences of discrete tokens, such as ICD codes, that also commonly appear in EHRs. Previous methods Yu et al. (2017); Choi et al. (2017) circumvent this problem by learning representations from the original data which further enables backpropagation in discrete settings, but there is still a lack of a generative approach for joint modelling of the mixed-type timeseries with heterogeneous nature. Second, although mixed-type clinical timeseries differ in syntax and distributions, they are highly correlated and inform one another of the underlying health of an individual Yu et al. (2019); Ghassemi et al. (2017); Wang et al. (2019). It is therefore important to capture the temporal correlations between them when generating the synthetic EHR data. For example, the medications (documented in the form of discrete data) prescribed to patients are based on measurements of patients’ physiological status (presented as continuous-valued signals). Concurrently, the efficacy of the medical treatments, affect the patient’s physiological condition directly. It is therefore critical to accurately capture the temporal correlation between the mixed-type patient trajectories simultaneously to improve clinical decision support.

To address the aforementioned limitations, for the first time, we propose a GAN framework for simultaneously synthesizing mixed-type longitudinal EHR data (denoted as EHR-M-GAN thereafter). Specifically, we focus on generating timeseries in the critical care setting, where the intensive care units (ICU) patients are continuously and closely monitored (see Fig. 1a). Patient trajectories with high-dimensionality and heterogeneous data types (both continuous-valued and discrete-valued timeseries) are generated while the underlying temporal dependencies are captured. The main contributions of our work are as follows:

1. A novel GAN model entitled EHR-M-GAN is proposed for simultaneously generating mixed-type multivariate EHR timeseries with high fidelity, and overcoming the challenges when extending GANs into the mixed-type data settings (see Fig. 1b). First, to jointly model the underlying distributions of the heterogeneous features, EHR-M-GAN first maps data from different observational spaces into a reversible, lower-dimensional, shared latent space through a *dual variational autoencoder* (dual-VAE). Then, to capture the correlated temporal dynamics of the mixed-type data, a sequentially coupled generator that is built upon a *coupled recurrent network* (CRN) is employed. In addition, a conditional version of our model — EHR-M-GAN<sub>cond</sub> — is also

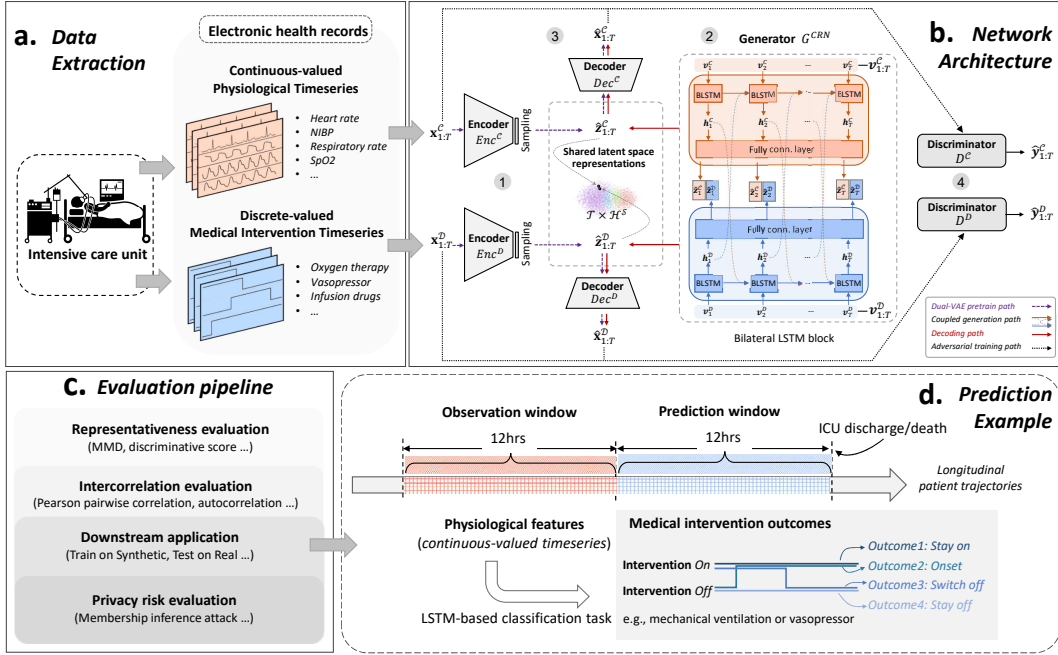


Figure 1: **Overall schematics.** **a. Data extraction.** Electronic health records (EHRs) data of mixed-type are routinely collected for patients in intensive care units (ICUs). **b. Network architecture.** EHR-M-GAN contains two key components — *Dual-VAE* and *Coupled Recurrent Network (CRN)*. **Step 1:** *Dual-VAE* is first pretrained for mapping heterogeneous data ( $x_t^c, x_t^d$ ) into shared latent representations ( $z_t^c, z_t^d$ ). Multiple objective loss constraints are used to bridge the domain/distribution gap. The training process for Step 1 is indicated in the *Dual-VAE pretrain path* (dashed purple line). **Step 2:** Then, a CRN is established as the generator based on the parallel bilateral LSTM block, which takes the random noise vectors ( $v_t^c, v_t^d$ ) as inputs (see the *Coupled generation path*). **Step 3:** The synthetic latent representations ( $\hat{z}_t^c, \hat{z}_t^d$ ) provided by CRN are then decoded into synthetic samples ( $\hat{x}_t^c, \hat{x}_t^d$ ) using the pretrained decoder in *Dual-VAE*, which is indicated in the *Decoding path* (solid red line). **Step 4:** Finally, the adversarial loss is derived from the discriminators and backpropagated to update the network, which is indicated in the *Adversarial training path* (dotted black line). **c. Evaluation pipeline.** The pipeline includes metrics for evaluating the synthetic data. **d. Prediction example.** Data within 24-hours prior to the patient’s endpoints in the ICU (discharge or mortality) is extracted. Both the observation window and prediction window are fixed as 12 hours. The classification task is to use patients’ continuous-valued physiological measurements within the observation window as input, to predict the forthcoming discrete-valued medical intervention status in the prediction window. The four outcomes of the intervention status can be categorized as follows: **Stay On:** The intervention begins with *on* and *stays on* within the prediction window; **Onset:** The intervention begins with *off* and is *turned on* within the prediction window; **Switch off:** The intervention begins with *on* and is *stopped* within the prediction window; **Stay Off:** The intervention begins with *off* and *stays off* within the prediction window.

implemented, which is capable of synthesizing condition-specific EHR patient data, such as those result in *ICU mortality* or *hospital readmission*. The code of our proposed work is publicly available on GitHub<sup>1</sup>.

- Evaluations are performed based on three publicly available ICU datasets: MIMIC-III Johnson et al. (2016), eICU Pollard et al. (2018) and HiRID Yèche et al. (2021) from a total of 141,488 patients. Standardized preprocessing pipelines are applied for the three ICU datasets to provide generalizable machine learning benchmarks. The code for the end-to-end preprocessing pipelines is also available on GitHub<sup>2</sup>.

<sup>1</sup><https://github.com/jli0117/ehrMGAN>

<sup>2</sup>[https://github.com/jli0117/preprocessing\\_physionet](https://github.com/jli0117/preprocessing_physionet)

3. Our EHR-M-GAN outperforms the state-of-the-art benchmarks on a diverse spectrum of evaluation metrics. When compared to real EHR data, both qualitative and quantitative metrics are used to assess the representativeness of the mixed-type data and their inter-dependencies. We further demonstrate the advantages offered by EHR-M-GAN in augmenting clinical timeseries for downstream tasks under various clinical scenarios.
4. In the evaluation of privacy risks, we perform an empirical analysis on EHR-M-GAN based on membership inference attack Shokri et al. (2017). We then further evaluate the performance of EHR-M-GAN under the framework of differential privacy for its application in downstream task Dwork (2006).

## 2 METHODS

In this section, we first formulate the problem based on the *mixed-type* temporal EHR data and its corresponding mathematical notation. Then, we discuss the challenges of synthesizing mixed-type EHR timeseries and the intuition behind the proposed model. Finally, we introduce the proposed EHR-M-GAN model in detail.

### 2.1 PROBLEM FORMULATION

The longitudinal patient EHR dataset is denoted as  $\mathcal{D} = \{(\mathbf{x}_{i,1:T_i})\}_{i=1}^N$ , with each record (e.g., individual patient) being indexed by  $i \in \{1, 2, \dots, N\}$ . Here we consider the  $i$ -th instance tuple  $\mathbf{x}_{i,1:T_i} = \{\mathbf{x}_{i,1:T_i}^C, \mathbf{x}_{i,1:T_i}^D\}$  consists of two components (i.e., two types of data). Let  $\mathbf{x}_{i,1:T_i}^C \in \mathbb{R}^{|J|}$  denote the  $|J|$ -dimensional continuous-valued timeseries, such as physiological signals from real-time bedside monitors. And  $\mathbf{x}_{i,1:T_i}^D \in \mathbb{Z}^{|K|}$  denotes the  $|K|$ -dimensional discrete-valued timeseries, such as life-support interventions with the categorical value indicate its status (presence or absence).

### 2.2 CHALLENGES IN MIXED-TYPE TIMESERIES GENERATION

There are two main challenges when synthesizing mixed-type EHR timeseries. First, GANs have serious limitations on the type of data they can model Hjelm et al. (2017). Specifically, as GANs require generators and discriminators to be both fully differentiable, generating discrete-valued timeseries using traditional GANs architectures would raise problems during backpropagation as no direct gradient can be provided Choi et al. (2017); Yu et al. (2017). Therefore, it is intrinsically difficult to model the underlying joint distribution of mixed data type timeseries using a single unified framework. Second, as the mixed-type timeseries are correlated (such as correlations between ICU patients’ physiological signals and treatment status in the critical care setting), it is therefore important to model the interdependencies among heterogeneous types of timeseries.

### 2.3 INTUITION BEHIND EHR-M-GAN

First, to jointly model the distribution of continuous-valued and discrete-valued timeseries using GANs, we build the generative model based on the latent space encoded by VAE networks. Instead of directly synthesizing discrete-valued timeseries that deactivate the backpropagation in GANs, the generator first synthesizes latent representations that allow the direct gradient in the network, therefore satisfying the prerequisite for GANs architecture to be fully differentiable. The synthetic latent representations for both types of data can be further decoded into raw timeseries using the decoders in VAEs.

Even though the aforementioned network architectures enable the joint modelling of mixed-type data distribution, it still lacks the capability of capturing the inter-dependencies in heterogeneous data. In order to address the second issue, we devised *dual-VAE* module for pretraining step and *sequentially coupled generator* module for generation step. The *dual-VAE* incorporates multiple loss constraints, which were previously adopted in domains such as self-supervised learning (SSL), timeseries representation learning, and domain adaptation (DA), to extract useful hierarchical representations from heterogeneous but correlated data types. The *sequentially coupled generator* module replaces the traditional LSTM cell with the novel bilateral LSTM (BLSTM) cell we propose, where the “communication” of the two types of information are introduced into the networks. Therefore, the temporal dynamics between the mixed-type data can be preserved during the iteration.

## 2.4 PROPOSED MODEL

As illustrated above, EHR-M-GAN can be factorized into two key components (see Fig. 1b): (1) a *dual-VAE* framework for learning the shared latent space representations; (2) an RNN-based *sequentially coupled generator* and its corresponding sequence discriminators.

As shown in Fig. 1b, during the *pretrain* stage, both continuous-valued and discrete-valued temporal trajectories are first jointly mapped into a shared latent space using the *dual-VAE* component (Step 1). Then, the *sequentially coupled generator* in EHR-M-GAN produces the synthetic latent representations (Step 2), which further can be recovered into features in the observational space by the pretrained decoders in the *dual-VAE* (Step 3). Finally, the adversarial loss is provided based on discriminative results and backpropagated to update the network (Step 4). The following sections discuss them in turn.

### 2.4.1 DUAL-VAE PRETRAINING FOR SHARED LATENT SPACE REPRESENTATIONS

One premise of successfully training EHR-M-GAN to generate reversible latent codes is to meet the assumption that for the *same* patient indexed with  $i$ , both  $\mathbf{x}_{i,1:T_i}^C$  and  $\mathbf{x}_{i,1:T_i}^D$  can be encoded into the *same* latent space  $\mathcal{H}^S \subset \mathbb{R}^{|S|}$ , where  $|S|$  denotes its spatial dimension. For the sake of simplicity, the subscripts  $i$  are omitted throughout most of the paper. To achieve this, we propose to use a *dual-VAE* framework, which exploits two VAE networks to encode both continuous and discrete multivariate timeseries into dense representations within  $\mathcal{H}^S$  based on multiple constraints.

Fig. S2 (see Section S.1.C in Supplementary materials) diagrams the details of the proposed *dual-VAE* framework for learning the shared latent representations. We start with training two encoders, i.e.,  $Enc^C: \phi_{\mathcal{T} \times \mathcal{X}^C} \rightarrow \phi_{\mathcal{T} \times \mathcal{H}^S}$  and  $Enc^D: \phi_{\mathcal{T} \times \mathcal{X}^D} \rightarrow \phi_{\mathcal{T} \times \mathcal{H}^S}$ , with the embedding functions:

$$\mathbf{z}_{1:T}^C = Enc^C(\mathbf{x}_{1:T}^C) \quad \mathbf{z}_{1:T}^D = Enc^D(\mathbf{x}_{1:T}^D) \quad (1)$$

After passing data from  $\mathcal{X}^C$  and  $\mathcal{X}^D$  through two encoders, a pair of embedding vectors  $(\mathbf{z}_{1:T}^C, \mathbf{z}_{1:T}^D)$  in the shared latent space  $\mathcal{H}^S$  can be obtained. Then the decoders for both domains  $Dec^C: \psi_{\mathcal{T} \times \mathcal{H}^S} \rightarrow \psi_{\mathcal{T} \times \mathcal{X}^C}$  and  $Dec^D: \psi_{\mathcal{T} \times \mathcal{H}^S} \rightarrow \psi_{\mathcal{T} \times \mathcal{X}^D}$  further reconstruct features based on the latent embeddings using mapping functions that operate in the opposite direction:

$$\tilde{\mathbf{x}}_{1:T}^C = Dec^C(\mathbf{z}_{1:T}^C) \quad \tilde{\mathbf{x}}_{1:T}^D = Dec^D(\mathbf{z}_{1:T}^D) \quad (2)$$

Also, to incentivize dual-VAE to better bridge the gap between domains of mixed-type timeseries, we enforce a weight-sharing constraint Liu et al. (2017); Liu and Tuzel (2016) within specific layers of both the encoders pairs and the decoders pairs (See Section S.1.B for details).

In the following subsections, we define multiple loss constraints for the optimization of *dual-VAE*, including *ELBO loss*, *matching loss*, *contrastive loss*, as well as *semantic loss* for EHR-M-GAN<sub>cond</sub>. Among these losses, *ELBO loss* ensures that the mixed-type timeseries can be successfully reconstructed after being encoded into latent representations. The *matching loss* ensures that heterogeneous types of features from a single patient share contexts during representation learning (instance-wise). The goal of *contrastive loss* is to ensure that patients with similar trajectories stay close to each other in the latent space (population-wise). And *semantic loss* used in EHR-M-GAN<sub>cond</sub> encourages patients with the same conditional labels (e.g., outcomes) to share similar latent representations. Intuitions and descriptions behind the objectives are discussed in turn.

**Evidence Lower Bound (ELBO).** We first incorporate the standard VAE loss, with the optimization objective as the evidence lower bound (ELBO). VAE holds the assumption of spherical Gaussian prior for the distribution of latent embeddings, where features can then be reconstructed by sampling from that space. The re-parameterization tricks enable differentiable stochastic sampling and network optimization. For encoder and decoder in the dual-VAE for domain  $d \in \{\mathcal{C}, \mathcal{D}\}$ , the objective function is defined as:

$$\mathcal{L}_d^{\text{ELBO}} = -\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\psi(\mathbf{x}|\mathbf{z})] + \beta_{\text{KL}} D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \| p_\psi(\mathbf{z})) \quad (3)$$

where  $\mathbf{z} \sim Enc(\mathbf{x}) \triangleq q_\phi(\mathbf{z}|\mathbf{x})$ ,  $\tilde{\mathbf{x}} \sim Dec(\mathbf{z}) \triangleq p_\psi(\mathbf{x}|\mathbf{z})$ , and  $D_{\text{KL}}$  is the Kullback-Leibler divergence. The first term in Eq. (3) is the expected log-likelihood term that penalizes the deviations

in reconstructing the inputs, while the second term of KL-divergence is the regularization imposed over the latent distribution from its Gaussian prior (normally chosen to be  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ ).  $\beta_{\text{KL}}$  is the hyperparameter for balancing the weights between two terms.

**Matching loss.** Typically, representations derived from the *same* patient are assumed to capture the shared context. Therefore, embedding vectors ( $\mathbf{z}_{i,1:T_i}^{\mathcal{C}}, \mathbf{z}_{i,1:T_i}^{\mathcal{D}}$ ) projected from the *same* patient  $i$ , are supposed to be positioned closely in the shared latent space (See Fig. S2 in Supplementary materials). Therefore, in this study, we borrow the concept of matching loss from domain alignment in DA, which enables efficient representation learning crossing domains/modalities Wan et al. (2020). In this study, the matching loss ensures that low-dimensional latent space can be shared between heterogeneous features. Hence, the pairwise matching loss is incorporated to motivate the encoders to minimize the distance within the corresponding representation pairs. In the low-dimensional Euclidean space, we optimize the network by using the following objective:

$$\mathcal{L}^{\text{Match}} = \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} \left[ \sum_{t \in \mathcal{T}} \|\mathbf{z}_t^{\mathcal{C}} - \mathbf{z}_t^{\mathcal{D}}\|^2 \right] \quad (4)$$

The pairwise matching loss achieve its optimal when the distance proxy  $\mathcal{L}^{\text{Match}}$  becomes zero.

**Contrastive loss.** On the flip side, pairwise reconstruction error (i.e., intra-correlations within one instance) measured by *matching loss* neglects the commonalities present across patients (inter-correlations of data) Kiyasseh et al. (2021). In order to guarantee sufficient bound for representation learning, we incorporate *contrastive loss* as another distance metric to capture the inter-correlations among the population.

Contrastive learning is a concept that has recently been popularized in self-supervised learning (SSL) Liu et al. (2021), which aims to capture intrinsic patterns from input data without human annotations. In this study, we instantiate the contrastive loss via *NT-Xent*, which is proposed by Chen et al. in their work SimCLR Chen, Kornblith, Norouzi and Hinton (2020). The core of contrastive learning is to encourage networks to attract positive pairs closer and repulse negative pairs apart in the latent space. In this study, we adapt the corresponding auxiliary tasks for calculating contrastive loss to the scenario of learning representations from mixed-type timeseries. The objective of the task is to determine whether a set of representations transformed from the observational space belong to the *same* patient. And this leads to the corresponding positive pairs (true) and negative pairs (false).

For patient data of  $N$  records, we can obtain  $N$  pairs of latent representations from the encoders in *dual-VAE*. For patient indexed with  $i$ ,  $\mathbf{h}_i^{\mathcal{C}}$  and  $\mathbf{h}_i^{\mathcal{D}}$  denotes the embeddings derived from the continuous-valued and discrete-valued observational space, respectively. Due to the symmetric architecture of *dual-VAE*, here we use  $d$  and  $d'$  to represent one of each different domain, i.e.,  $d, d' \in \{\mathcal{C}, \mathcal{D}\}$  and  $d \neq d'$ . Therefore, the positive pairs for patient  $i$  can be referred as  $(i^d, i^{d'})$ , while the other  $2(N - 1)$  samples are regarded as negative pairs. Then the contrastive loss for a positive pair  $(i^d, i^{d'})$  is defined as:

$$\mathcal{L}_{i^d, i^{d'}}^{\text{Contra}} = -\log \frac{\exp(\text{sim}(\mathbf{h}_{i^d}, \mathbf{h}_{i^{d'}})/\tau)}{\sum_{i^{dd'}=1}^{2N} \mathbb{1}_{[i^{dd'} \neq i^d]} \exp(\text{sim}(\mathbf{h}_{i^d}, \mathbf{h}_{i^{dd'}})/\tau)} \quad (5)$$

where  $\text{sim}(u, v) = u^T v / \|u\| \|v\|$  denotes the cosine similarity between two vectors.  $\tau > 0$  denotes a temperature hyperparameter.  $\mathbb{1}_{[n \neq m]} \in \{0, 1\}$  is an indicator evaluating to 1 iff  $n \neq m$ . And  $i^{dd'} \in \{1, 2, \dots, 2N\}$  represents the index of latent embeddings from *both* data types. The final contrastive loss is computed across the total number of  $|i^d - i^{d'}| = N$  positive pairs for both  $(i^d, i^{d'})$  and  $(i^{d'}, i^d)$ , and is defined as:

$$\mathcal{L}^{\text{Contra}} = \frac{1}{2N} \sum_{i^d=1}^N \sum_{i^{d'}=1}^N [\mathcal{L}_{i^d, i^{d'}}^{\text{Contra}} + \mathcal{L}_{i^{d'}, i^d}^{\text{Contra}}] \quad (6)$$

**Semantic loss.** In  $\text{EHR-M-GAN}_{\text{cond}}$ , semantic loss is imposed to better align patients with same labels (conditions) into the same latent space clusters. For example, if the label of *severe clinical deterioration* in the ICU is given for conditional data generation, the corresponding synthetic continuous-valued timeseries (e.g., severely deranged vital signs) is supposed to be strongly associated with the discrete-valued timeseries (e.g., intensive medical interventions) under the same label.

For both domains, additional linear classifiers are trained to classify the latent embeddings based on their corresponding conditions in the observational space. We implement logistic regression as the linear classifiers and calculate the cross entropy as the semantic losses for both domains. For  $d \in \{\mathcal{C}, \mathcal{D}\}$ , given the latent embedding vector  $\mathbf{z}^d$  and the conditional information vector  $\mathbf{y}$ :

$$\mathcal{L}_d^{\text{Class}} = \mathbb{E}_{\mathbf{z}^d \in \mathcal{H}^S} \text{CE} (f_{\text{linear}}^d(\mathbf{z}^d), \mathbf{y}) \quad (7)$$

where  $f_{\text{linear}}^d$  denotes the linear classifier for the corresponding domain. And  $\text{CE} = -\sum_j y_j \log(\hat{y}_j)$ , ( $j = 1, 2, \dots, |L|$ ) denotes the cross entropy loss, where  $\hat{y}_j$  is the output of the linear classifier, and  $y_j$  is the ground truth label for class  $j$  in condition vector  $\mathbf{y}$ .

In summary, to train the *dual*-VAE for learning the shared latent space representation, the total objective function for  $d \in \{\mathcal{C}, \mathcal{D}\}$  is:

$$\mathcal{L}_d = \beta_0 \mathcal{L}_d^{\text{ELBO}} + \beta_1 \mathcal{L}^{\text{Match}} + \beta_2 \mathcal{L}^{\text{Contra}} \quad (8)$$

Under the conditional learning scenario of EHR-M-GAN<sub>cond</sub>, the total loss becomes:

$$\mathcal{L}_d = \beta_0 \mathcal{L}_d^{\text{ELBO}} + \beta_1 \mathcal{L}^{\text{Match}} + \beta_2 \mathcal{L}^{\text{Contra}} + \beta_3 \mathcal{L}_d^{\text{Class}} \quad (9)$$

where  $\beta_0, \beta_1, \beta_2$ , and  $\beta_3$  are scalar loss weights used to balance the loss terms.

To validate the effectiveness of multiple losses and the weight-sharing constraint in the proposed *dual*-VAE network, we perform the corresponding ablation study using MIMIC-III dataset as an example. The results can be found in S.3.B in the Supplementary materials. As shown in Table S7, all the components in the proposed *dual*-VAE network contribute to the improvement of EHR-M-GAN's performance when generating mixed-type timeseries data.

#### SEQUENTIALLY COUPLED GENERATOR BASED ON CRN

We propose the *sequentially coupled generator* for generating latent representations for mixed-type timeseries, which is built based on the network architecture of *coupled recurrent network (CRN)*. Specifically, a CRN exploits bilateral long short-term memory (BLSTM) cells as its recurrent layer to preserve the temporal dependencies between the continuous and discrete-valued sequences. The novel network architecture of bilateral-LSTM we proposed can extract and transmit the correlations between the mixed-type timeseries, as opposed to vanilla-LSTM which has only one recursive connection. In the following section, we first discuss the structure of BLSTM in detail as its essential recurrent layer of CRN, and then build the *sequentially coupled generator* based on CRN.

**Bilateral long short-term memory.** As the traditional LSTM only considers temporal dynamics from single-type timeseries, therefore is incapable to extract and transmit temporal correlation from heterogeneous features. Therefore, we propose the novel bilateral-LSTM cell with two network connections to characterize the correlations between two types of data. Given  $d, d' \in \{\mathcal{C}, \mathcal{D}\}$ ,  $\mathbf{v}_t^d$  and  $\mathbf{h}_t^d$  denotes the input vector (i.e., the random noise during GANs' training) and hidden state vector for domain  $d$  at time step  $t$ , respectively. An additional set of weights for introducing hidden states representations  $\mathbf{h}_t^{d'}$  from domain  $d'$  is included when updating the input gate  $\mathbf{i}_t^d$ , forget gate  $\mathbf{f}_t^d$ , output gate  $\mathbf{o}_t^d$ , and cell memory  $\tilde{\mathbf{c}}_t^d$ . The state transition functions for BLSTM are:

$$\begin{aligned} \mathbf{i}_t^d &= \sigma \left( \mathbf{W}_{idv} \mathbf{v}_t^d + \mathbf{W}_{idh^{d'}} \mathbf{h}_{t-1}^{d'} + \mathbf{W}_{idh^d} \mathbf{h}_{t-1}^d + \mathbf{b}_{id} \right) \\ \mathbf{f}_t^d &= \sigma \left( \mathbf{W}_{fdv} \mathbf{v}_t^d + \mathbf{W}_{fdh^{d'}} \mathbf{h}_{t-1}^{d'} + \mathbf{W}_{fdh^d} \mathbf{h}_{t-1}^d + \mathbf{b}_{fd} \right) \\ \mathbf{o}_t^d &= \sigma \left( \mathbf{W}_{odv} \mathbf{v}_t^d + \mathbf{W}_{odh^{d'}} \mathbf{h}_{t-1}^{d'} + \mathbf{W}_{odh^d} \mathbf{h}_{t-1}^d + \mathbf{b}_{od} \right) \\ \tilde{\mathbf{c}}_t^d &= \tanh \left( \mathbf{W}_{cdv} \mathbf{v}_t^d + \mathbf{W}_{cdh^{d'}} \mathbf{h}_{t-1}^{d'} + \mathbf{W}_{cdh^d} \mathbf{h}_{t-1}^d + \mathbf{b}_{cd} \right) \\ \mathbf{c}_t^d &= \mathbf{f}_t^d \odot \mathbf{c}_{t-1}^d + \mathbf{i}_t^d \odot \tilde{\mathbf{c}}_t^d \\ \mathbf{h}_t^d &= \mathbf{o}_t^d \odot \tanh(\mathbf{c}_t^d) \end{aligned} \quad (10)$$

As indicated by Eq. 10, the proposed BLSTM network overcomes the limitation of vanilla-LSTM network on modelling the correlation between the mixed-type timeseries by establishing the supplemental recursive connection. The new connection facilitates the model to intrinsically

decide how much information it should pass through the gates from its counterpart. A diagram of the BLSTM cell in contrast to vanilla-LSTM cell can be found in the Supplementary materials (see Fig. S3).

**Coupled recurrent network.** The architecture of CRN consists of three layers: the *input layers*, the *recurrent layers*, and the *fully connected layers*. First, the random noise vectors  $\mathbf{v}_t^d$  and  $\mathbf{v}_t^{d'}$  for two domains, which are sampled from uniform distributions (i.e.,  $\mathbf{v}_t^d, \mathbf{v}_t^{d'} \in \mathcal{U}(0, 1)$ ), are separately fed into the *input layers*. Then the *recurrent layers*  $f_{\text{rec}}$ , which are built based on two streams of BLSTM, one for each data type, are used to recursively iterate hidden states from both branches. Finally, the *fully connected layers*  $f_{\text{conn}}^d$  and  $f_{\text{conn}}^{d'}$  produce the generated latent vectors  $\hat{\mathbf{z}}_t^d$  and  $\hat{\mathbf{z}}_t^{d'}$  for the decoding stage in *dual-VAE*. At time step  $t$ , CRN can be formulated as:

$$\begin{aligned} (\mathbf{h}_t^d, \mathbf{h}_t^{d'}) &= f_{\text{rec}}((\mathbf{v}_t^d, \mathbf{v}_t^{d'}), (\mathbf{h}_{t-1}^d, \mathbf{h}_{t-1}^{d'})) \\ \hat{\mathbf{z}}_t^d &= f_{\text{conn}}^d(\mathbf{h}_t^d) \\ \hat{\mathbf{z}}_t^{d'} &= f_{\text{conn}}^{d'}(\mathbf{h}_t^{d'}) \end{aligned} \quad (11)$$

In summary, heterogeneous timeseries that exhibits mutual influence on each other are integrated into CRN to model their interdependencies. By exploiting the BLSTM cell as its recurrent layer, two streams of the inputs in the generator are encouraged to “communicate” with each other. CRN is therefore capable of exploiting the interplay between mixed-type data that correlates over time.

#### JOINT TRAINING AND OPTIMIZATION

The overall architecture of EHR-M-GAN is shown in Fig. 1. In this section, we give a detailed description of the entire network’s structure and the optimization objective of the model. The steps for the training and optimization of EHR-M-GAN are as follows:

- The *pretraining of dual-VAE*: First, a dual-VAE network which consists of a pair of encoders ( $Enc^C, Enc^D$ ) and decoders ( $Dec^C, Dec^D$ ) is pretrained with both continuous and discrete data. Based on multiple objective constraints in Eq. 8 (for EHR-M-GAN<sub>cond</sub> the objective function can be referred in Eq. 9), a shared latent space is learnt using *dual-VAE*, where the gap between the embedding representations ( $\mathbf{z}_{1:T}^C, \mathbf{z}_{1:T}^D$ ) from both domains is minimized.
- The *generation of latent representations based on CRN*: Then, during the joint training stage, the *sequentially coupled generator* which is built based on CRN, takes the random noise vector ( $\hat{\mathbf{z}}_{1:T}^C, \hat{\mathbf{z}}_{1:T}^D$ ) as inputs and iterating across the timesteps  $t \in \{1, 2, \dots, T\}$  by the internal transition functions. Therefore, the synthetic latent embedding representations ( $\hat{\mathbf{z}}_{1:T}^C, \hat{\mathbf{z}}_{1:T}^D$ ) for both continuous and discrete data can be obtained.
- The *decoding for the mixed-type timeseries*: Next, the generated latent embeddings ( $\hat{\mathbf{z}}_{1:T}^C, \hat{\mathbf{z}}_{1:T}^D$ ) are further fed into the pretrained decoders ( $Dec^C, Dec^D$ ) and decoded into the corresponding synthetic patient trajectories ( $\hat{\mathbf{x}}_{1:T}^C, \hat{\mathbf{x}}_{1:T}^D$ ) in the observational space.
- The *adversarial loss update based on the discriminators*: Finally, the adversarial loss can be calculated from the LSTM network-based discriminators  $D^C$  and  $D^D$  by distinguishing between the real samples and synthetic timeseries for both data types.

The mathematical expression for the min-max objectives in EHR-M-GAN is provided as follows:

$$\begin{aligned} \min_G \max_D V_{\text{EHR-M-GAN}} &= \\ &\mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}} [\log D^C(\mathbf{x}^C) + \log D^D(\mathbf{x}^D)] + \\ &\mathbb{E}_{\mathbf{v} \sim p_{\mathbf{v}}} [\log(1 - D^C(\hat{\mathbf{x}}^C)) + \log(1 - D^D(\hat{\mathbf{x}}^D))] \end{aligned} \quad (12)$$

#### 2.5 CONDITIONAL VERSION OF EHR-M-GAN

For the conditional extension of EHR-M-GAN<sub>cond</sub>, the auxiliary label information is first used during the pretraining step of *dual-VAE*. Both the encoders and decoders condition on the auxiliary (one-hot) labels from  $\mathcal{L}$ , to make the model better adapted to particular contexts. In *dual-VAE*, the additional semantic loss is also incorporated during the optimization for the shared latent space (see Eq. 9).



Meanwhile, the same conditional labels are also applied in the sequentially coupled generator and discriminators, where the classes are fed as additional inputs through concatenation, as in the original CGAN architecture proposed by Mirza et al Mirza and Osindero (2014).

The t-SNE visualisation of the latent embeddings induced from *dual-VAE* can be found in Supplementary materials (see Section S.4.C), which indicates that the conditional information carried into EHR-M-GAN<sub>cond</sub> can be beneficial when synthesizing patient trajectories under specific medical conditions. Overall, the adversarial loss for EHR-M-GAN<sub>cond</sub> can be denoted as follows:

$$\begin{aligned} \min_G \max_D V_{\text{EHR-M-GAN}_{\text{cond}}} = & \\ & \mathbb{E}_{\mathbf{y}, \mathbf{x} \sim p_{\mathbf{y}, \mathbf{x}}} [\log D^C(\mathbf{x}^C | \mathbf{y}) + \log D^D(\mathbf{x}^D | \mathbf{y})] + \\ & \mathbb{E}_{\mathbf{y} \sim p_{\mathbf{y}}, \mathbf{v} \sim p_{\mathbf{v}}} [\log(1 - D^C(\hat{\mathbf{x}}^C | \mathbf{y})) + \log(1 - D^D(\hat{\mathbf{x}}^D | \mathbf{y}))] \end{aligned} \quad (13)$$

The pseudocodes for dual-VAE and EHR-M-GAN are provided in the Supplementary materials (see Section S.1.E).

### 3 DATA AND EVALUATION

#### 3.1 DATASET DESCRIPTION

The following three publicly accessible ICU datasets are used for evaluating the performance of EHR-M-GAN in generating the longitudinal EHR data:

- **MIMIC-III** (Medical Information Mart for Intensive Care) Johnson et al. (2016) — a freely accessible database that comprises de-identified EHRs associated with approximately 60,000 ICU admitted patients and 312 million observations to Beth Israel Deaconess Medical Center.
- **eICU** (eICU Collaborative Research Database) Pollard et al. (2018) — a multi-center critical care database containing data for over 200,000 admissions and 827 million observations to ICUs from 208 hospitals located throughout the United States.
- **HiRID** (High time-resolution ICU dataset) Yèche et al. (2021) — a high-resolution ICU dataset relating to more than 3 billion observations from almost 34,000 ICU patient admissions, monitored at the Department of Intensive Care Medicine, Bern University Hospital, Switzerland.

All these critical care databases include vital sign measurements, laboratory tests, treatment information, survival records, and other routinely collected data from hospital EHR systems. From these clinical observations, we featurize the patient trajectories as the combination of continuous-valued physiological timeseries (such as heart rate, oxygen saturation, and measurements from blood gas tests) and discrete-valued medical intervention timeseries (such as the usage of therapeutic devices or intravenous medications). Temporal trajectories **24 hours prior to patients' ICU endpoints** (discharge or death) are extracted for the three critical care databases. Data are preprocessed following an open-source framework — MIMIC-Extract Wang et al. (2020), where the patients' physiological and intervention signals are hourly aggregated for denser representations. Details on data curation, including the cohort selection criteria, full list of features, and imputation method, are explained in Supplementary Materials (see S.2 Datasets). Overall, the summarising statistics of the finalised cohorts for three databases are shown in Table 1.

#### 3.2 BASELINE MODELS

We compare the performance of EHR-M-GAN with eight state-of-the-art generative methods in literature. However, as these benchmarks typically face challenges when modeling mixed-type EHR timeseries Hjelm et al. (2017) and can only synthesize single-type EHRs, we draw the comparison between EHR-M-GAN and the benchmark models using the corresponding partial component of our synthetic results, i.e., either the continuous-valued part or the discrete-valued part. For continuous-valued timeseries generation, benchmark GAN models include C-RNN-GAN Mogren (2016), R(C)GAN Esteban et al. (2017) and TimeGAN Yoon et al. (2019). For discrete-valued timeseries generation, classic medGAN Choi et al. (2017), seqGAN Yu et al. (2017), and two

Table 1: **Summary of the cohorts after preprocessing on three critical care databases.** Number of patients and ICU admissions, as well as the dimensions of continuous-valued and discrete-valued variables, are provided for each dataset. Temporal trajectories 24 hours prior to patients’ ICU endpoints are extracted for the three critical care databases. Note that only the first ICU admission is selected for each patient. The dimension of the continuous- and discrete-valued data are provided. The conditional labels for training  $\text{EHR-M-GAN}_{\text{cond}}$  and the corresponding counts for each class are also listed.

	Number of patients	Number of ICU admissions	Dimension of continuous-valued variables	Dimension of discrete-valued variables	Conditional labels	Counts
MIMIC-III	28,344	28,344	78	20	ICU mortality	1,870 (6.59%)
					Hospital mortality	911 (3.21%)
					30-day readmission	1,122 (3.95%)
					No 30-day readmission	24,441 (86.22%)
eICU	99,015	99,015	55	19	ICU mortality	4,500 (4.54%)
					Hospital mortality	3,291 (3.32%)
					Hospital discharge	91,224 (92.13%)
HiRID	14,129	14,129	50	39	ICU mortality	1,266 (8.96%)
					ICU discharge	12,963 (91.74%)

recently proposed work — SynTEG Zhang et al. (2021) and DualAAE Lee et al. (2020) are used for comparison. Apart from these GAN-based models, we also incorporate PrivBayes Zhang et al. (2017) to synthesize discrete-valued timeseries, which falls in the class of non-GAN generative approaches using a Bayesian framework Tucker et al. (2020). As the original paper of PrivBayes focuses on data anonymization using differential privacy, we therefore implemented its ‘Non-Private’ version for a fair comparison with other baselines (see Section 4.1 Non-Private Methods in Zhang et al. (2017)). For medGAN and PrivBayes, we feed the flattened temporal sequence as the input since the models cannot produce timeseries data.

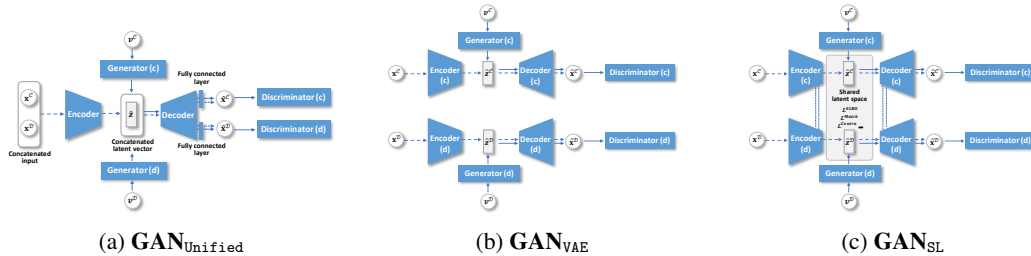


Figure 2: **The network architectures in the ablation study.** Three variants of EHR-M-GAN are implemented in the ablation study. Compared with the full model of EHR-M-GAN,  $\text{GAN}_{\text{Unified}}$  learns the joint representations of heterogeneous types of data in a unified network;  $\text{GAN}_{\text{VAE}}$  maintains the basic architecture of EHR-M-GAN, but ignore the dependency learning (i.e., separate networks for two streams of inputs are trained in parallel);  $\text{GAN}_{\text{SL}}$  constructs the shared latent space using the *dual-VAE* module but omit the *sequentially coupled generator* for learning the temporal correlations in the mixed-type timeseries.

We further perform the **ablation study** to investigate whether our introduced novel components in the proposed model have advantages over its variants that also model mixed-type EHRs. First, as EHR-M-GAN learns the joint representations from heterogeneous types of data using separate (but inherently correlated and weights-sharing) VAE networks, we compare it with a variant that jointly models the mixed-type data using a single unified VAE network (denoted as  $\text{GAN}_{\text{Unified}}$ ). Then, we test the variant that encodes the mixed-type inputs separately with two independent VAE networks, and then combines the resulted synthesis of different data types as outputs (denoted as  $\text{GAN}_{\text{VAE}}$ ). Lastly, we assess the effectiveness of the proposed *dual-VAE* component in our model alone by implementing  $\text{GAN}_{\text{SL}}$ . The architectures of different variants of EHR-M-GAN in the ablation study are detailed as follows (also see Fig. 2 for illustration):

- $\text{GAN}_{\text{Unified}}$ : It contains a unified VAE module and two separate GANs. The continuous-valued and discrete-valued timeseries is concatenated together, via normalization and one-hot encoding, as input to the encoder in the unified VAE network. The decoder receives the concatenation of the generated latent vectors as the input, and then decodes it into synthetic timeseries with the corresponding data types using the separate fully connected

---

layers. Each component in the architecture of  $\text{GAN}_{\text{Unified}}$  (unified encoder and decoder, separate generators and discriminators) is implemented with LSTMs, which are the same as EHR-M-GAN.

- $\text{GAN}_{\text{VAE}}$ : It is composed of a pair of VAE networks and GANs (one for each type of inputs). The continuous-valued timeseries and discrete-valued timeseries from the same patients are separately fed into the corresponding paths in  $\text{GAN}_{\text{VAE}}$ , and then run in parallel. The synthetic outputs of each data type are then combined as the final results. It maintains the basic structure of EHR-M-GAN but lacks the latent space sharing with *dual-VAE* and the *sequentially coupled generator* in the original EHR-M-GAN.
- $\text{GAN}_{\text{SL}}$ : In addition to  $\text{GAN}_{\text{VAE}}$ , it learns the shared latent space representations through *dual-VAE* by adding the corresponding loss functions in EHR-M-GAN, including *ELBO loss*, *Matching loss* and *Contrastive loss*. This model lacks the *sequentially coupled generator*.
- EHR-M-GAN: In addition to  $\text{GAN}_{\text{SL}}$ , it incorporates the *sequentially coupled generator* for the learning the correlated temporal dynamics in timeseries of different data types. This is the proposed full model.
- EHR-M-GAN<sub>cond</sub>: This version is implemented on the basis of conditional GAN Mirza and Osindero (2014), where the conditional inputs are fed into EHR-M-GAN to generate patients under specific labels.

For training EHR-M-GAN<sub>cond</sub>, auxiliary information from the patient status is used as conditional input. These conditional inputs are selected since synthesizing EHR information of patient subgroups with potential outcomes would be valuable for clinicians in their decision-making process. Other conditional labels (such as patient demographics in the categorized format) can also be used in the proposed conditional synthesizer for other research purposes. For MIMIC-III dataset, the classes are (1) *ICU mortality*: patient died within the ICU; (2) *Hospital mortality*: patient discharged alive from the ICU, and died within the hospital; (3) *30-day readmission*: patient discharged alive from the hospital, and readmitted to the hospital within 30 days; (4) *No 30-day readmission*: patient discharged alive from the hospital, and had no readmission record to the hospital within 30 days. For eICU and HiRID datasets, the corresponding labels are also extracted based on the availability of the patient outcomes (see Table 1).

### 3.3 EVALUATION METRICS

Evaluating GAN models is a notoriously challenging task. Advantages and pitfalls of commonly used evaluation metrics for GANs are discussed in Borji (2021). In this work, a systematic evaluation framework is adopted to assess the quality of synthetic patient EHRs with respect to its *fidelity*, *correlation*, *utility*, and *privacy* (see Table 2). We first individually assess the representativeness of the synthetic continuous-valued and discrete-valued timeseries. This includes measuring the distance between underlying data distributions (such as *Maximum mean discrepancy* and *Dimension-wise probability*), comparing the feature-level statistics between the real and synthetic data (*Patient trajectories*), and assessing the indistinguishability of the synthetic data to the true data (i.e., *Discriminative score*). Secondly, we evaluate to which extent our model can reconstruct the interdependency between different features (*Pearson pairwise correlations*), and the temporal dynamics in the patient trajectories (*Autocorrelation function*), by using a set of qualitative and quantitative metrics. Thirdly, we introduce data augmentation by incorporating synthesized EHR timeseries under various settings, and quantitatively assess the improvement provided by EHR-M-GAN in the *Downstream tasks* for medical intervention prediction in the ICU (i.e., the utility of the synthetic data). Lastly, we measure the attribute of patient privacy-preserving of EHR-M-GAN under *Membership inference attack*. We also evaluate the performance of the same downstream tasks under *Differential privacy* guarantees (See Fig. 1c and Table 2 for the evaluation pipeline).

## 4 RESULTS

### 4.1 MAXIMUM MEAN DISCREPANCY

To measure the similarity between the continuous-valued synthetic data and the real data, maximum mean discrepancy (MMD) is used. MMD can assess whether two sets of samples are from the same

Table 2: **Summary of the evaluation protocol in this study.** A comprehensive set of evaluation metrics are used to test the *Fidelity*, *Correlation*, *Utility* and *Privacy* of the synthetic EHR data. Definitions of evaluation metrics for corresponding data types are explained. The last column illustrates when the corresponding evaluation metric indicates better performance.

	Evaluation metric	Data type	Definition	Better performance
<b>Fidelity</b>	Maximum mean discrepancy	Continuous	A kernel-based statistic is calculated to determine whether the real and synthetic data are from the same distribution.	Lower value
	Dimension-wise probability	Discrete	The Bernoulli success probability of each feature dimension (i.e., the probability of whether the treatment is active) at the given timestamp is calculated. Probabilities from real and synthetic data are represented on the x and y axis in a single plot to compare the consistency.	Scatters closer to the diagonal line (lower RMSE and CC)
	Discriminative score	Continuous and discrete	A classifier is trained post-hoc to tell the difference between the real and synthetic data with its accuracy calculated.	Lower accuracy
	Patient trajectories	Continuous and discrete	The mean and standard deviation per time point of real and synthetic patient trajectories are compared and visualized.	Similar distributions between the real and synthetic data
<b>Correlation</b>	Pearson pairwise correlations	Continuous and discrete	The correlation between different features is calculated and visualized in a heatmap for both real and synthetic data.	Heatmaps corresponding to real and synthetic data more similar (lower CorAcc and $\mu_{abs}$ )
	Autocorrelation function	Continuous and discrete	The correlation between the timeseries and its lagged version is calculated and visualized as an ACF curve for both real and synthetic data.	ACF curves corresponding to real and synthetic data more similar
<b>Utility</b>	TSTR (downstream task)	Continuous and discrete	The downstream classifier is trained which uses synthetic data as training set, and (hold-out) real data as test set. The result is compared with TRTR to see whether it can maintain the same.	Higher AUROCs (with TRTR as baselines)
	TSRTR (downstream task)	Continuous and discrete	The downstream classifier is trained which uses real data and synthetic data as training set, and (hold-out) real data as test set. The result is compared with TRTR to see whether the performance can be improved.	Higher AUROCs (with TRTR as baselines)
<b>Privacy</b>	Membership inference attack	Continuous and discrete	A threat model is trained under the black-box setting to determine whether a record is used for training GANs. This quantifies the risk of sensitive information from real data being revealed by synthetic data.	Lower accuracy or recall
	Differential privacy	Continuous and discrete	The downstream classifier is trained with differential privacy guarantee. The result is compared with TRTR to see whether it can maintain the same.	Higher AUROCs (with TRTR as baselines)

distributions, and in our case, one from the true data  $x$  and one from synthetic data  $x'$  generated by GANs. To calculate the statistics, a kernel function  $K : X \times X' \rightarrow \mathbb{R}$  is used to quantify the similarity between the two distributions. In this study, a sum of Gaussian kernel sets is adopted following the implementations in Sutherland et al. (2016), which can be expressed as:

$$K(\mathbf{x}, \mathbf{x}') = \sum_i \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_F^2}{\sigma_i^2}\right) \quad (14)$$

where  $\sigma_i$  is the value of the  $i$ -th selected bandwidth for calculating MMD. As in our study, the real and synthetic samples are multivariate timeseries aligned along the fixed time axis (i.e., 24 data points per patient), we therefore handle these multivariate timeseries as matrices and use the kernel function to calculate the Frobenius norm ( $\|\cdot\|_F$ ) between them Esteban et al. (2017).

Finally, given samples  $\{\mathbf{x}_i\}_{i=1}^N$  from real distributions, and samples  $\{\mathbf{x}'_j\}_{j=1}^M$  from the synthetic distributions (with  $N$  and  $M$  denoting the corresponding sample sizes), the estimation of MMD can be defined as:

$$\begin{aligned} \widehat{\text{MMD}}^2 &= \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n K(\mathbf{x}_i, \mathbf{x}_j) - \frac{2}{mn} \sum_{i=1}^n \sum_{j=1}^m K(\mathbf{x}_i, \mathbf{x}'_j) \\ &+ \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m K(\mathbf{x}'_i, \mathbf{x}'_j) \end{aligned} \quad (15)$$

It can be inferred from the equation [15] that higher similarity between the two distributions leads to the lower MMD value, with the lower bound value zero indicating that the two distributions are identical.

Table 3: **Maximum mean discrepancy (MMD) of continuous-valued synthetic data.** Lower values of MMD indicate models which can better learn the distribution of the real data.

	C-RNN-GAN	R(C)GAN	TimeGAN	GAN <sub>Unified</sub>	GAN <sub>VAE</sub>	GAN <sub>SL</sub>	EHR-M-GAN	EHR-M-GAN <sub>cond</sub>
MIMIC-III	1.038 ± 0.013	0.971 ± 0.029	0.694 ± 0.025	0.893 ± 0.027	0.926 ± 0.038	0.745 ± 0.040	<b>0.692 ± 0.034</b>	<b>0.604 ± 0.027</b>
eICU	1.139 ± 0.023	1.106 ± 0.042	0.672 ± 0.038	0.850 ± 0.032	0.842 ± 0.029	0.670 ± 0.034	<b>0.651 ± 0.026</b>	<b>0.540 ± 0.018</b>
HiRID	0.982 ± 0.017	0.865 ± 0.020	0.470 ± 0.024	0.518 ± 0.030	0.532 ± 0.035	0.508 ± 0.028	<b>0.428 ± 0.015</b>	<b>0.389 ± 0.024</b>

As indicated in Table 3, EHR-M-GAN outperforms the state-of-the-art benchmarks among all three datasets in synthesizing continuous-valued timeseries. The conditional version — EHR-M-GAN<sub>cond</sub> further boosts the performance of the model by leveraging the information of the condition-specific inputs. Furthermore, as shown in the ablation study, EHR-M-GAN and EHR-M-GAN<sub>cond</sub> produce smaller MMD values when compared to their variants. Using MIMIC-III as an example, compared with the basic model GAN<sub>VAE</sub>, by integrating the shared latent space learning using *dual-VAE* under multiple loss constraints, the performance of GAN<sub>SL</sub> significantly improves (GAN<sub>SL</sub> vs. GAN<sub>VAE</sub>, 0.745 to 0.926,  $p$ -value < 0.05 from  $t$ -test<sup>3</sup>). By further building the sequentially coupled generator based on BLSTMs and exploiting the information within mixed-type data, the MMD of EHR-M-GAN shows a nearly 24% improvement over GAN<sub>VAE</sub>. When synthesizing mixed-type timeseries using the unified network, the performance of GAN<sub>Unified</sub> for generating continuous-valued timeseries lags behind the proposed EHR-M-GAN. It therefore can be inferred that, compared with EHR-M-GAN which extracts useful hierarchical representations for each data type using tailored encoding layers, it is quite challenging for GAN<sub>Unified</sub> to learn marginal distributions from raw mixed-type timeseries with a unified architecture.

## 4.2 DIMENSION-WISE PROBABILITY

To evaluate the representativeness of the synthetic discrete-valued timeseries, the dimension-wise probability test is employed. To test the probability distributions between the real and synthetic binary features, the Bernoulli success probability  $p \in [0, 1]$  is calculated for the discrete-valued timeseries, and is visualized through scatterplot. As a sanity check, it investigates if the probability of the medical intervention being active at the given timestamps is matched between the real data ( $x$ -axis) and synthetic data ( $y$ -axis). The correlation coefficients (CCs) and root-mean-square errors (RMSEs) are also adopted Baowaly et al. (2019) based on the Bernoulli success probabilities to quantitatively measure the distribution divergence between real and synthetic data.

As shown in Fig. 3 (see Fig. S4 and S5 for more results on eICU and HiRID datasets), the optimal results are provided by EHR-M-GAN and EHR-M-GAN<sub>cond</sub>. The close-to-real probability distributions that appear along the diagonal line indicate the remarkable similarity between the real data and the synthetic data provided by our models. The quantified CC and RMSE also correspond with the visualisation results, which are close to the highest mark (EHR-M-GAN: RMSE = 0.0095, CC = 0.9973). Similar to the results in MMD, the dimensional-wise distributions are better captured when modules such as *dual-VAE* and *sequentially coupled generator* are introduced in EHR-M-GAN. GAN<sub>Unified</sub> suffers from mode collapse (the generator fails to produce outputs with sufficient diversity), and therefore shows poor performance compared with other variants when synthesizing discrete-valued timeseries. As the mixed-type features are treated as unimodal input without differentiating their heterogeneous nature, no marginal representations are explicitly learned.

Among all state-of-the-art benchmark models, DualAAE shows the best result but is slightly sub-optimal when compared to EHR-M-GAN. In contrast, both skewed distribution and low performance scores are observed in medGAN, as it lacks the ability to capture the temporal correlations within timeseries. SynTEG shows improved performance over medGAN, as it is capable of synthesizing discrete-valued features in EHRs with timestamps. The non-GAN generative method PrivBayes also shows good performance among all the benchmark synthesizers when modeling the underlying probability distribution of the discrete-valued EHR timeseries. On the other hand, despite the well-known performance of SeqGAN in natural language generation, it is not quite applicable in synthesizing sequential clinical EHRs.

<sup>3</sup>Unpaired (two-sample)  $t$ -test with a significance level of 0.05 is used throughout the paper unless specified otherwise.

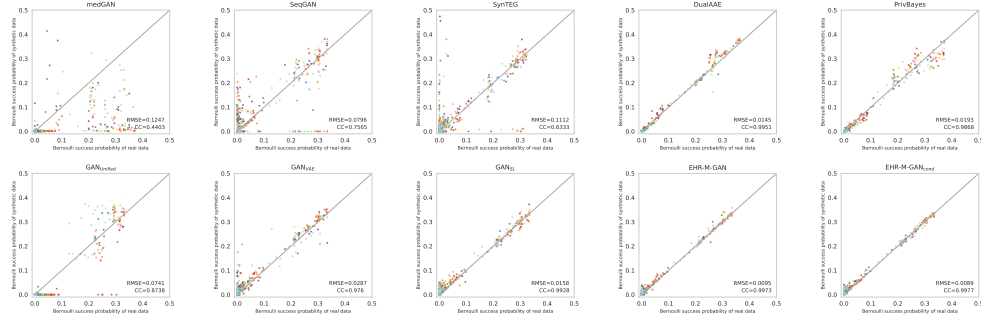


Figure 3: **Scatterplot of the dimension-wise probability test on MIMIC-III dataset.** Dimension-wise probability calculates the Bernoulli success probability of each dimension, i.e., the probability of the treatment being active at a particular time. The x-axis and y-axis represent dimension-wise probability for the real data and synthetic data generated from different models, respectively. The same color indicates the same treatment (but with varying timestamps). The optimal performance appears along the diagonal line. The corresponding CCs ( $[0, 1]$ , the higher the better) and RMSEs ( $[0, +\infty)$ , the lower the better) are also calculated to quantify the probability distribution similarities between the real and synthetic EHRs timeseries. Dimension-wise probability plot for eICU and HiRID dataset can be found in Supplementary materials (see S.4.A).

Generating discrete-valued features are known to be problematic for traditional GANs. Due to their limitation in passing the gradients from the critic models, vanilla GANs cannot update their generators efficiently based on the adversarial loss Yu et al. (2017); Choi et al. (2017). However, the result of EHR-M-GAN shows its superiority in explicitly capturing each dimension of the discrete-valued sequences. EHR-M-GAN mitigates this problem by learning the shared latent representations using *dual-VAE*. Discrete-valued timeseries are encoded into a gradient-differentiable space for further optimizing the generators and thus solving the problem.

### 4.3 PATIENT TRAJECTORIES

We compare the distribution of patient trajectories per timepoint between the real data and synthetic data generated by EHR-M-GAN for the MIMIC-III dataset. Five commonly measured vital sign and laboratory features — *Oxygen Saturation*, *Systolic Blood Pressure*, *Respiratory Rate*, *Heart Rate*, *Temperature*, as well as two medical intervention features — *Mechanical Ventilation* and *Vasopressor*, are considered and compared as an exemplar in Fig. 5. It can be inferred that the proposed model can accurately capture the statistical distribution (mean and standard deviation) of both continuous-valued and discrete-valued features. The temporal dynamics are well-preserved in the synthetic timeseries. For example, the variance of *Oxygen Saturation* gradually increases towards the ICU endpoints in the real data, and is closely reflected in the synthetic timeseries. Furthermore, EHR-M-GAN<sub>cond</sub> shows superior performance as it can generate correct trajectories with specific patient conditions (see section S.4.D in Supplementary Materials for results).

### 4.4 DISCRIMINATIVE SCORE

For both continuous-valued and discrete-valued data, the discriminative score is measured as the accuracy of a discriminator trained post-hoc to separate real from generated samples. Synthetic data are generated with the same amount of the hold-out test set from the original data, and are labeled as *synthetic* and *real* correspondingly to train the binary classifier. In this study, the classifier (critic) is implemented with a single-layered Bi-directional Long Short-Term Memory (Bi-LSTM) model (i.e., *many-to-one*), with its parameters randomly initialized (as opposed to critic built upon representations from the trained generative model Zhang et al. (2022)). The critic trained from the supervised learning task can be used to characterize the temporal correlations across the patient EHR timeseries.

As indicated from the results in Table 4, it appears that EHR-M-GAN and EHR-M-GAN<sub>cond</sub> can produce synthetic data that are less distinguishable from real data than the benchmarked models.

Table 4: **Discriminative score of synthetic data.** A discriminative model is trained post-hoc to discriminate between synthetic samples and real samples. The accuracy from the discriminative classifier is used as the discriminative score, where the lower value indicates better performance. The result is bounded by 0.5 when the classifier cannot distinguish between two distributions.

	Method	MIMIC-III	eICU	HiRID
Continuous-valued synthetic data	C-RNN-GAN	0.825 ± 0.013	0.876 ± 0.010	0.774 ± 0.022
	R(C)GAN	0.833 ± 0.028	0.850 ± 0.021	0.742 ± 0.016
	TimeGAN	0.763 ± 0.018	0.790 ± 0.013	<b>0.716 ± 0.024</b>
	GAN <sub>Unified</sub>	0.809 ± 0.023	0.863 ± 0.027	0.749 ± 0.014
	GAN <sub>VAE</sub>	0.842 ± 0.020	0.871 ± 0.014	0.802 ± 0.017
	GAN <sub>SL</sub>	0.786 ± 0.016	0.813 ± 0.023	0.752 ± 0.021
	EHR-M-GAN	<b>0.746 ± 0.018</b>	<b>0.776 ± 0.015</b>	0.724 ± 0.015
	EHR-M-GAN <sub>cond</sub>	<b>0.729 ± 0.025</b>	<b>0.745 ± 0.017</b>	<b>0.693 ± 0.012</b>
Discrete-valued synthetic data	medGAN	0.903 ± 0.027	0.915 ± 0.034	0.896 ± 0.021
	seqGAN	0.937 ± 0.025	0.924 ± 0.023	0.913 ± 0.027
	SynTEG	0.879 ± 0.021	0.902 ± 0.030	0.878 ± 0.025
	DualAAE	0.847 ± 0.029	0.860 ± 0.033	0.829 ± 0.024
	PrivBayes	0.859 ± 0.036	0.883 ± 0.034	0.832 ± 0.017
	GAN <sub>Unified</sub>	0.890 ± 0.022	0.907 ± 0.026	0.849 ± 0.015
	GAN <sub>VAE</sub>	0.862 ± 0.024	0.881 ± 0.029	0.824 ± 0.018
	GAN <sub>SL</sub>	0.829 ± 0.032	0.844 ± 0.028	0.816 ± 0.025
	EHR-M-GAN	<b>0.813 ± 0.026</b>	<b>0.831 ± 0.024</b>	<b>0.802 ± 0.020</b>
	EHR-M-GAN <sub>cond</sub>	<b>0.784 ± 0.024</b>	<b>0.803 ± 0.022</b>	<b>0.778 ± 0.019</b>

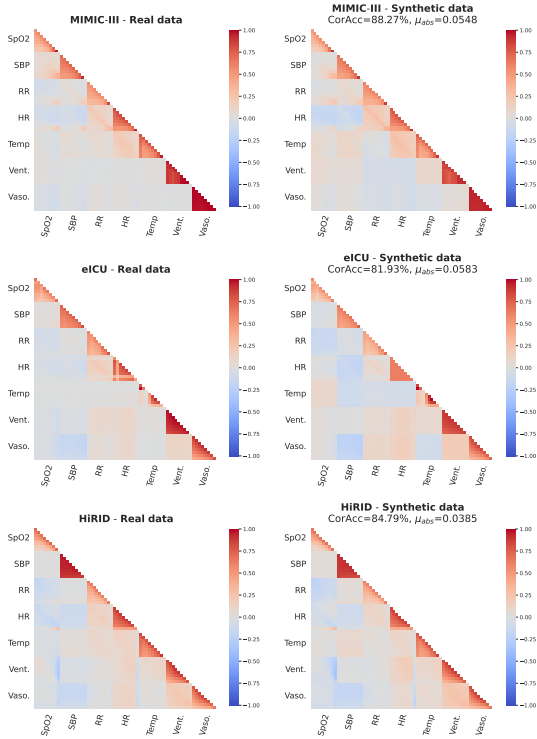
Especially for EHR-M-GAN<sub>cond</sub>, it achieves the optimal discriminative scores consistently against other benchmarks for both continuous-valued and discrete-valued timeseries. For discrete-valued data generation, EHR-M-GAN-generated samples achieve the discriminative score of 0.813 on the MIMIC-III dataset, which has a 4% statistically significant improvement over the best performing benchmark (EHR-M-GAN vs. DualAAE: 0.813 to 0.847,  $p < 0.05$ ). The overall discriminative scores produced by PrivBayes on three ICU databases are comparable with the GAN models such as SynTEG and DualAAE. For continuous-valued timeseries generation, the discriminative score of TimeGAN on HiRID dataset outperforms the other models as well as EHR-M-GAN, though not statistically significant (EHR-M-GAN vs. TimeGAN: 0.724 to 0.716,  $p = 0.4374$ ). By leveraging the additional information from the conditional inputs, EHR-M-GAN<sub>cond</sub> can provide significantly better result than TimeGAN (EHR-M-GAN<sub>cond</sub> vs. TimeGAN: 0.693 to 0.716,  $p < 0.05$ ).

The ablation study has proved the effectiveness of EHR-M-GAN for generating high quality EHR timeseries. The shared latent space representation learning in the *dual-VAE* (i.e., GAN<sub>SL</sub>) have shown remarkable success as making the synthetic data more realistic than separately generating the latent embeddings based on VAEs (as in GAN<sub>VAE</sub>). The *sequentially coupled generator* further improves the model by capturing the dynamics between mixed-type data and iterating over time, therefore enabling the synthetic timeseries to become more indistinguishable from the original. Further compared with GAN<sub>Unified</sub> that models the mixed-type data in a unified network, our proposed model enables effective learning for the marginal distributions from each data type. More importantly, EHR-M-GAN can leverage its dependency learning components to explicitly capture the correlations between heterogeneous types of data.

#### 4.5 INTERDEPENDENCY CHARACTERISTICS

In this section, we first employ Pearson pairwise correlation (PPC), which ranges from -1 to 1, to evaluate how closely the synthetic data can model the correlations between continuous-valued and discrete-valued timeseries. Timestamps of the patient trajectories are extracted with every 3 hours interval in a total of 24 hours ICU stay, to explore the temporal dependencies within different variables. To further quantitatively measure the difference between heatmaps generated from real and synthetic samples, we calculate the mean value of the absolute difference between the two PCC matrices ( $\mu_{abs}$ ). We also adopted correlation accuracy (*CorAcc*) Tao et al. (2021) which quantifies the similarity of two heatmaps within the range of 0 to 1. We discretize the correlation coefficients into 6 correlation levels: *strong negative* ( $[-1, -0.5)$ ), *middle negative* ( $[-0.5, -0.3)$ ), *low negative* ( $[-0.3, -0.1)$ ), *no correlation* ( $[-0.1, 0.1)$ ), *low positive* ( $[0.1, 0.3)$ ), *middle positive* ( $[0.3, 0.5)$ ), and

strong positive ( $(0.5, 1)$ ). Then,  $CorAcc$  can be calculated as the percentage of pairs where the real and synthetic data is assigned to the same correlation level.



**Figure 4: Pearson pairwise correlation (PPC) between continuous-valued and discrete-valued timeseries.** The plots contrast the PPC calculated within the real data (left column) and the synthetic data generated by EHR-M-GAN (right column). Besides the visual inspection, the similarity between two heatmaps are quantified by  $CorAcc$  and  $\mu_{obs}$ . These metrics indicate how well the synthetic data reconstruct the correlations observed in the real patient trajectories. As shown in this figure,  $SpO_2$ ,  $SBP$ ,  $RR$ ,  $HR$ ,  $Temp$  represents *Oxygen Saturation*, *Systolic Blood Pressure*, *Respiratory Rate*, *Heart Rate*, *Temperature*, respectively. And  $Vent.$  and  $Vaso.$  corresponds to *Vasopressor* and *Mechanical Ventilation*. PPC is calculated every 3 hours over the total 24 hours of ICU stay (ticks of the timestamps are omitted).

*Rate*, *Oxygen Saturation*, and *Systolic Blood Pressure*, the positive ACF coefficients rapidly decrease within the period of first few hours, followed by the growing trends of negative temporal correlation. The lag with the lowest correlation coefficient is identified at approximately 4 hours. Specifically, global peaks appear roughly at the 12-hour ticks of *Temperature* for both real and synthetic data on three critical databases. Meanwhile, the negative correlation strengthens as the time lag increase for *Mechanical ventilation* in the original timeseries. Since these behaviours can be reproduced by EHR-M-GAN, therefore they demonstrate that our model can effectively capture the temporal characteristics in the original timeseries.

As observed, correlation trends over distinctive features are closely reflected by the synthetic data, with the quantitative measure  $CorAcc$  consistently exceed 0.8 on three critical care databases. It is also worth noticing that EHR-M-GAN can successfully recover temporal dependencies with a high granularity from real patient trajectories. For example, synchronized correlations across timestamps are observed between *Respiratory Rate* and *Heart Rate* in the MIMIC-III dataset. Such trends are preserved in synthetic data. This can be explained by the common regulation of these two features by the autonomic nervous system and their synchronized increase in cases of physiological stress, such as hypoxemia. In summary, the proposed EHR-M-GAN can reconstruct the temporal dynamics and correlations between features in the real data, which is valuable for downstream ML-based classification and prediction applications.

Then, autocorrelation functions (ACF) Benedetti et al. (2020) and the corresponding root-mean-square errors (RMSEs) are calculated to show how EHR-M-GAN can capture the temporal correlations among the timeseries. ACF measures the relationship between the timeseries and its lagged version. Fig. S6 - S8 in the Supplementary materials shows the ACF calculated for selected continuous-valued and discrete-valued variables (same as Pearson pairwise plot) on real and synthetic timeseries. The time lags are specified as the hourly intervals up to 24 hours before patients' ICU endpoints (ICU discharge or death). Additionally, RMSEs are calculated to quantitatively evaluate the similarity between the corresponding two curves produced by real data and synthetic data.

Similar patterns are presented between the ACF calculated for real data and their synthetic counterparts, while the quantitative statistics also correspond with the observation. Moreover, overlapping confidence intervals indicate that the synthetic data is able to consistently capture the underlying temporal distributions within the real timeseries. For variables such as *Heart*



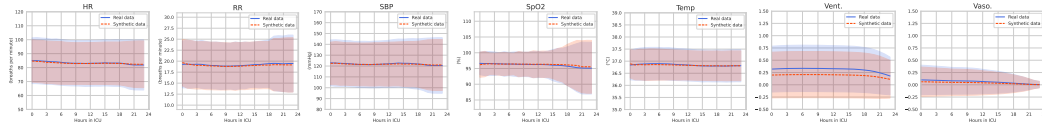


Figure 5: **Comparison of the distribution of values at each timepoint (mean and standard deviation) between real and synthetic patient trajectory produced by EHR-M-GAN.** Multivariate timeseries 24 hours before patients’ ICU endpoints are generated, including *Heart Rate*, *Respiratory Rate*, *Systolic Blood Pressure*, *Oxygen Saturation*, *Temperature*, *Mechanical Ventilation* and *Vasopressor*. The mean value of the real/synthetic feature at each timepoint is plotted by the solid/dotted line, with the shaded area indicating  $\pm 1$  standard deviation. For *Mechanical Ventilation* and *Vasopressor*, the y-axis indicates the probability distribution of such intervention being applied ("On") at a given time. The synthetic patient trajectories generated by EHR-M-GAN<sub>cond</sub> under different conditions can be found in Supplementary materials.

#### 4.6 DOWNSTREAM TASKS

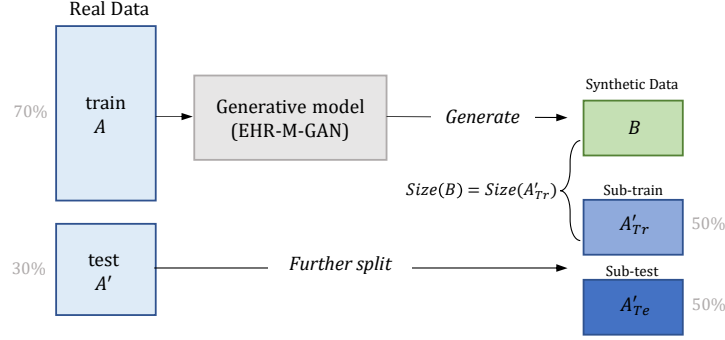
As previously discussed, one of the most prominent goals for GANs is to benefit the future downstream analyses in the real clinical application. A relevant question in the ICU is whether specialized medical treatments, such as therapeutic interventions or organ support, are required for critically ill patients during the admission. Accurate predictions on such tasks can help clinicians to provide actionable, in-time interventions in the resource-intensive ICU. Therefore in this section, **clinical intervention prediction** tasks are implemented to evaluate the potential of EHR-M-GAN and EHR-M-GAN<sub>cond</sub> in synthesizing high-fidelity synthetic data to further boost the performance of ML classifiers. In line with prior work Wang et al. (2020); Wu et al. (2017); Suresh et al. (2017), we establish LSTM-based classifiers to predict the status of *mechanical ventilation* and *vasopressors* using continuous-valued multivariate physiological signals as the predictors. A fixed duration of 12 hours is used for both observation window and prediction window (see Fig. 1). Four outcomes of medical intervention status are defined as: *Stay on*, *Onset*, *Switch off*, *Stay off* (detailed descriptions can be found in Fig. 1).

We partition the dataset as illustrated in Figure 6a, and the performances are assessed from two aspects (see Figure 6b): **(i) Traditional approach:** To explore whether the synthetic data can represent the real data accurately, we compare *Train on Real, Test on Real* (TRTR) with *Train on Synthetic, Test on Real* (TSTR), to show whether the performance of a classifier trained on synthetic data from EHR-M-GAN or EHR-M-GAN<sub>cond</sub> can be generalized to real data. In addition to the proposed models, synthetic data produced by the baseline models are also used to train the downstream classifiers for comparison. Other than a measurement of data utility where the downstream task is to predict discrete-valued medical intervention (described as outcomes in this scenario) using continuous-valued physiological features (denoted as predictors), TSTR can also be used to assess data synthesizers’ ability to capture the interdependencies between the mixed-type features. **(ii) Data augmentation approach:** As data augmentation is employed as a means of circumventing the issue caused by the under-resourced EHR data, here we explore whether synthetic data can be used to improve the existing ML algorithms through data augmentation. Therefore, *Train on Synthetic and Real, Test on Real* (TSRTR) is compared with TRTR to measure the improvement of the classifier’s performance when trained on the augmented data Esteban et al. (2017); Kiyasseh et al. (2020). The augmentation ratio  $\alpha$  or  $\beta$  is applied on sub-train data  $A'_{Tr}$ , or synthetic data  $B$ , in two different scenarios of TSRTR, respectively. Details are explained as follows (also see Figure 6b for illustration).

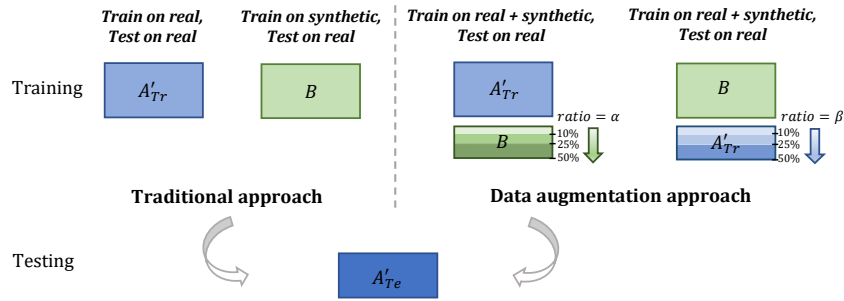
Firstly, as the dearth of data potentially degrades the performance of downstream classifiers, given that the real data has a limited and fixed sample size, we investigate whether adding synthetic EHR data provided by EHR-M-GAN and EHR-M-GAN<sub>cond</sub> can improve the training of downstream classifiers. **Ratio**  $\alpha$  indicates the portion of synthetic data (see Figure 6b) being used to augment the real data to improve the quality and robustness of the downstream classifiers.  $\alpha$  is set to be 10%, 25%, and 50%, representing the availability of synthetic samples provided for augmentation.

Secondly, the acquisition of healthcare data is generally time-consuming and expensive, therefore another overarching goal for the generative model is to minimize the efforts in collecting data. In this section, we investigate whether high-fidelity synthetic data can offer a viable solution for boosting

the downstream classifiers’ performance when the availability of real data is limited. This allows us to understand if the sample size required for real data collection can be reduced while maintaining sufficient predictive power through the use of synthetic data. During the experiment, the synthetic data  $B$  is given (to emulate the scenario where synthetic datasets are available for a particular clinical research purpose), which further is combined with limited real data (collected during clinical trial), to train the downstream classifiers (i.e., augment synthetic data with limited real data). Then by implementing EHR-M-GAN or EHR-M-GAN<sub>cond</sub> in *TSRTR*, we investigate the proportion of the real data  $A'_{Tr}$  (**ratio**  $\beta$ ) required to maintain the same performance as in *TRTR* based on the entire synthetic dataset  $B$  (see Figure 6b).



(a) Data splitting.



(b) Data augmentation scenarios.

**Figure 6: Downstream intervention prediction experimental setup. a. Data splitting.** During training stage, the real data is split into two sets with 70% training data  $A$  and 30% test data  $A'$ . The test data  $A'$  is further split into sub-train data  $A'_{Tr}$  and sub-test data  $A'_{Te}$  with equal size. Then, the synthetic data  $B$ , with size equal to the sub-train data  $A'_{Tr}$ , is synthesized by EHR-M-GAN (or EHR-M-GAN<sub>cond</sub>) trained on the real training data  $A$ . **b. Data augmentation scenarios.** Subsequent experiments are trained on set  $A'_{Tr}$ , or  $B$ , or  $A'_{Tr} \cup B$  and then tested on  $A'_{Te}$ . In traditional approach, results based on *Train on Real, Test on Real (TRTR)* and *Train on Synthetic, Test on Real (TSTR)* are compared to assess the generalisability of the synthetic data. In data augmentation approach, i.e., *Train on Synthetic and Real, Test on Real (TSRTR)*, we either augment real data  $A'_{Tr}$  with  $\alpha$  (augmentation ratio, 0 to 50%) of the synthetic samples  $B$ , or augment synthetic samples  $B$  with  $\beta$  (0 to 50%) of the real data  $A'_{Tr}$ .

**Traditional approach.** Table 5 compares the classification performances of predicting forthcoming medical interventions in the ICUs under the experimental setting of *TRTR* and *TSTR*. It is expected that the optimal AUROCs are achieved by the classifiers that are trained on real data. In comparison, the classifiers trained on the synthetic data provided by proposed models can achieve similar performances. More specifically, synthetic data generated by EHR-M-GAN<sub>cond</sub> demonstrates better generalisability when compared with EHR-M-GAN in the downstream application, such as the task of predicting *mechanical ventilation* on the HiRID dataset (*TRTR* vs. *TSTR* from EHR-M-GAN<sub>cond</sub>: 0.867 to 0.856, with  $p=0.3906$ ).

Table 5: **Downstream task evaluation.** Downstream tasks are evaluated under the training scenarios of *Train on Real, Test on Real (TRTR)* and *Train on Synthetic, Test on Real (TSTR)*. Prediction of two outcomes of interest – intervention by *Mechanical ventilation (Vent.)* and *Vasopressors (Vaso.)* are selected as exemplary tasks. Macro-AUROC is used to score the performance of the LSTM-based classifiers on the multi-class prediction tasks (labeled as *Stay on, Onset, Switch off, Stay off*).

Dataset	Treatments	Real data	GAN <sub>Unified</sub>	GAN <sub>VAE</sub>	GAN <sub>SL</sub>	EHR-M-GAN	EHR-M-GAN <sub>cond</sub>
MIMIC-III	Vent.	<b>0.894 ± 0.016</b>	0.724 ± 0.015	0.701 ± 0.018	0.728 ± 0.010	0.740 ± 0.009	0.823 ± 0.020
	Vaso.	<b>0.841 ± 0.009</b>	0.694 ± 0.012	0.651 ± 0.015	0.679 ± 0.009	0.725 ± 0.015	0.810 ± 0.019
eICU	Vent.	<b>0.868 ± 0.015</b>	0.697 ± 0.014	0.702 ± 0.009	0.718 ± 0.012	0.745 ± 0.008	0.795 ± 0.015
	Vaso.	<b>0.813 ± 0.018</b>	0.648 ± 0.011	0.657 ± 0.012	0.665 ± 0.014	0.706 ± 0.014	0.748 ± 0.017
HiRID	Vent.	<b>0.867 ± 0.012</b>	0.765 ± 0.014	0.747 ± 0.013	0.803 ± 0.008	0.825 ± 0.019	0.856 ± 0.033
	Vaso.	<b>0.878 ± 0.010</b>	0.754 ± 0.018	0.752 ± 0.020	0.779 ± 0.013	0.814 ± 0.015	0.844 ± 0.018

Table 6: **Downstream task evaluation with data augmentation.** Downstream tasks are evaluated under the training scenarios of *Train on Synthetic and Real, Test on Real (TSRTR)*. The predictive tasks and evaluation metrics are in accordance with Table 5. The upper arrow (↑) indicates that the AUROC value under *TSRTR* is higher than *TRTR* in Table 5 for the corresponding task, while the bold arrow (⤴) indicates that the value is significantly improved using *t*-test ( $p \leq 0.05$ ).

(a) Performance of the downstream LSTM-based classifier under *TSRTR* with data augmentation ratio  $\alpha$ . All data from sub-train data  $A'_{Tr}$  concatenated with  $\alpha$  of the synthetic data  $B$  (augmentation ratio  $\alpha = 10\%$ ,  $25\%$  or  $50\%$ ) is used as the training set.

Dataset	Treatments	EHR-M-GAN			EHR-M-GAN <sub>cond</sub>		
		$\alpha = 10\%$	$\alpha = 25\%$	$\alpha = 50\%$	$\alpha = 10\%$	$\alpha = 25\%$	$\alpha = 50\%$
MIMIC-III	Vent.	0.828 ± 0.013	0.877 ± 0.014	0.912 ± 0.015 (⤴)	0.845 ± 0.022	0.896 ± 0.013 (↑)	<b>0.923 ± 0.018 (⤴)</b>
	Vaso.	0.816 ± 0.015	0.834 ± 0.023	0.859 ± 0.013 (⤴)	0.848 ± 0.012 (↑)	0.876 ± 0.017 (⤴)	<b>0.896 ± 0.015 (⤴)</b>
eICU	Vent.	0.858 ± 0.008	0.862 ± 0.012	0.873 ± 0.014 (↑)	0.865 ± 0.009	0.879 ± 0.014 (↑)	<b>0.883 ± 0.016 (⤴)</b>
	Vaso.	0.798 ± 0.015	0.805 ± 0.020	0.821 ± 0.028 (↑)	0.813 ± 0.016 (↑)	0.834 ± 0.019 (⤴)	<b>0.839 ± 0.014 (⤴)</b>
HiRID	Vent.	0.871 ± 0.025 (↑)	0.882 ± 0.021 (↑)	0.913 ± 0.019 (⤴)	0.894 ± 0.015 (⤴)	0.906 ± 0.018 (⤴)	<b>0.923 ± 0.021 (⤴)</b>
	Vaso.	0.850 ± 0.016	0.874 ± 0.022	0.894 ± 0.018 (⤴)	0.883 ± 0.017 (↑)	0.908 ± 0.024 (↑)	<b>0.913 ± 0.019 (⤴)</b>

(b) Performance of the downstream LSTM-based classifier under *TSRTR* with data augmentation ratio  $\beta$ . All data from synthetic data  $B$  concatenated with  $\beta$  of the sub-train data  $A'_{Tr}$  (augmentation ratio  $\beta = 10\%$ ,  $25\%$  or  $50\%$ ) is used as the training set.

Dataset	Treatments	EHR-M-GAN			EHR-M-GAN <sub>cond</sub>		
		$\beta = 10\%$	$\beta = 25\%$	$\beta = 50\%$	$\beta = 10\%$	$\beta = 25\%$	$\beta = 50\%$
MIMIC-III	Vent.	0.757 ± 0.016	0.824 ± 0.010	0.885 ± 0.009	0.847 ± 0.017	0.903 ± 0.014 (↑)	<b>0.915 ± 0.009 (⤴)</b>
	Vaso.	0.786 ± 0.019	0.810 ± 0.020	0.849 ± 0.017 (↑)	0.823 ± 0.014	0.851 ± 0.019 (↑)	<b>0.873 ± 0.017 (⤴)</b>
eICU	Vent.	0.761 ± 0.011	0.822 ± 0.012	0.870 ± 0.019 (↑)	0.816 ± 0.016	0.845 ± 0.018	<b>0.872 ± 0.020 (⤴)</b>
	Vaso.	0.742 ± 0.014	0.797 ± 0.013	0.846 ± 0.018 (⤴)	0.785 ± 0.022	0.819 ± 0.021 (↑)	<b>0.834 ± 0.013 (⤴)</b>
HiRID	Vent.	0.856 ± 0.012	0.879 ± 0.019 (↑)	0.895 ± 0.021 (⤴)	0.874 ± 0.016 (↑)	0.896 ± 0.018 (⤴)	<b>0.904 ± 0.012 (⤴)</b>
	Vaso.	0.826 ± 0.024	0.859 ± 0.013	0.893 ± 0.018 (↑)	0.865 ± 0.025	0.897 ± 0.021 (⤴)	<b>0.914 ± 0.018 (⤴)</b>

Compared with the baseline models, the proposed EHR-M-GAN shows improved performance in *TSTR*, as it can model the distribution of mixed-type EHRs more accurately, while preserving the temporal correlations in the heterogeneous timeseries through the dependency learning components. The results indicate that interdependency between the mixed-type EHRs is weakly captured by GAN<sub>VAE</sub>, as the two streams of inputs are trained in parallel and separately. GAN<sub>Unified</sub> attempts to capture the temporal correlations of mixed-type EHRs through jointly modeling their underlying distribution in a unified network. However, its unified architecture limits the model’s capacity to learn the marginal distribution of each data type, the resulted quality of the synthetic EHRs is impaired and so is its performance in *TSTR*.

**Data augmentation approach (with ratio  $\alpha$ ).** The results in Table 6a demonstrate that classifiers boosted by EHR-M-GAN can consistently outperform TRTR (see Table 5) at the augmentation ratio of 50%. In comparison, only 25% of augmentation ratio is needed to achieve improved results for EHR-M-GAN<sub>cond</sub>. For example, the classifier trained on MIMIC-III to predict the status of *Vasopressor* with augmentation ratio  $\alpha$  set as 50%, significantly increase the AUROC by 6% when compared to the classifier trained using only the real data (EHR-M-GAN<sub>cond</sub> vs. TRTR: 0.896 to 0.841,  $p < 0.05$ ). Our experiment results have demonstrated that the proposed models can be used for data augmentation to overcome the issue of data scarcity and subsequently improve the classifiers’ performance.

**Data augmentation approach (with ratio  $\beta$ ).** On the other hand, as shown in Table 6b, by augmenting with the synthetic data provided by EHR-M-GAN, only approximately 50% of the real data is required to keep the classification AUROCs on par with, or even significantly better than fully exploiting the real data under *TRTR*. For EHR-M-GAN<sub>cond</sub>, the ratio needed for real data to maintain the comparable predictive power is further reduced to 25%, which equivalently indicates a 75% reduction of sample size required in real data collection. Overall, results presented in Table 6b demonstrate that by exploiting only a limited ratio of the real data, EHR-M-GAN and EHR-M-GAN<sub>cond</sub> can robustly maintain the level of prediction performance, therefore alleviating the necessity for acquiring clinical data at scale.

#### 4.7 PRIVACY RISK EVALUATION

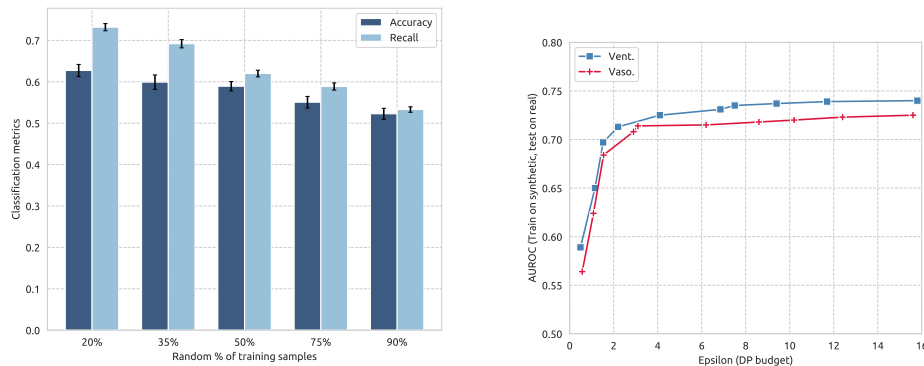
Patient privacy is a major concern with regards to sharing electronic health records in any means. Generative models can overcome the explicit one-to-one mapping towards the underlying original data in contrast to data anonymisation. However, GAN could potentially raise privacy concerns of information leakage if they simply “memorise” the training data, or synthesize samples nearly identical to the real samples (often due to mode collapse). In that case, sensitive medical information (e.g. national insurance number) belonging to a specific patient used in training GANs can be retrieved during the generative stage, thus posing challenges for preserving privacy in downstream applications.

In this section, we first quantify the vulnerability of EHR-M-GAN to adversary’s **membership inference attacks**, also known as presence disclosure Hayes et al. (2019); Chen, Yu, Zhang and Fritz (2020). The threat model is implemented under the membership inference for GANs in the *black-box settings* Hayes et al. (2019). The attacker is assumed to possess complete knowledge of all the patient records set  $P$ , where a subset from  $P$  further is used to train GANs. During the experiment, the number of samples in the subset for training EHR-M-GAN are varied to investigate the impact of the availability of training data on the success of the attacker (see Figure 7a). By observing the synthetic patient records from EHR-M-GAN, the adversary’s goal is to determine whether a single known record  $x$  in the patient record set  $P$  is from the data used in training EHR-M-GAN. If EHR-M-GAN simply “memorises” the training data and can only generate synthetic samples (nearly) identical to the real samples, it would be straightforward for the adversary to identify samples that are used as training data. Determined by whether the attacker can correctly infer a given record is *in* or *not in* GAN’s training, the accuracy and recall can be calculated.

As shown in Figure 7a, when 90% of the training data is used for developing EHR-M-GAN, the attacker had a recall of 0.533 and accuracy of 0.527 to recover which training data are considered. This is eminently close to flipping a coin with random guess (i.e., 0.5), indicating EHR-M-GAN is sufficiently robust against the membership inference attack. In other words, patient samples used in EHR-M-GAN’s training are not recoverable by the threat model. On the other hand, as the percentage of the training data reduces, both accuracy and recall for membership inference attacks rise. An accuracy of 0.624 and recall of 0.732 are reached with 20% of training data. This offers a guideline for future application in developing GANs that incorporating more training data can make the generator less susceptible to such attack. This is also consistent with the conclusion derived from the experiment on membership inference attacks in the prior research Lin et al. (2020).

The concept differential privacy (DP) Dwork (2008), which is a rigorous mathematical definition of privacy, has emerged to be the prevailing notion in the context of statistically analyzing data privacy. The  $(\epsilon, \delta)$ -differential privacy is guaranteed for model  $M$ , if given any pair of *adjacent* datasets  $D$  and  $D'$  (differing on a single patient record), it holds:  $P[\mathcal{M}(D) \in S] \leq e^\epsilon P[\mathcal{M}(D') \in S] + \delta$ . In our case,  $\mathcal{M}(\cdot)$  is the GAN model trained based on  $D$  or  $D'$ , and  $S$  is the subset of any possible outcomes of the generative process. By perturbing the underlying data distribution, DP bounds the maximum variations of the algorithm when *any* single individual is included or excluded from the dataset. In practice, recent works on developing differentially private deep learning models have benefited from differential private stochastic gradient descent (DP-SGD) algorithm. DP-SGD operates DP by gradient clipping and noise adding during SGD, thereby ensuring that the impact of single record in the training dataset on algorithm parameters is limited within DP’s extend. In this section,  $(\epsilon, \delta)$ -differential privacy is implemented in EHR-M-GAN using TensorFlow Privacy<sup>4</sup>.

<sup>4</sup><https://github.com/tensorflow/privacy>



(a) Membership inference attack.

(b) Differential privacy.

**Figure 7: Privacy risk evaluation of EHR-M-GAN on MIMIC-III dataset. a. Membership inference attack.** Membership inference attack against EHR-M-GAN vs. the percentage of the training data. Accuracy and recall are used to evaluate the success rate of such attacks. Lower accuracy or recall indicates less privacy information disclosed by the attacker from the generative model (0.5 can be seen as the *random guess* baseline where strong privacy guarantees are provided by GANs). Recall indicates the ratio of samples that are successfully claimed by the attacker among all the real data that are used in training GAN models. **b. Differential privacy.** Performance of medical intervention prediction tasks, under various differential privacy (DP) budgets, measured by Macro-AUROC.

We then perform the same downstream tasks on medical intervention prediction using synthetic data generated from DP-guaranteed EHR-M-GAN, and compare its performance with *TSTR* (as shown in Table 5).

Figure 7b shows the *TSTR* performance of EHR-M-GAN under differential privacy guarantee with varying budgets  $\epsilon$  ( $\delta$  fixed at  $\leq 0.001$ ). The value  $\epsilon$  determines how strict the privacy is, where the smaller value indicates a stronger privacy restriction. As suggested in Figure 7b, the performance of the downstream tasks operated based on the synthetic data generated by EHR-M-GAN improves as the DP budget relaxes ( $\epsilon$  increases). We observe that the AUROC of DP-bounded EHR-M-GAN can maintain at an acceptable level even under strict privacy setting. For example, the AUROC for predicting the treatment of *Vasopressor* can maintain at 0.714 (AUROC = 0.725 under *TRTR*) even when the  $\epsilon$  decrease to 4, which is an empirically reasonable value for implementing DP in practice Differential Privacy Team (2017). Future work that focuses on privacy-preserving GAN under DP-guarantee is expected, where the fidelity of the synthetic data can be restored without compromising its privacy.

## 5 DISCUSSION AND CONCLUSIONS

In this study, we propose a generative adversarial network entitled EHR-M-GAN, aiming at mitigating the challenge of synthesizing longitudinal EHR with mixed data types. A comprehensive list of evaluation metrics is introduced for the systematic assessment, in terms of the fidelity, correlation, utility, and privacy of the synthesis model. First, both EHR-M-GAN and its conditional version, EHR-M-GAN<sub>cond</sub>, demonstrate consistent improvements against the state-of-the-art benchmark GANs in synthesizing timeseries data with high-fidelity. This indicates that the distributional characteristics of the EHR timeseries can be well-preserved in the synthetic data provided by EHR-M-GAN, therefore ensuring its usability during clinical data sharing. Second, as opposed to previous models which were confined to synthesizing only one specific type of data, EHR-M-GAN can produce mixed-type timeseries and successfully capture the temporal dynamics and correlation between features. By accurately reconstructing the interdependencies and complex clinical relationships between features, downstream studies such as association analysis and outcome prediction can be supported. Notably, the proposed models also outperform the GAN variants that allow mixed-type inputs in the ablation study, indicating that the components in EHR-M-GAN are effective in synthesizing mixed-type

---

timeseries with high fidelity, while successfully reconstructing the interdependencies between them. Then, during downstream task evaluation, given the prediction of medical interventions in fast-paced critical care environments as an exemplar, the results demonstrate the broad applicability of our model in developing ML algorithm-based decision support tools by data augmentation. Lastly, the assessment of privacy risks further demonstrates the synthetic data provided by EHR-M-GAN can preserve the sensitive information in patient records while maintaining an acceptable level of data utility.

The results in our study have several notable implications with respect to the synthesis of EHR data. First, as the proposed model can be used to provide synthetic longitudinal EHRs for various data types while preserving their underlying correlations, it is now feasible to use the synthesized data to improve the performance of ML models for downstream applications such as the prediction of next intervention, or understanding the disease dynamics and patient phenotyping, based on both the continuous and discrete components of EHR timeseries Alaa and van der Schaar (2019); Lee and Van Der Schaar (2020). Second, the experimental results indicate that the quality of the synthetic EHR data can be improved by the integration of mixed-type information, in contrast to the benchmarks that utilize single-type data for learning. This also enables us to mimic how information is presented in clinical practices. Furthermore, we can generate condition/outcome-specific patient trajectories along with corresponding interventions, to facilitate clinical prediction and decision-making. Third, though facing the privacy-utility tradeoffs, the synthetic EHRs data provided by the proposed model leads to negligible privacy risks under the membership inference attacks. This paves the way for a series of applications in clinical research, including but not limited to, enabling the development of ML models by accessing the synthetic data, overcoming the paucity of medical data and improving the robustness of ML algorithms through data augmentation.

Due to the heterogeneous nature of EHR data, besides the ICU setting in our empirical evaluation, there are needs for synthesizing mixed-type EHR timeseries in various clinical scenarios. For example, patients' encounters in hospitals are documented as structured EHRs recorded in the temporal order. Each visit is typically associated with the corresponding medical events presented in the form of discrete-valued ICD codes Zhang et al. (2021), and continuous-valued measurements. These mixed-type EHR timeseries capture a patient's health status and better align with clinical decision-making process than those using the single-type data alone. Therefore, developing GANs targeting mixed-type EHRs generation have the potential to pave the way for complex deep-learning systems that are capable of integrating information from various sources. However, it is worth noting that the validation of our proposed model is based on critical care settings with limited feature dimensions, can only serve as a proof of concept. When extending the proposed model to other clinical settings, such as synthesizing ICD codes with hundreds or thousands of feature dimensions Zhang et al. (2021), the scalability and utility of our proposed model when dealing with the enlarged, sparse feature space needs further investigation.

There are limitations in the current work. First, data curation strategies on clinical timeseries, including truncating, smoothing and imputation, are applied before the EHR timeseries are used for the training of generative models. As during the data preprocessing, we first extract the timeseries with a fixed duration (i.e., 24 hours before the ICU clinical endpoints), and then hourly aggregate patients' physiological and intervention signals based on their mean statistics, followed by completing the missing value in the timeseries through the "Simple Imputation" approach Che et al. (2018). Although these preprocessing steps are commonly used in clinical research under the critical care settings Wang et al. (2020), the proposed model cannot model the irregular time intervals between signals nor missing values within the timeseries. However, dealing with irregularity of the timestamps when synthesizing clinical events in EHRs could be useful for predicting outcomes that are time-aware in the downstream tasks Zhang et al. (2021). Modeling such time intervals could be non-trivial as the determinative perspectives sometimes go beyond the scope of inferring patients physiological status such as resource allocations within hospitals. Also, synthesizing timeseries while incorporating the missing values could also be beneficial in the real-world application scenarios. As ML models are sometimes sensitive to the data missingness, imputing the incomplete data in EHRs using generative approaches could improve the performance of ML models, and has become an area of active research Yoon, Jordon and Schaar (2018). Furthermore, as evaluations are performed based on clinical timeseries with a fixed length, no comparisons are made between the model's scalability when dealing with timeseries with varying lengths. Recent studies have found the quality of the synthetic longitudinal data degenerates over time, also called as the "drift problem" Zhang et al. (2022). Such

---

problems when dealing with long sequences should be recognised and mitigated with techniques such as conditional fuzzing and regularization methods Zhang et al. (2022), in both the generation and evaluation steps.

The evaluation of GANs is still a challenging task. Recent findings have suggested that systematical assessment for EHR synthesizers is critical before their applications in different use cases Yan et al. (2022). In this study, a comprehensive evaluation list is provided with regards to the fidelity, correlation, utility and privacy of the synthesis models. It is also worth noting that evaluation metrics should be properly chosen and implemented based on the purpose of the task, otherwise may lead to biased results. For example, recent findings Zhang et al. (2022) have reported that the traditional implementation of the discriminative score which trains the critic using the randomly initialised parameters, though widely used Yoon et al. (2019), may lead to unreliable results. Improvement has been made to this evaluation metric for a more robust assessment, where the parameters of the trained generative model can be used for the critic’s initialization.

Finally, the conditional aspect of our model is currently limited as it can not generate patient-specific EHRs conditioning on information at a more granular level. Even though the proposed conditional GANs can synthesize a subgroup of patients with target outcomes or statuses that clinicians specify, it is still limited in incorporating personalised information during the conditional generation. Future work for developing GANs in healthcare data can be extended to patient-level EHRs generation, such as synthesizing counterfactual information of a target patient for treatment effect estimation Yoon, Jordon and Van Der Schaar (2018); Qian et al. (2021). Ultimately, by constructing the “synthetic twin” of patients using GANs, the synthesis tool can become more generalisable for precision medicine and support the clinical decision making in delivering personalized healthcare.

Synthetic data provides an alternative to sharing real patient data while preserving patient privacy. Results in our study demonstrate that the proposed EHR-M-GAN and EHR-M-GAN<sub>cond</sub> can generate realistic longitudinal EHR timeseries with mixed data types. By providing synthetic EHR data with higher fidelity and more variety, the proposed model can therefore enable faster development in AI-driven clinical tools with increased robustness and adaptability. In addition to the improved performance against the existing state-of-the-art benchmark models, augmentation provided by synthetic data during training boosts the predictive performance in downstream clinical tasks. EHR-M-GAN can help eliminate the barriers to data acquisition for healthcare studies, therefore overcoming the challenges posed by the paucity of medical data available and approved for research use. Despite the novelty of this study in filling the research gap for synthesizing longitudinal EHRs in mixed-type settings, we acknowledge that there is still a gap between the real EHRs data and its synthetic counterparts produced by current generative methods. Therefore developing advanced EHR synthesizers especially in mixed-type settings still requires active research in the future study.

## REFERENCES

- Alaa, A. M. and van der Schaar, M. (2019), ‘Attentive state-space modeling of disease progression’, *Advances in neural information processing systems* **32**.
- Artzi, N. S., Shilo, S., Hadar, E., Rossman, H., Barbash-Hazan, S., Ben-Haroush, A., Balicer, R. D., Feldman, B., Wiznitzer, A. and Segal, E. (2020), ‘Prediction of gestational diabetes based on nationwide electronic health records’, *Nature medicine* **26**(1), 71–76.
- Baowaly, M. K., Lin, C.-C., Liu, C.-L. and Chen, K.-T. (2019), ‘Synthesizing electronic health records using improved generative adversarial networks’, *Journal of the American Medical Informatics Association* **26**(3), 228–241.
- Benedetti, J. d., Oues, N., Wang, Z., Myles, P. and Tucker, A. (2020), Practical lessons from generating synthetic healthcare data with bayesian networks, in ‘Joint European Conference on Machine Learning and Knowledge Discovery in Databases’, Springer, pp. 38–47.
- Borji, A. (2021), ‘Pros and cons of gan evaluation measures: New developments’, *arXiv preprint arXiv:2103.09396*.
- Che, Z., Purushotham, S., Cho, K., Sontag, D. and Liu, Y. (2018), ‘Recurrent neural networks for multivariate time series with missing values’, *Scientific reports* **8**(1), 1–12.

- 
- Chen, D., Yu, N., Zhang, Y. and Fritz, M. (2020), Gan-leaks: A taxonomy of membership inference attacks against generative models, in ‘Proceedings of the 2020 ACM SIGSAC conference on computer and communications security’, pp. 343–362.
- Chen, R. J., Lu, M. Y., Chen, T. Y., Williamson, D. F. and Mahmood, F. (2021), ‘Synthetic data in machine learning for medicine and healthcare’, *Nature Biomedical Engineering* pp. 1–5.
- Chen, T., Kornblith, S., Norouzi, M. and Hinton, G. (2020), A simple framework for contrastive learning of visual representations, in ‘International conference on machine learning’, PMLR, pp. 1597–1607.
- Choi, E., Biswal, S., Malin, B., Duke, J., Stewart, W. F. and Sun, J. (2017), Generating multi-label discrete patient records using generative adversarial networks, in ‘Machine learning for healthcare conference’, PMLR, pp. 286–305.
- Differential Privacy Team, A. (2017), Learning with privacy at scale.
- Dinov, I. D. (2016), ‘Methodological challenges and analytic opportunities for modeling and interpreting big healthcare data’, *Gigascience* **5**(1), s13742–016.
- Dwork, C. (2006), Differential privacy, in ‘International Colloquium on Automata, Languages, and Programming’, Springer, pp. 1–12.
- Dwork, C. (2008), Differential privacy: A survey of results, in ‘International conference on theory and applications of models of computation’, Springer, pp. 1–19.
- El Emam, K., Mosquera, L., Jonker, E. and Sood, H. (2021), ‘Evaluating the utility of synthetic covid-19 case data’, *JAMIA open* **4**(1), ooab012.
- Esteban, C., Hyland, S. L. and Rätsch, G. (2017), ‘Real-valued (medical) time series generation with recurrent conditional gans’, *arXiv preprint arXiv:1706.02633* .
- Esteva, A., Chou, K., Yeung, S., Naik, N., Madani, A., Mottaghi, A., Liu, Y., Topol, E., Dean, J. and Socher, R. (2021), ‘Deep learning-enabled medical computer vision’, *NPJ digital medicine* **4**(1), 1–9.
- Frid-Adar, M., Diamant, I., Klang, E., Amitai, M., Goldberger, J. and Greenspan, H. (2018), ‘Gan-based synthetic medical image augmentation for increased cnn performance in liver lesion classification’, *Neurocomputing* **321**, 321–331.
- Futoma, J., Simons, M., Panch, T., Doshi-Velez, F. and Celi, L. A. (2020), ‘The myth of generalisability in clinical research and machine learning in health care’, *The Lancet Digital Health* **2**(9), e489–e492.
- Ghassemi, M., Wu, M., Hughes, M. C., Szolovits, P. and Doshi-Velez, F. (2017), ‘Predicting intervention onset in the icu with switching state space models’, *AMIA Summits on Translational Science Proceedings* **2017**, 82.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. and Bengio, Y. (2014), ‘Generative adversarial nets’, *Advances in neural information processing systems* **27**.
- Hayes, J., Melis, L., Danezis, G. and De Cristofaro, E. (2019), Logan: Membership inference attacks against generative models, in ‘Proceedings on Privacy Enhancing Technologies (PoPETs)’, Vol. 2019, De Gruyter, pp. 133–152.
- Hjelm, R. D., Jacob, A. P., Che, T., Trischler, A., Cho, K. and Bengio, Y. (2017), ‘Boundary-seeking generative adversarial networks’, *arXiv preprint arXiv:1702.08431* .
- Johnson, A. E., Pollard, T. J., Shen, L., Li-Wei, H. L., Feng, M., Ghassemi, M., Moody, B., Szolovits, P., Celi, L. A. and Mark, R. G. (2016), ‘Mimic-iii, a freely accessible critical care database’, *Scientific data* **3**(1), 1–9.



- 
- Jordon, J., Wilson, A. and van der Schaar, M. (2020), ‘Synthetic data: Opening the data floodgates to enable faster, more directed development of machine learning methods’, *arXiv preprint arXiv:2012.04580* .
- Jordon, J., Yoon, J. and Van Der Schaar, M. (2018), Pate-gan: Generating synthetic data with differential privacy guarantees, in ‘International conference on learning representations’.
- Kearney, V., Chan, J. W., Wang, T., Perry, A., Descovich, M., Morin, O., Yom, S. S. and Solberg, T. D. (2020), ‘Dosegan: a generative adversarial network for synthetic dose prediction using attention-gated discrimination and generation’, *Scientific reports* **10**(1), 1–8.
- Kim, J., Neumann, L., Paul, P., Day, M. E., Aratow, M., Bell, D. S., Doctor, J. N., Hinske, L. C., Jiang, X., Kim, K. K. et al. (2021), ‘Privacy-protecting, reliable response data discovery using covid-19 patient observations’, *Journal of the American Medical Informatics Association* **28**(8), 1765–1776.
- Kiyasseh, D., Tadesse, G. A., Thwaites, L., Zhu, T., Clifton, D. et al. (2020), ‘Plethaugment: Gan-based ppg augmentation for medical diagnosis in low-resource settings’, *IEEE journal of biomedical and health informatics* **24**(11), 3226–3235.
- Kiyasseh, D., Zhu, T. and Clifton, D. A. (2021), Clocs: Contrastive learning of cardiac signals across space, time, and patients, in ‘International Conference on Machine Learning’, PMLR, pp. 5606–5615.
- Lee, C. and Van Der Schaar, M. (2020), Temporal phenotyping using deep predictive clustering of disease progression, in ‘International Conference on Machine Learning’, PMLR, pp. 5767–5777.
- Lee, D., Yu, H., Jiang, X., Rogith, D., Gudala, M., Tejani, M., Zhang, Q. and Xiong, L. (2020), ‘Generating sequential electronic health records using dual adversarial autoencoder’, *Journal of the American Medical Informatics Association* **27**(9), 1411–1419.
- Lin, Z., Jain, A., Wang, C., Fanti, G. and Sekar, V. (2020), Using gans for sharing networked time series data: Challenges, initial promise, and open questions, in ‘Proceedings of the ACM Internet Measurement Conference’, pp. 464–483.
- Liu, M.-Y., Breuel, T. and Kautz, J. (2017), Unsupervised image-to-image translation networks, in ‘Advances in neural information processing systems’, pp. 700–708.
- Liu, M.-Y. and Tuzel, O. (2016), ‘Coupled generative adversarial networks’, *Advances in neural information processing systems* **29**, 469–477.
- Liu, X., Zhang, F., Hou, Z., Mian, L., Wang, Z., Zhang, J. and Tang, J. (2021), ‘Self-supervised learning: Generative or contrastive’, *IEEE Transactions on Knowledge and Data Engineering* .
- Marouf, M., Machart, P., Bansal, V., Kilian, C., Magruder, D. S., Krebs, C. F. and Bonn, S. (2020), ‘Realistic in silico generation and augmentation of single-cell rna-seq data using generative adversarial networks’, *Nature communications* **11**(1), 1–12.
- Menger, V., Spruit, M., Van Est, R., Nap, E. and Scheepers, F. (2019), ‘Machine learning approach to inpatient violence risk assessment using routinely collected clinical notes in electronic health records’, *JAMA network open* **2**(7), e196709–e196709.
- Miotto, R., Wang, F., Wang, S., Jiang, X. and Dudley, J. T. (2018), ‘Deep learning for healthcare: review, opportunities and challenges’, *Briefings in bioinformatics* **19**(6), 1236–1246.
- Mirza, M. and Osindero, S. (2014), ‘Conditional generative adversarial nets’, *arXiv preprint arXiv:1411.1784* .
- Mogren, O. (2016), ‘C-rnn-gan: Continuous recurrent neural networks with adversarial training’, *arXiv preprint arXiv:1611.09904* .
- N3C. *Synthetic Data Workstream* (n.d.), [https://covid.cd2h.org/N3C\\_synthetic\\_data](https://covid.cd2h.org/N3C_synthetic_data). Accessed: 2021-12-02.

- 
- Pollard, T. J., Johnson, A. E., Raffa, J. D., Celi, L. A., Mark, R. G. and Badawi, O. (2018), ‘The eicu collaborative research database, a freely available multi-center database for critical care research’, *Scientific data* **5**(1), 1–13.
- Qian, Z., Zhang, Y., Bica, I., Wood, A. and van der Schaar, M. (2021), ‘Synctwin: Treatment effect estimation with longitudinal outcomes’, *Advances in Neural Information Processing Systems* **34**.
- Rajkomar, A., Dean, J. and Kohane, I. (2019), ‘Machine learning in medicine’, *New England Journal of Medicine* **380**(14), 1347–1358.
- Raket, L. L., Jaskolowski, J., Kinon, B. J., Brasen, J. C., Jönsson, L., Wehnert, A. and Fusar-Poli, P. (2020), ‘Dynamic electronic health record detection (detect) of individuals at risk of a first episode of psychosis: a case-control development and validation study’, *The Lancet Digital Health* **2**(5), e229–e239.
- Shokri, R., Stronati, M., Song, C. and Shmatikov, V. (2017), Membership inference attacks against machine learning models, in ‘2017 IEEE Symposium on Security and Privacy (SP)’, IEEE, pp. 3–18.
- Simon, G. E., Shortreed, S. M., Coley, R. Y., Penfold, R. B., Rossom, R. C., Waitzfelder, B. E., Sanchez, K. and Lynch, F. L. (2019), ‘Assessing and minimizing re-identification risk in research data derived from health care records’, *eGEMs* **7**(1).
- Suresh, H., Hunt, N., Johnson, A., Celi, L. A., Szolovits, P. and Ghassemi, M. (2017), Clinical intervention prediction and understanding with deep neural networks, in ‘Machine Learning for Healthcare Conference’, PMLR, pp. 322–337.
- Sutherland, D. J., Tung, H.-Y., Strathmann, H., De, S., Ramdas, A., Smola, A. and Gretton, A. (2016), ‘Generative models and model criticism via optimized maximum mean discrepancy’, *arXiv preprint arXiv:1611.04488*.
- Synthetic data at CPRD, howpublished = <https://www.cprd.com/content/synthetic-data> (n.d.). Accessed: 2021-12-02.*
- Tao, Y., McKenna, R., Hay, M., Machanavajjhala, A. and Miklau, G. (2021), ‘Benchmarking differentially private synthetic data generation algorithms’, *arXiv preprint arXiv:2112.09238*.
- Tucker, A., Wang, Z., Rotalinti, Y. and Myles, P. (2020), ‘Generating high-fidelity synthetic patient data for assessing machine learning healthcare software’, *NPJ digital medicine* **3**(1), 1–13.
- Wan, Z., Zhang, B., Chen, D., Zhang, P., Chen, D., Liao, J. and Wen, F. (2020), ‘Old photo restoration via deep latent space translation’, *arXiv preprint arXiv:2009.07047*.
- Wang, L., Zhang, W. and He, X. (2019), Continuous patient-centric sequence generation via sequentially coupled adversarial learning, in ‘International Conference on Database Systems for Advanced Applications’, Springer, pp. 36–52.
- Wang, S., McDermott, M. B., Chauhan, G., Ghassemi, M., Hughes, M. C. and Naumann, T. (2020), Mimiextract: A data extraction, preprocessing, and representation pipeline for mimiciii, in ‘Proceedings of the ACM Conference on Health, Inference, and Learning’, pp. 222–235.
- Watson, D. S., Krutzinna, J., Bruce, I. N., Griffiths, C. E., McInnes, I. B., Barnes, M. R. and Floridi, L. (2019), ‘Clinical applications of machine learning algorithms: beyond the black box’, *Bmj* **364**.
- Wilkinson, J., Arnold, K. F., Murray, E. J., van Smeden, M., Carr, K., Sippy, R., de Kamps, M., Beam, A., Konigorski, S., Lippert, C. et al. (2020), ‘Time to reality check the promises of machine learning-powered precision medicine’, *The Lancet Digital Health*.
- Wirth, F. N., Meurers, T., Johns, M. and Prasser, F. (2021), ‘Privacy-preserving data sharing infrastructures for medical research: systematization and comparison’, *BMC Medical Informatics and Decision Making* **21**(1), 1–13.
- Wu, M., Ghassemi, M., Feng, M., Celi, L. A., Szolovits, P. and Doshi-Velez, F. (2017), ‘Understanding vasopressor intervention and weaning: risk prediction in a public heterogeneous clinical time series database’, *Journal of the American Medical Informatics Association* **24**(3), 488–495.

- 
- Yan, C., Yan, Y., Wan, Z., Zhang, Z., Omberg, L., Guinney, J., Mooney, S. D. and Malin, B. A. (2022), ‘A multifaceted benchmarking of synthetic electronic health record generation models’, *Nature communications* **13**(1), 1–18.
- Yang, Q., Yan, P., Zhang, Y., Yu, H., Shi, Y., Mou, X., Kalra, M. K., Zhang, Y., Sun, L. and Wang, G. (2018), ‘Low-dose ct image denoising using a generative adversarial network with wasserstein distance and perceptual loss’, *IEEE transactions on medical imaging* **37**(6), 1348–1357.
- Yèche, H., Kuznetsova, R., Zimmermann, M., Hüser, M., Lyu, X., Faltys, M. and Ratsch, G. (2021), ‘Hirid-icu-benchmark—a comprehensive machine learning benchmark on high-resolution icu data.’.
- Yoon, J., Jarrett, D. and Van der Schaar, M. (2019), ‘Time-series generative adversarial networks’.
- Yoon, J., Jordon, J. and Schaar, M. (2018), Gain: Missing data imputation using generative adversarial nets, in ‘International conference on machine learning’, PMLR, pp. 5689–5698.
- Yoon, J., Jordon, J. and Van Der Schaar, M. (2018), Ganite: Estimation of individualized treatment effects using generative adversarial nets, in ‘International Conference on Learning Representations’.
- Yu, C., Liu, J. and Zhao, H. (2019), ‘Inverse reinforcement learning for intelligent mechanical ventilation and sedative dosing in intensive care units’, *BMC medical informatics and decision making* **19**(2), 111–120.
- Yu, L., Zhang, W., Wang, J. and Yu, Y. (2017), Seqgan: Sequence generative adversarial nets with policy gradient, in ‘Proceedings of the AAAI conference on artificial intelligence’, Vol. 31.
- Zhang, J., Cormode, G., Procopiuc, C. M., Srivastava, D. and Xiao, X. (2017), ‘Privbayes: Private data release via bayesian networks’, *ACM Transactions on Database Systems (TODS)* **42**(4), 1–41.
- Zhang, Z., Yan, C., Lasko, T. A., Sun, J. and Malin, B. A. (2021), ‘Synteg: a framework for temporal structured electronic health data simulation’, *Journal of the American Medical Informatics Association* **28**(3), 596–604.
- Zhang, Z., Yan, C. and Malin, B. A. (2022), ‘Keeping synthetic patients on track: feedback mechanisms to mitigate performance drift in longitudinal health data simulation’, *Journal of the American Medical Informatics Association* **29**(11), 1890–1898.

# Appendix

## 1. METHODOLOGY.

### A. Related work.

Generative adversarial networks (GANs) have been used in EHR data synthesis, which can augment limited clinical data or even replace sensitive patient information. As longitudinal EHR data can capture patients' status over time, generative approaches for EHR data synthesis in the previous literature have been extended from static data to clinical timeseries generation. EHRs consist of a set of heterogeneous data types, such as continuous-valued and discrete-valued features. When generating continuous-valued timeseries such as heart rate and respiratory rate in the critical care database, models such as C-RNN-GAN [1], R(C)GAN [2] and TimeGAN [3] can be adopted. In order to model the temporal dynamics in the real-valued timeseries, recurrent neural networks (RNNs) such as long short-term memory (LSTM) are used as the generator and discriminator in their architectures. For synthesizing discrete-valued timeseries data such as diagnostic ICD-codes, GANs variants such as SynTEG [4], LS-EHR [5], and DualAAE [6] models are proposed. For example, SynTEG generates time-stamped clinical events across patients' multiple visits. Its amended version — LS-EHR [5] model enhances the longitudinal EHR data synthesis by overcoming the performance drift through feedback mechanisms (including condition fuzzing, regularization and rejection sampling).

As EHRs is an amalgamation of heterogeneous data types, previous work has demonstrated the importance of synthesizing mixed-type EHRs for various clinical applications. Several models have been proposed to generate static EHRs of mixed data types, such as discrete-valued medical concepts and continuous-valued measurements [7–9]. However, for synthesizing clinical timeseries, most of the proposed models have been capable of synthesizing only a single data type (either continuous or discrete-valued timeseries separately). Consequently, previous work has tended to ignore the inter-dependencies among different data types, as shown in Fig. S1. In contrast, our proposed model can simultaneously generate both continuous-valued and discrete-valued timeseries, while capturing the inter-dependencies between the mixed-type data.

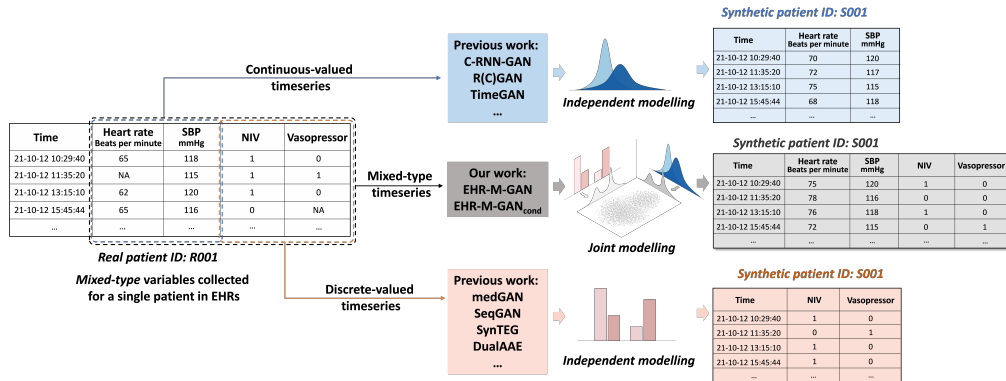


Fig. S1. A comparison of the proposed model with existing generative models.

### B. Implementation of GANs using LSTMs

GAN consists of two networks that are adversarially trained to compete against each other. The recurrent neural networks (RNNs) are instantiated considering simulating the temporal structure for generating the sequential data. As shown in Fig. 1 (b. Network architecture) in the main article, the generator  $G$  accepts  $v_{1:T} \in \mathcal{T} \times \mathcal{V}$  as the input, which is a sequence of length  $T$  sampled independently from a prior distribution [2], such as Gaussian distribution or uniform distribution. In this study, uniform distribution on the unit interval is chosen as the prior for sampling the random noise. Then  $G$  is optimized to approximate the distribution of true data,  $p_x$ ,

by generating samples  $\hat{\mathbf{x}}_{1:T}$  that are hard for the discriminator to distinguish from. Meanwhile, the discriminator  $D$  is optimized to distinguish real samples  $\mathbf{x}_{1:T}$  from synthetic samples  $\hat{\mathbf{x}}_{1:T}$ . Overall, the training of GAN is a minmax game with the following objective function:

$$\min_G \max_D V_{\text{GAN}} = \mathbb{E}_{\mathbf{x} \sim p_x} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{v} \sim p_v} [\log(1 - D(G(\mathbf{v})))] \quad (\text{S1})$$

Conditional GAN is the extension case of GAN, where both the generator  $G$  and discriminator  $D$  receive conditional information  $\mathbf{y} \in \mathcal{L} = \{1, 2, \dots, |L|\}$  from  $|L|$  classes [2]. In other words, the inputs are augmented by being concatenated with  $\mathbf{y}$  at each timestamp, i.e.,  $\mathbf{x}_{1:T} \rightarrow [\mathbf{y}; \mathbf{x}_{1:T}]$ . This formulation allows  $G$  to generate samples conditioned on the auxiliary information of  $|L|$ -dimensional categorical labels. In this case, the objective function becomes:

$$\begin{aligned} \min_G \max_D V_{\text{CGAN}} = & \mathbb{E}_{\mathbf{y}, \mathbf{x} \sim p_{\mathbf{y}, \mathbf{x}}} [\log D(\mathbf{x}|\mathbf{y})] \\ & + \mathbb{E}_{\mathbf{y} \sim p_y, \mathbf{v} \sim p_v} [\log(1 - D(G(\mathbf{y}, \mathbf{v})|\mathbf{y}))] \end{aligned} \quad (\text{S2})$$

### C. Shared latent space learning using dual-VAE.

As shown in Fig. S2, the shared latent space is learnt by a dual-VAE network, which contains a pair of encoders (parameterized as  $\phi_{\text{Enc}^c}$  and  $\phi_{\text{Enc}^d}$ ), and a pair of decoders (parameterized as  $\psi_{\text{Dec}^c}$  and  $\psi_{\text{Dec}^d}$ ) of VAE networks, one for each type of timeseries. We found VAE preferable to vanilla autoencoder in our case, considering that (1) the KL regularization in VAE strengthens the learning of the compressed latent representations, which further narrows the domain gap for mixed-type features [10]; (2) VAE can be easily extended to the conditional learning scenario in EHR-M-GAN<sub>cond</sub>. The encoders map the observations into the latent space with  $\text{Enc}(\mathbf{x}) \triangleq q_\phi(\mathbf{z}|\mathbf{x})$ , while the decoders further map the representations into the reconstructed input with  $\text{Dec}(\mathbf{z}) \triangleq p_\psi(\mathbf{x}|\mathbf{z})$ . During the implementation, we found that except for pretraining the dual-VAE, integrating the optimization for decoders during the joint training stage also benefit the generative model from learning an improved representations in the shared latent space.

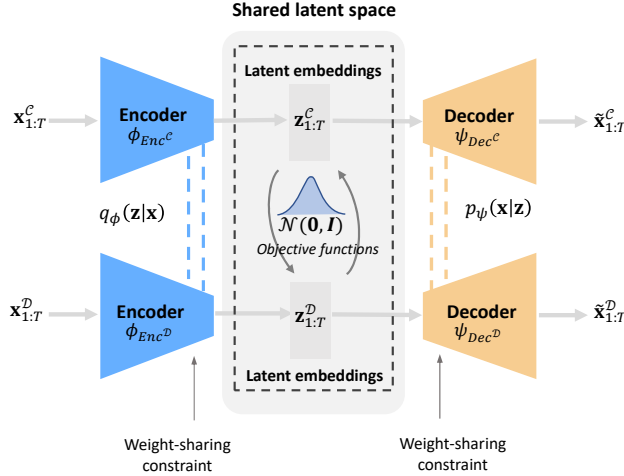


Fig. S2. The network architecture of dual-VAE during the pretraining stage.

In dual-VAE, we enforce a weight-sharing constraint [11] across certain layers within both the encoders pairs and decoders pairs to further eliminate the gap between domains (see Fig. S2). To be specific, only weights of the last few layers of the encoders and the first few layers of the decoders are shared [12]. This forces the encoders to derive the same high-level representations while maintaining different low-level realizations. Meanwhile, it forces the decoders to share the same high-level semantics and decode them into different low-level feature space observations.

### D. Comparison between LSTM and Bilateral-LSTM

To better compare with BLSTMs, we elaborate the architecture of the LSTM network. LSTM utilizes three gates to control the cell state in order to mitigate the problems of gradient vanishing and exploding that appears in the recurrent neural network (RNN) — an input gate  $\mathbf{i}_t$  that controls

the amount of input information to be passed along into the memory cell, a forget gate  $\mathbf{f}_t$  which controls the amount of past information to be neglected, and an output gate  $\mathbf{o}_t$  which controls the update of the new memory cell. The range of outputs from  $\mathbf{i}_t$ ,  $\mathbf{f}_t$  and  $\mathbf{o}_t$  are limited by  $[0, 1]$  due to the sigmoid activation function. At each time step  $t$ , the transition functions in LSTM are as follows:

$$\begin{aligned}
 \mathbf{i}_t &= \sigma(\mathbf{W}_{iv}v_t + \mathbf{W}_{ih}h_{t-1} + \mathbf{b}_i) \\
 \mathbf{f}_t &= \sigma(\mathbf{W}_{fv}v_t + \mathbf{W}_{fh}h_{t-1} + \mathbf{b}_f) \\
 \mathbf{o}_t &= \sigma(\mathbf{W}_{ov}v_t + \mathbf{W}_{oh}h_{t-1} + \mathbf{b}_o) \\
 \tilde{\mathbf{c}}_t &= \tanh(\mathbf{W}_{cv}v_t + \mathbf{W}_{ch}h_{t-1} + \mathbf{b}_c) \\
 \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t \\
 \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t)
 \end{aligned} \tag{S3}$$

where  $\mathbf{c}_t$  denotes the context vector,  $\sigma$  denotes the sigmoid activation function, and  $\odot$  denotes the operation of element-wise multiplication.

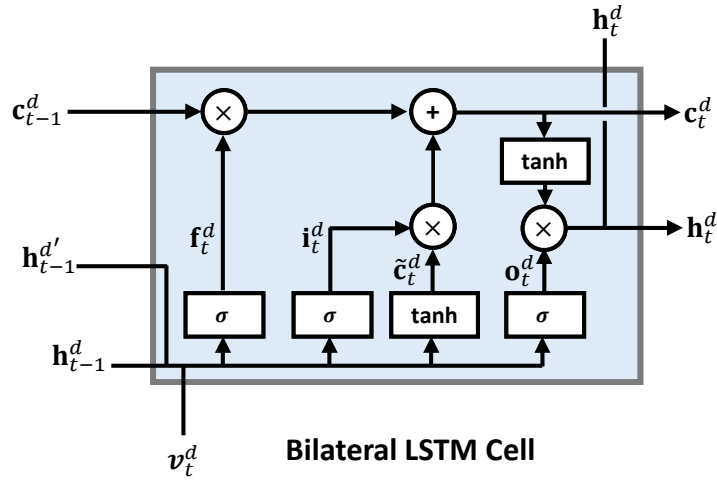


Fig. S3. Illustration of BLSTM cell.

Based on the basic structure of LSTM, the Bilateral Long Short-Term Memory (BLSTM) network is proposed (see Fig. S3). Equations that demonstrate the calculation of BLSTM units can be found in *Methodology* section in the main article.

## E. Algorithms.

**Algorithm S1.** Algorithm of dual-VAE for pretraining.

- 
- 1: **Input:**  $\mathcal{D} = \{(\mathbf{x}_{i,1:T}^{\mathcal{C}}, \mathbf{x}_{i,1:T}^{\mathcal{D}})\}_{i=1}^N$ , learning rate  $\eta_{\text{VAE}}$ , scalar loss weights  $\beta_0, \beta_1, \beta_2, \beta_3$  (if conditional), minibatch size  $n_{mb}$ .
  - 2: Initialize parameters:  $\phi_{\text{Enc}}^{\mathcal{C}}, \phi_{\text{Enc}}^{\mathcal{D}}, \psi_{\text{Dec}}^{\mathcal{C}}, \psi_{\text{Dec}}^{\mathcal{D}}$
  - 3: **for** number of pretrain iterations **do**
  - 4:   Sample a minibatch of  $n_{mb}$  data samples:  $\{(\mathbf{x}_{i,1:T}^{\mathcal{C}}, \mathbf{x}_{i,1:T}^{\mathcal{D}})\}_{i=1}^{n_{mb}} \stackrel{i.i.d.}{\sim} \mathcal{D}$   
    *// Map between features and latent representations:*
  - 5:   **for**  $i = 1, 2, \dots, n_{mb}, t = 1, 2, \dots, T$  **do**
  - 6:      $(\mathbf{z}_{i,t}^{\mathcal{C}}, \mathbf{z}_{i,t}^{\mathcal{D}}) = (\text{Enc}^{\mathcal{C}}(\mathbf{x}_{i,t}^{\mathcal{C}}, \mathbf{z}_{i,t-1}^{\mathcal{C}}), \text{Enc}^{\mathcal{D}}(\mathbf{x}_{i,t}^{\mathcal{D}}, \mathbf{z}_{i,t-1}^{\mathcal{D}}))$
  - 7:      $(\tilde{\mathbf{x}}_{i,t}^{\mathcal{C}}, \tilde{\mathbf{x}}_{i,t}^{\mathcal{D}}) = (\text{Dec}^{\mathcal{C}}(\mathbf{z}_{i,t}^{\mathcal{C}}), \text{Dec}^{\mathcal{D}}(\mathbf{z}_{i,t}^{\mathcal{D}}))$   
    *// Estimate the loss terms:*
  - 8:     **for**  $d \in \{\mathcal{C}, \mathcal{D}\}$  **do**
  - 9:        $\mathcal{L}_d^{\text{ELBO}} = \frac{1}{n_{mb}} \sum_{i=1}^{n_{mb}} [-\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\psi}(\mathbf{x}|\mathbf{z})] + \beta_{\text{KL}} D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{\psi}(\mathbf{z}))]$
  - 10:        $\mathcal{L}^{\text{Match}} = \frac{1}{n_{mb}} \sum_{i=1}^{n_{mb}} [\mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\sum_{t \in \mathcal{T}} \|\mathbf{z}_t^{\mathcal{C}} - \mathbf{z}_t^{\mathcal{D}}\|^2]]$
  - 11:        $\mathcal{L}^{\text{Contra}} = \frac{1}{2n_{mb}} \sum_{i^d=1}^{n_{mb}} \sum_{i^{d'}=1}^{n_{mb}} [\mathcal{L}_{i^d, i^{d'}}^{\text{Contra}} + \mathcal{L}_{i^{d'}, i^d}^{\text{Contra}}]$
  - 12:        $\mathcal{L}_d = \beta_0 \mathcal{L}_d^{\text{ELBO}} + \beta_1 \mathcal{L}^{\text{Match}} + \beta_2 \mathcal{L}^{\text{Contra}}$
  - 13:       **if** conditional:
  - 14:          $\mathcal{L}_d^{\text{Class}} = \frac{1}{n_{mb}} \sum_{i=1}^{n_{mb}} [\mathbb{E}_{\mathbf{z}^d \in \mathcal{H}^S} \text{CE}(f_{\text{linear}}^d(\mathbf{z}^d), \mathbf{y})]$
  - 15:          $\mathcal{L}_d = \beta_0 \mathcal{L}_d^{\text{ELBO}} + \beta_1 \mathcal{L}^{\text{Match}} + \beta_2 \mathcal{L}^{\text{Contra}} + \beta_3 \mathcal{L}_d^{\text{Class}}$   
       *// Update the network weights:*
  - 16:        $\phi_{\text{Enc}}^{\mathcal{C}} = \text{Adam}\left(\frac{\partial \mathcal{L}_{\text{Enc}}^{\mathcal{C}}}{\partial \phi_{\text{Enc}}^{\mathcal{C}}}, \eta_{\text{VAE}}\right), \psi_{\text{Dec}}^{\mathcal{C}} = \text{Adam}\left(\frac{\partial \mathcal{L}_{\text{Dec}}^{\mathcal{C}}}{\partial \psi_{\text{Dec}}^{\mathcal{C}}}, \eta_{\text{VAE}}\right)$
  - 17:        $\phi_{\text{Enc}}^{\mathcal{D}} = \text{Adam}\left(\frac{\partial \mathcal{L}_{\text{Enc}}^{\mathcal{D}}}{\partial \phi_{\text{Enc}}^{\mathcal{D}}}, \eta_{\text{VAE}}\right), \psi_{\text{Dec}}^{\mathcal{D}} = \text{Adam}\left(\frac{\partial \mathcal{L}_{\text{Dec}}^{\mathcal{D}}}{\partial \psi_{\text{Dec}}^{\mathcal{D}}}, \eta_{\text{VAE}}\right)$
  - 18: **Return:**  $\psi_{\text{Dec}}^{\mathcal{C}}, \psi_{\text{Dec}}^{\mathcal{D}}$
-

**Algorithm S2.** Algorithm of EHR-M-GAN.

- 
- 1: **Input:**  $\mathcal{D} = \{(\mathbf{x}_{i,1:T}^C, \mathbf{x}_{i,1:T}^D)\}_{i=1}^N$ , pretrained decoder in dual-VAE for both domains  $\psi_{\text{Dec}}^C, \psi_{\text{Dec}}^D$   
learning rate  $\eta_{\text{GAN}}$ , minibatch size  $n_{mb}$
  - 2: Initialize parameters:  $\theta_G^{\text{CRN}}, \mu_D^C, \mu_D^D$ .
  - 3: **for** number of training iterations **do**
  - 4:   Sample a minibatch of  $n_{mb}$  random noise samples:  $\{(\mathbf{v}_{i,1:T}^C, \mathbf{v}_{i,1:T}^D)\}_{i=1}^{n_{mb}} \stackrel{i.i.d.}{\sim} \mathcal{V}$
  - 5:   **for**  $i = 1, 2, \dots, n_{mb}, t = 1, 2, \dots, T$  **do**  
     *// Generate synthetic latent codes using coupled-generator:*  
 6:      $(\hat{\mathbf{z}}_{i,t}^C, \hat{\mathbf{z}}_{i,t}^D) = G^{\text{CRN}}((\mathbf{v}_{i,t}^C, \mathbf{v}_{i,t}^D), (\mathbf{h}_{i,t-1}^C, \mathbf{h}_{i,t-1}^D))$   
     *// Decode generated latent codes into observational space :*  
 7:      $(\hat{\mathbf{x}}_{i,t}^C, \hat{\mathbf{x}}_{i,t}^D) = (\text{Dec}^C(\hat{\mathbf{z}}_{i,t}^C), \text{Dec}^D(\hat{\mathbf{z}}_{i,t}^D))$
  - 8:   Sample a minibatch of  $n_{mb}$  real data samples:  $\{(\mathbf{x}_{i,1:T}^C, \mathbf{x}_{i,1:T}^D)\}_{i=1}^{n_{mb}} \stackrel{i.i.d.}{\sim} \mathcal{D}$ , and a minibatch  
of  $n_{mb}$  synthetic data samples  $\{(\hat{\mathbf{x}}_{i,1:T}^C, \hat{\mathbf{x}}_{i,1:T}^D)\}_{i=1}^{n_{mb}} \stackrel{i.i.d.}{\sim} \mathcal{D}$   
     *// Distinguish real and fake samples using discriminators and estimate loss :*  
 9:      $\mathcal{L}_{\text{GAN}} = \frac{1}{n_{mb}} \sum_{n=1}^{n_{mb}} [\log(D^C(\mathbf{x}_i^C)) + \log(D^D(\mathbf{x}_i^D))] +$   
 10:      $[\log(1 - D^C(\hat{\mathbf{x}}_i^C)) + \log(1 - D^D(\hat{\mathbf{x}}_i^D))]$   
     *// Update network weights via Adam optimizer :*  
 11:      $\theta_G^{\text{CRN}} = \text{Adam}\left(\frac{\partial \mathcal{L}_{\text{GAN}}}{\partial \theta_G^{\text{CRN}}}, \eta_{\text{GAN}}\right)$   
 12:      $\mu_D^C = \text{Adam}\left(\frac{\partial \mathcal{L}_{\text{GAN}}}{\partial \mu_D^C}, \eta_{\text{GAN}}\right), \mu_D^D = \text{Adam}\left(\frac{\partial \mathcal{L}_{\text{GAN}}}{\partial \mu_D^D}, \eta_{\text{GAN}}\right)$   
     *// Synthesize M pairs of coupled mixed-types of features for M patients:*  
 13:     Sample  $\{(\mathbf{v}_{i,1:T}^C, \mathbf{v}_{i,1:T}^D)\}_{i=1}^M \stackrel{i.i.d.}{\sim} \mathcal{V}$   
 14:     **for**  $i = 1, 2, \dots, M, t = 1, 2, \dots, T$  **do**  
 15:        $(\hat{\mathbf{z}}_{i,t}^C, \hat{\mathbf{z}}_{i,t}^D) = G^{\text{CRN}}((\mathbf{v}_{i,t}^C, \mathbf{v}_{i,t}^D), (\mathbf{h}_{i,t-1}^C, \mathbf{h}_{i,t-1}^D))$   
 16:        $(\hat{\mathbf{x}}_{i,t}^C, \hat{\mathbf{x}}_{i,t}^D) = (\text{Dec}^C(\hat{\mathbf{z}}_{i,t}^C), \text{Dec}^D(\hat{\mathbf{z}}_{i,t}^D))$
  - 17: **Return:**  $\hat{\mathcal{D}} = \{(\hat{\mathbf{x}}_{i,1:T}^C, \hat{\mathbf{x}}_{i,1:T}^D)\}_{i=1}^M$
-



## 2. DATASETS.

We construct the pipeline of data preprocessing based on the work of MIMIC-Extract [13]. Three large-scale, publicly available datasets — MIMIC-III, eICU, and HiRID are processed based on the standard pipeline. The complete steps for data preprocessing include:

- Cohort selection: In cohort selection, patients in three ICU databases are selected based on the same predefined criteria (see Section 2.A for details).
- Timeseries features extraction: Then, the timeseries features are extracted based on the lists provided in Section 2.C. Both continuous-valued and discrete-valued features are selected accordingly.
- Unit conversion and outlier filtering: Due to the fact that clinical data is often measured in different units, unit conversion are applied (such as converting Fahrenheit to Celsius for *Temperature*). For outlier filtering, a reasonable physiologically valid range are applied for different measurements (see [13] for details).
- Semantic grouping: Next, semantically similar variables are grouped based on clinical concepts (such as *Heart Rate* is recorded as ItemID 211 in CareVUE EHR systems and ItemID 220045 under MetaVision EHR systems). A clinical taxonomy are used to aggregate features that are semantically equivalent [13].
- Hourly aggregation: Timestamps with different granularity are provided for different in three databases. Time-varying physiological signals such as *Heart Rate* are frequently measured (e.g., most parameters under bedside monitoring are recorded every 2 minutes in HiRID dataset [14]). While other features such as laboratory test results are measured infrequently. Therefore, we hourly aggregate the timeseries further into a uniform hourly bucket.
- Imputation and normalization: Finally, imputation method in Section B are used and normalization are applied to obtain the final result of the data matrix.

### A. Cohort selection criteria.

In line with the previous literature [13, 15], the cohort are selected based on the following criteria: (1) Only the first known ICU admission of the patient is selected. This is because the patients who have multiple ICU admission records typically require specific treatments for life-support intervention; (2) Patient has to be an adult at the time of ICU admission (at least 15); (3) The duration of the patients' ICU stay is at least 12 hours and less than 10 days. This is because the treatment for patients who have longer hours in the ICU stay usually indicates their physiological changes can not be directly linked to the positive effect of the treatment (but compensating for the life support treatment being taken off) [15].

### B. Imputation method.

For continuous-valued timeseries, missing data is imputed based on method of *Simple Imputation* [16]. The missing timeseries data is imputed as the last observed value, or individual-specific mean if no previous observation is provided. Else, if there is no observation for the subject, the imputation value is set to the global mean of the entire cohort. Compared to imputation methods developed upon customized RNN models or explicitly designed for the applied domains, it does not rely on additional information such as the prediction labels therefore more generalizable. Though simple, such method has been widely applied in clinical timeseries analysis [17] including MIMIC-III datasets [13, 16, 18]. For discrete-valued timeseries, we followed the preprocessing rules in MIMIC-EXTRACT. For intermittent interventions such as oral antibiotics, its status is regarded as 'not applied' when missing. For intervention with multi-hour continuous duration, such as mechanical ventilation, the missed status is considered to be consistent with the previous status until the new administration occurs. Therefore, the imputation method was not applied to the discrete-valued data.

### C. Timeseries features extraction.

Features of continuous-valued and discrete-valued timeseries are extracted for three critical care databases based on the following lists (for MIMIC-III dataset, see Table S1, S2; for eICU dataset see Table S3, S4; for HiRID dataset, see Table S5, S6).

**Table S1.** List of vital sign and laboratory test features for MIMIC-III dataset. Features are further extracted based on the preprocessed results of MIMIC-Extract (see Appendix A. Feature set in [13]). The dimension of continuous-valued features for MIMIC-III dataset during model’s training is 78.

Measurement			
heart rate	respiratory rate	systolic blood pressure	diastolic blood pressure
mean blood pressure	oxygen saturation	temperature	glucose
central venous pressure	hematocrit	potassium	sodium
chloride	pulmonary artery pressure systolic	hemoglobin	ph
creatinine	blood urea nitrogen	bicarbonate	platelets
anion gap	co2 (etco2, pco2, etc.)	partial pressure of carbon dioxide	magnesium
white blood cell count	positive end-expiratory pressure set	calcium	fraction inspired oxygen set
red blood cell count	mean corpuscular hemoglobin concentration	mean corpuscular hemoglobin	mean corpuscular volume
tidal volume observed	partial thromboplastin time	prothrombin time inr	prothrombin time pt
phosphate	phosphorous	peak inspiratory pressure	calcium ionized
respiratory rate set	fraction inspired oxygen	tidal volume set	partial pressure of oxygen
cardiac index	co2	systemic vascular resistance	potassium serum
tidal volume spontaneous	plateau pressure	pulmonary artery pressure mean	cardiac output thermodilution
lactate	lactic acid	bilirubin	asparate aminotransferase
alanine aminotransferase	alkaline phosphate	positive end-expiratory pressure	albumin
troponin-t	neutrophils	lymphocytes	monocytes
ph urine	fibrinogen	lactate dehydrogenase	basophils
cardiac output fick	creatinine urine	pulmonary capillary wedge pressure	red blood cell count urine
white blood cell count urine	cholesterol	cholesterol hdl	post void residual
cholesterol ldl	chloride urine		

**Table S2.** List of medical intervention features for MIMIC-III dataset, where **Features** indicates the name of the intervention features during model’s training, **Category of treatment** shows the category of treatment that the specific intervention feature belongs to, and **Source** is the corresponding chart(s) where the variable is extracted based on<sup>1</sup>. The dimension of discrete-valued features for MIMIC-III dataset during model’s training is 20.

Category of treatment	Features	Source
Oxygen therapy	supplemental oxygen mechanical ventilation	chartevents, procedureevents_mv
Vasopressor	adenosine dobutamine dopamine epinephrine isuprel milrinone norepinephrine phenylephrine vasopressin	inputevents_cv, inputevents_mv
Antibiotics	antibiotics	prescriptions
Renal therapy	continuous renal replacement therapy	chartevents
Invasive lines	arterial line central line	procedureevents_mv, chartevents
Colloid bolus	colloid bolus	inputevents_mv, inputevents_cv, chartevents
Crystalloid bolus	crystalloid bolus	inputevents_mv, inputevents_cv

<sup>1</sup><https://github.com/MIT-LCP/mimic-code>

**Table S3.** List of vital sign and laboratory test features for eICU dataset. Features are selected base on the recommendation from Rocheteau et al [19]. The dimension of continuous-valued features for eICU dataset during model’s training is 55.

Measurement			
Hct	calcium	anion gap	MCH
troponin - I	MCHC	PT	PT - INR
-eos	potassium	-basos	albumin
-polys	lactate	glucose	creatinine
AST (SGOT)	Hgb	MPV	WBC × 1000
ALT (SGPT)	HCO3	MCV	-lymphs
Exhaled MV	RDW	chloride	sodium
bicarbonate	pH	urinary specific gravity	SaO2
Tidal Volume (set)	-monos	Heart Rate	BUN
platelets × 1000	total bilirubin	Exhaled TV (patient)	alkaline phos
Noninvasive BP Diastolic	Noninvasive BP Mean	Noninvasive BP Systolic	Base Excess
paO2	FiO2	Temperature	RBC
PTT	magnesium	RR	SpO2
total protein	paCO2	phosphate	

**Table S4.** List of medical intervention features for eICU dataset, where **Features** indicates the name of the intervention features during model’s training, **Category of treatment** shows the category of treatment that the specific intervention feature belongs to, and **Source** is the corresponding chart(s) where the variable is extracted based on<sup>2</sup>. The dimension of discrete-valued features for eICU dataset during model’s training is 19.

Category of treatment	Features	Source
Oxygen therapy	supplemental oxygen mechanical ventilation	respiratorycharting, nursecharting, treatment
Vasopressor	dopamine epinephrine norepinephrine phenylephrine vasopressin milrinone dobutamine	infusionDrug
Anesthesia	fentanyl propofol midazolam dexmedetomidine	infusionDrug
Anticoagulants	heparin	infusionDrug
Insulin	insulin	infusionDrug
Antibiotics	antibiotics	medication

<sup>2</sup><https://github.com/MIT-LCP/eicu-code>

**Table S5.** List of vital sign and laboratory test features for HiRID dataset. Features are extracted based on the official HiRID preprocessing codes (meta-variables from Merging stage<sup>3</sup>) [14]. The dimension of continuous-valued features for HiRID dataset during model’s training is 50.

Measurement			
HR	T Central	ABPs	ABPd
ABPm	NIBPs	NIBPd	NIBPm
PAPm	PAPs	PAPd	CO
SvO2(m)	ZVD	ST1	ST2
ST3	SpO2	ETCO2	RR
OUTurine/h	ICP	Liquor/h	a-BE
a_COHb	a_Hb	a_HCO3-	a_Lac
a_MetHb	a_pH	a_pCO2	a_PO2
a_SO2	K+	Na+	Cl-
Ca2+ ionized	phosphate	Mg_lab	Urea
creatinine	INR	glucose	Hb
MCHC	MCV	platelet count	MCH
C-reactive protein	total white blood cell count		

<sup>3</sup><https://github.com/ratschlab/HiRID-ICU-Benchmark>

**Table S6.** List of medical intervention features for HiRID dataset, where **Features** indicates the name of the intervention features during model’s training, **Category of treatment** shows the category of treatment that the specific intervention feature belongs to, and **Source** is the corresponding feature names in the official HiRID preprocessing codes (meta-variables from Merging stage <sup>4</sup>) [14] that we extracted based on. The dimension of discrete-valued features for HiRID dataset during model’s training is 39.

Category of treatment	Features	Source
Oxygen therapy	supplemental oxygen	vm23
	mechanical ventilation	vm60
Crystalloids	crystalloids	vm33
Colloids	colloids	vm34
Renal therapy	haemofiltration	vm72
Blood transfusion	packed red blood cells	pm35
	FFP	pm36
	platelets	pm37
Vaspressor/inotropes	norepinephrine	pm39
	epinephrine	pm40
	dobutamine	pm41
	milrinone	pm42
	levosimendan	pm43
	theophyllin	pm44
	vasopressin	pm45
	desmopressin	pm46
Vasodilators	vasodilators	pm47
Antihypertensives	ACE inhibitors	pm48
	Calcium channel blockers	pm50
	Beta-blocker	pm51
Antiarrhythmics	adenosine	pm53
	digoxin	pm54
	amiodarone	pm55
	atropine	pm56
Antibiotics	antibiotics	pm73
	antimycotic	pm74
	antiviral	pm75
Insulin	insulin	pm82, pm83
Pain killers	opioid	pm86
	non-opioid	pm87
Steroids	steroids	pm91
Anticoagulants	heparin	pm95

<sup>4</sup><https://github.com/ratschlab/HiRID-ICU-Benchmark>

### 3. MODEL TRAINING.

#### A. Implementation details

During the model training of EHR-M-GAN, the hyperparameters are optimized based on the comparison between the synthetic data and leave-out real data, estimated by mean maximum discrepancy (MMD) for continuous data and mean squared errors (MSEs) over the Bernoulli probability for discrete data, as the scoring functions. Visual inspection is also used during training to intuitively compare the resemblance between synthetic and real data. Table A shows the hyperparameter values of the network architecture for searching over. The optimal hyperparameters for GANs’ training is listed in our GitHub codebase (see *train\_config.py* file). The model which generates the best results is saved and used for the final results.

**Table S7.** List of hyperparameters of EHR-M-GAN.

Hyperparameters	Searching space
Batch size	{128, 256, 512}
Epochs for pretraining	{200, 500, 800}
Epochs for training GANs	{500, 800}
Rounds for jointly training $G/D/V$	{3/1/1, 5/1/1}
Learning rate for pretraining	{0.001, 0.0001, 0.0005}
Learning rate for training GANs	{0.001, 0.0001, 0.0005}
Depths for encoders and decoders	{3, 5}
Depths for generators	{3, 5}
Depths for discriminators	{1, 3, 5}
Sizes for encoders and decoders	{64, 128, 256}
Sizes for generators	{256, 512}
Sizes for discriminators	{256, 512}
Weight scalar for pretraining	{0.01, 0.1, 0.25, 0.5, 1, 2, 5}
Weight scalar for training GANs	{0.1, 0.5, 1, 5, 10, 20}
Optimizer	Adam

During the pretraining stage of dual-VAE module, we implemented the VAEs with recurrent neural network based on Google DeepMind’s “DRAW” — Deep Recurrent Attentive Writer [20]. Instead of automatically generating the entire images/timeseries at once, it utilizes a sequential variational auto-encoding framework that enables the iterative generation of multivariate time-series. The reconstruction loss on the leave-out validation set (i.e., the “one-to-one” mapping) is used for optimizing the hyperparameters in dual-VAEs (see Table A).

Furthermore, to stabilize GANs’ training and overcome the problem of mode collapse, training strategies such as feature matching loss is utilized [21]. Feature matching is a regularizing objective that prevents the generator in GANs from overtraining on the current discriminator. It has been shown effective to stabilize the GANs’ training as it calculates the *statistics* of the real data per minibatch, instead of directly maximizing the output of the discriminator. The formal definition of feature matching loss is described as follows:

$$L = \left\| \mathbb{E}_{x \sim p_{\text{data}}} \mathbf{f}(x) - \mathbb{E}_{z \sim p_z} \mathbf{f}(G(z)) \right\|_2^2$$

where  $f(x)$  is the feature representation of the intermediate layer of the discriminator (layer before the final classification).

#### B. Ablation study for training dual-VAE

Multiple losses are placed when optimizing the shared latent space in the dual-VAE module. Except for the standard evidence lower bound (ELBO) loss in VAE, external losses, namely (1)



Matching loss; (2) Contrastive loss, and (3) Semantic loss (for the conditional variation of our proposed model) are used. Also, during the implementation, the weight-sharing constraint is adopted for specific layers in dual-VAE’s encoder and decoder pairs to extract the high-level representations from mixed-type inputs (see Section S.1.C *Shared latent space learning using dual-VAE* for details). In order to analyze the contribution of each aforementioned component when training dual-VAE, we perform an ablation study by varying the corresponding training configurations (see Table S8) using MIMIC-III dataset as an example. The performance for synthesizing continuous-valued timeseries is evaluated by maximum mean discrepancy (MMD) and discriminative score. For discrete-valued timeseries, the performance of GANs is evaluated by dimensional-wise probability (DWP) quantified by the averaged root mean squared errors (RMSEs) across all feature dimensions (see *Dimension-wise probability* section in the main text for details) and discriminative score. The results of the ablation study are shown in Table S8.

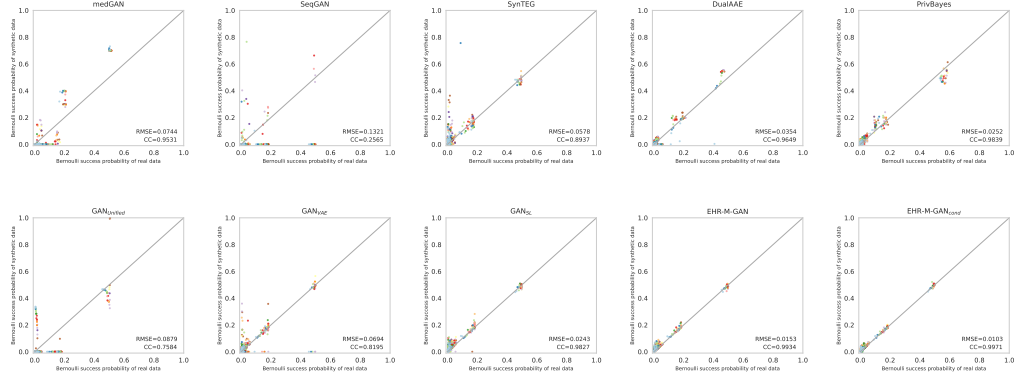
**Table S8.** The ablation study for components in Dual-VAE on MIMIC-III dataset. ‘Baseline’ represents the proposed GAN models (EHR-M-GAN or EHR-M-GANcond) with all components included. The quality of synthetic continuous-valued timeseries is evaluated by MMD and discriminative score (both the lower the better). The quality of synthetic discrete-valued timeseries is evaluated by averaged RMSEs in DWP and discriminative score (both the lower the better).

Model	Training configuration	Continuous-valued data		Discrete-valued data	
		MMD	Discriminative score	DWP (RMSEs)	Discriminative score
EHR-M-GAN	Baseline	<b>0.692 ± 0.034</b>	<b>0.746 ± 0.018</b>	<b>0.0104 ± 0.0006</b>	<b>0.813 ± 0.026</b>
	w/o Matching loss	0.722 ± 0.023	0.758 ± 0.015	0.0112 ± 0.0010	0.827 ± 0.019
	w/o Contrastive loss	0.719 ± 0.017	0.762 ± 0.012	0.0109 ± 0.0009	0.830 ± 0.023
	w/o Shared weights	0.704 ± 0.031	0.749 ± 0.019	0.0107 ± 0.0008	0.816 ± 0.035
EHR-M-GANcond	Baseline	<b>0.604 ± 0.027</b>	<b>0.729 ± 0.025</b>	<b>0.0093 ± 0.0005</b>	<b>0.784 ± 0.024</b>
	w/o Matching loss	0.634 ± 0.026	0.736 ± 0.017	0.0106 ± 0.0013	0.795 ± 0.022
	w/o Contrastive loss	0.629 ± 0.022	0.739 ± 0.020	0.0108 ± 0.0007	0.796 ± 0.028
	w/o Semantic loss	0.647 ± 0.034	0.743 ± 0.011	0.0114 ± 0.0004	0.798 ± 0.030
	w/o Shared weights	0.609 ± 0.035	0.732 ± 0.014	0.0097 ± 0.0012	0.786 ± 0.027

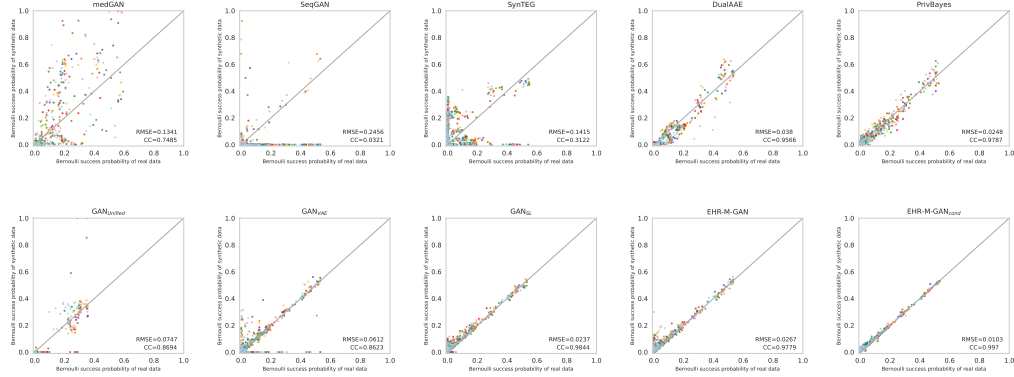
As shown in Table S8, both matching loss and contrastive loss contribute to the improvement of EHR-M-GAN’s performance when generating mixed-type timeseries data. For example, the absence of the contrastive loss leads to a noticeable degradation in the quality of the synthetic discrete-valued timeseries (evaluated by discriminative score). Also, removing the matching loss causes the increase of the MMD between real and synthetic continuous-valued timeseries. The weight-sharing scheme between the encoder and decoder architectures in the dual-VAE also boosts GANs’ performance but within a limited range. For EHR-M-GANcond model, the effectiveness of the components that appear in EHR-M-GAN can still be observed. On the other hand, semantic loss, which injects conditional information into the networks, plays a major role in synthesizing more realistic patient trajectories. The results in Table S8 show that the impact of the semantic loss exceeds the other two losses in learning the valid shared latent representations in dual-VAE.

## 4. RESULTS.

### A. Dimension-wise probability.



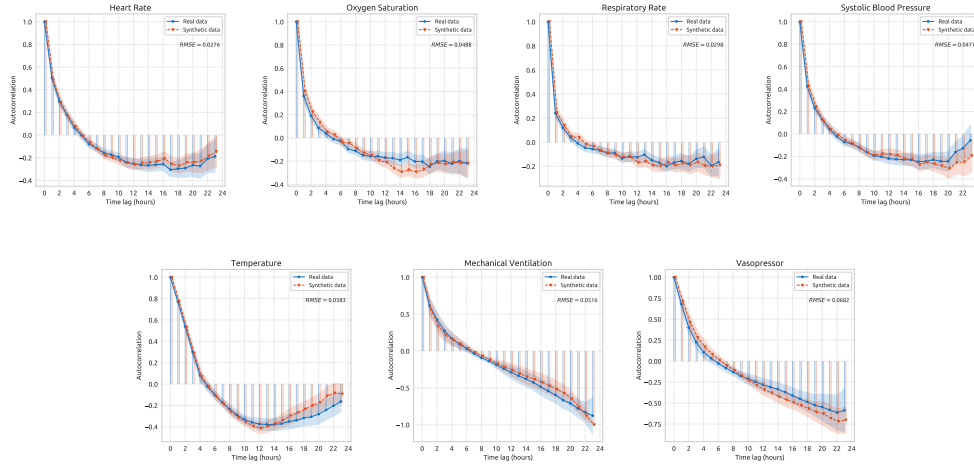
**Fig. S4. Scatterplot of the dimension-wise probability test on eICU dataset.** The x-axis and y-axis represents the probability distribution for the real data and synthetic data with same sample size, respectively. The optimal performance appears along the diagonal line. Each dot represents a treatment status at a particular time in the patient EHR data. The optimal performance appears along the diagonal line. The corresponding CCs ( $[0, 1]$ , the higher the better) and RMSEs ( $[0, +\infty)$ , the lower the better) are also calculated to quantify the probability distribution similarities between the real and synthetic EHRs timeseries.



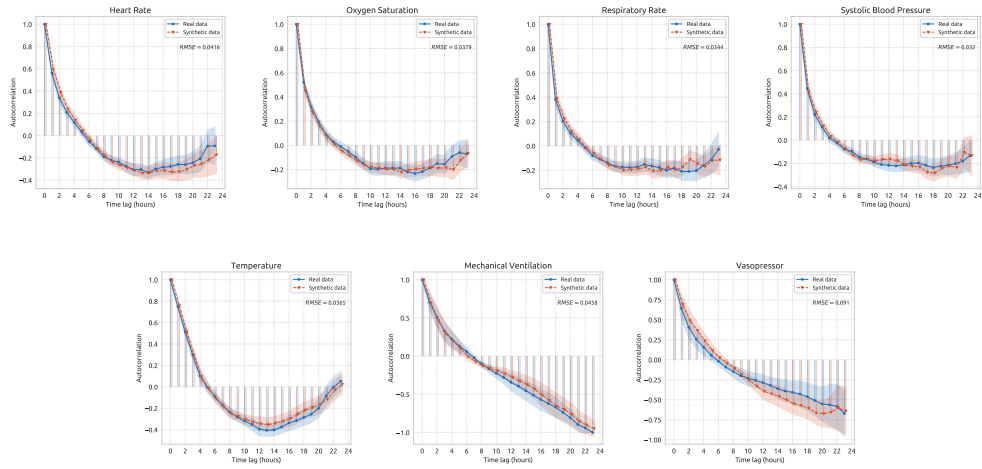
**Fig. S5. Scatterplot of the dimension-wise probability test on eICU dataset.** The x-axis and y-axis represents the probability distribution for the real data and synthetic data with same sample size, respectively. The optimal performance appears along the diagonal line. Each dot represents a treatment status at a particular time in the patient EHR data. The optimal performance appears along the diagonal line. The corresponding CCs ( $[0, 1]$ , the higher the better) and RMSEs ( $[0, +\infty)$ , the lower the better) are also calculated to quantify the probability distribution similarities between the real and synthetic EHRs timeseries.

## B. Temporal characteristics.

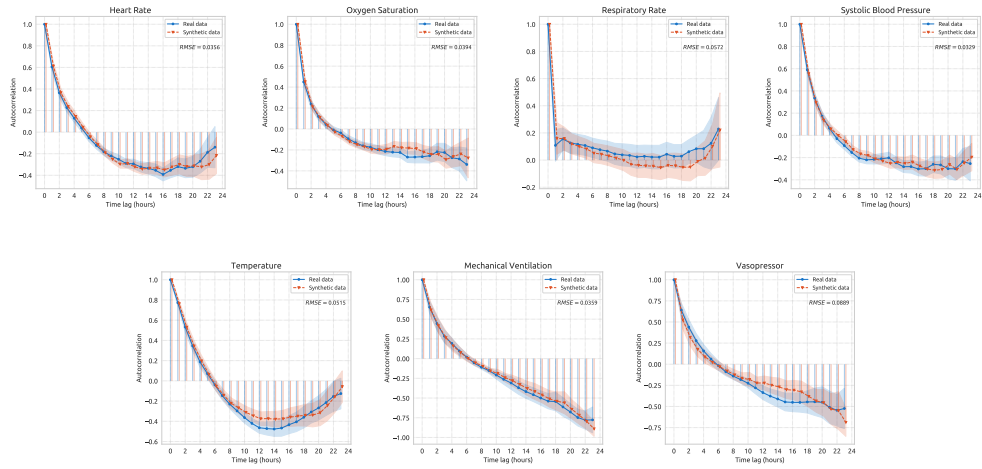
Fig. S6-S8 show the autocorrelation function (ACF) of real timeseries and synthetic timeseries generated by EHR-M-GAN on continuous-valued features (including *Heart Rate*, *Oxygen Saturation*, *Respiratory Rate*, *Systolic Blood Pressure*, and *Temperature*) and discrete-valued features (including *Vasopressor* and *Mechanical Ventilation*). The averaged ACF is calculated over the population sampled randomly from both real and synthetic patient data. The averaged autocorrelation for real patient trajectories (solid blue line) and synthetic patient trajectories (red dashed line) are calculated, with the light colored regions indicating the corresponding 95% confidence interval. The root-mean-square errors (RMSEs) are also calculated for the two curves on each variable to quantitatively evaluate the temporal characteristics captured by the synthetic data.



**Fig. S6.** Autocorrelation function (ACF) of real data and synthetic data generated by EHR-M-GAN on MIMIC-III dataset.



**Fig. S7.** Autocorrelation function (ACF) of real data and synthetic data generated by EHR-MGAN on eICU dataset.

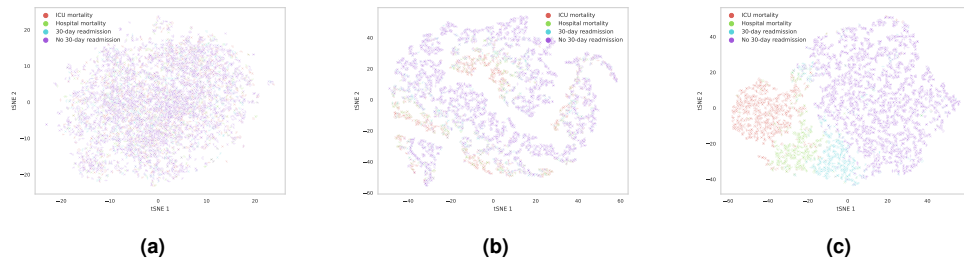


**Fig. S8.** Autocorrelation function (ACF) of real data and synthetic data generated by EHR-MGAN on HiRID dataset.

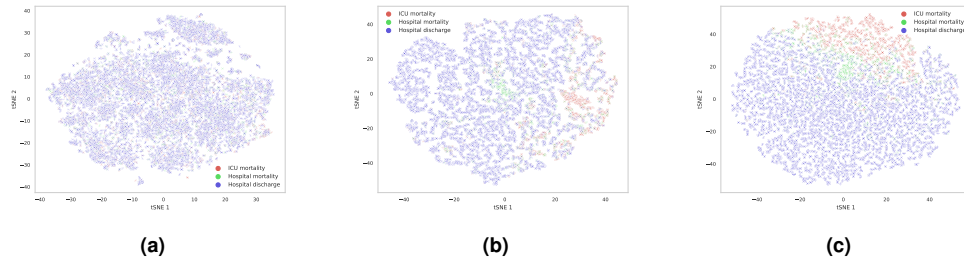
### C. Embedding visualisation.

We apply t-SNE to qualitatively visualise the latent representations generated by EHR-M-GAN and EHR-M-GAN<sub>cond</sub> on three critical care databases. The latent embedding vectors are induced by the encoders in the *dual*-VAE during learning the shared latent space representations (See Methods section, p12, for details). The t-SNE embedding results on raw timeseries are also included for comparison.

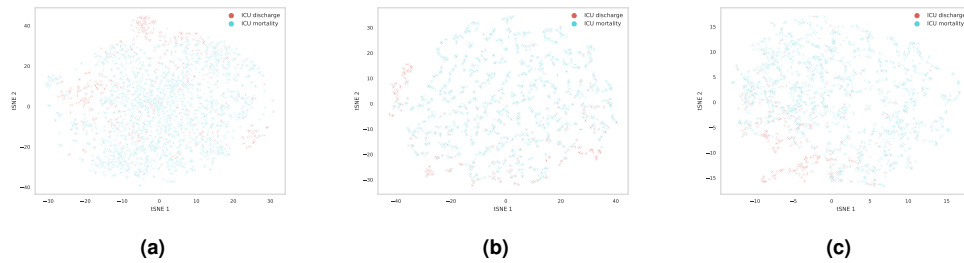
It can be seen that better separability of the representation clusters in the shared latent space is shown in the embeddings obtained from EHR-M-GAN<sub>cond</sub> compared with EHR-M-GAN and raw data. This illustrates the superiority of the EHR-M-GAN<sub>cond</sub> in terms of learning the contextual information from the patient trajectories. It therefore can be inferred that the conditional extension of the proposed model can further yield benefits by synthesizing condition-specific EHR timeseries with respect to distinctive patient health status.



**Fig. S9.** t-SNE embedding visualization from MIMIC-III dataset on (a) raw patient trajectories, (b) latent embeddings generated with EHR-M-GAN, and (c) latent embeddings generated with EHR-M-GAN<sub>cond</sub>.

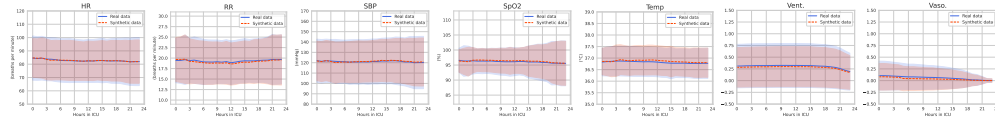


**Fig. S10.** t-SNE embedding visualization from eICU dataset on (a) raw patient trajectories, (b) latent embeddings generated with EHR-M-GAN, and (c) latent embeddings generated with EHR-M-GAN<sub>cond</sub>.

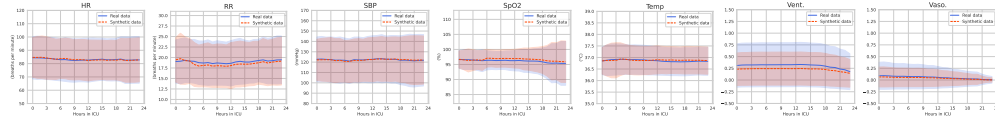


**Fig. S11.** t-SNE embedding visualization from HiRID dataset on (a) raw patient trajectories, (b) latent embeddings generated with EHR-M-GAN, and (c) latent embeddings generated with EHR-M-GAN<sub>cond</sub>.

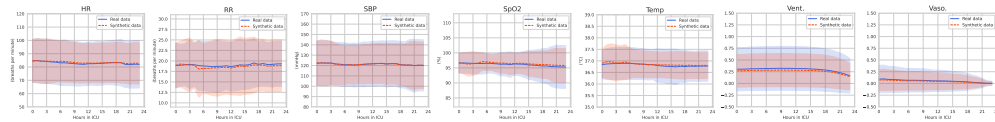
#### D. Patient trajectories visualisation.



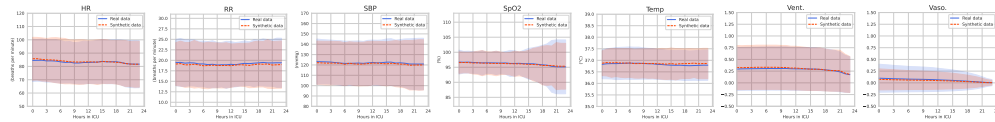
**Fig. S12. Comparison of patient trajectories.** The distribution of values at each timepoint (mean and standard deviation) are compared between real and synthetic patient trajectory produced by EHR-M-GAN<sub>cond</sub>, under the condition of ICU mortality.



**Fig. S13. Comparison of patient trajectories.** The distribution of values at each timepoint (mean and standard deviation) are compared between real and synthetic patient trajectory produced by EHR-M-GAN<sub>cond</sub>, under the condition of Hospital mortality.



**Fig. S14. Comparison of patient trajectories.** The distribution of values at each timepoint (mean and standard deviation) are compared between real and synthetic patient trajectory produced by EHR-M-GAN<sub>cond</sub>, under the condition of 30-day readmission.



**Fig. S15. Comparison of patient trajectories.** The distribution of values at each timepoint (mean and standard deviation) are compared between real and synthetic patient trajectory produced by EHR-M-GAN<sub>cond</sub>, under the condition of No 30-day readmission.

## REFERENCES

1. Olof Mogren. C-rnn-gan: Continuous recurrent neural networks with adversarial training. *arXiv preprint arXiv:1611.09904*, 2016.
2. Cristóbal Esteban, Stephanie L Hyland, and Gunnar Rätsch. Real-valued (medical) time series generation with recurrent conditional gans. *arXiv preprint arXiv:1706.02633*, 2017.
3. Jinsung Yoon, Daniel Jarrett, and Mihaela Van der Schaar. Time-series generative adversarial networks. *Advances in neural information processing systems*, 32, 2019.
4. Ziqi Zhang, Chao Yan, Thomas A Lasko, Jimeng Sun, and Bradley A Malin. Synteg: a framework for temporal structured electronic health data simulation. *Journal of the American Medical Informatics Association*, 28(3):596–604, 2021.
5. Ziqi Zhang, Chao Yan, and Bradley A Malin. Keeping synthetic patients on track: feedback mechanisms to mitigate performance drift in longitudinal health data simulation. *Journal of the American Medical Informatics Association*, 29(11):1890–1898, 2022.
6. Dongha Lee, Hwanjo Yu, Xiaoqian Jiang, Deevakar Rogith, Meghana Gudala, Mubeen Tejani, Qiuchen Zhang, and Li Xiong. Generating sequential electronic health records using dual adversarial autoencoder. *Journal of the American Medical Informatics Association*, 27(9):1411–1419, 2020.
7. Kieran Chin-Cheong, Thomas Sutter, and Julia E Vogt. Generation of heterogeneous synthetic electronic health records using gans. In *workshop on machine learning for health (ML4H) at the 33rd conference on neural information processing systems (NeurIPS 2019)*. ETH Zurich, Institute for Machine Learning, 2019.
8. Chao Yan, Ziqi Zhang, Steve Nyemba, and Bradley A Malin. Generating electronic health records with multiple data types and constraints. In *AMIA annual symposium proceedings*, volume 2020, page 1335. American Medical Informatics Association, 2020.
9. Shannon KS Kroes, Matthijs van Leeuwen, Rolf HH Groenwold, and Mart P Janssen. Generating synthetic mixed discrete-continuous health records with mixed sum-product networks. *Journal of the American Medical Informatics Association*, 2022.
10. Ziyu Wan, Bo Zhang, Dongdong Chen, Pan Zhang, Dong Chen, Jing Liao, and Fang Wen. Old photo restoration via deep latent space translation. *arXiv preprint arXiv:2009.07047*, 2020.
11. Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in neural information processing systems*, pages 700–708, 2017.
12. Ming-Yu Liu and Oncel Tuzel. Coupled generative adversarial networks. *Advances in neural information processing systems*, 29:469–477, 2016.
13. Shirly Wang, Matthew BA McDermott, Geeticka Chauhan, Marzyeh Ghassemi, Michael C Hughes, and Tristan Naumann. Mimic-extract: A data extraction, preprocessing, and representation pipeline for mimic-iii. In *Proceedings of the ACM Conference on Health, Inference, and Learning*, pages 222–235, 2020.
14. Hugo Yèche, Rita Kuznetsova, Marc Zimmermann, Matthias Hüser, Xinrui Lyu, Martin Faltys, and Gunnar Ratsch. Hirid-icu-benchmark—a comprehensive machine learning benchmark on high-resolution icu data. 2021.
15. Mike Wu, Marzyeh Ghassemi, Mengling Feng, Leo A Celi, Peter Szolovits, and Finale Doshi-Velez. Understanding vasopressor intervention and weaning: risk prediction in a public heterogeneous clinical time series database. *Journal of the American Medical Informatics Association*, 24(3):488–495, 2017.
16. Zhengping Che, Sanjay Purushotham, Kyunghyun Cho, David Sontag, and Yan Liu. Recurrent neural networks for multivariate time series with missing values. *Scientific reports*, 8(1):1–12, 2018.
17. Alexander Meyer, Dina Zverinski, Boris Pfahringer, Jörg Kempfert, Titus Kuehne, Simon H Sündermann, Christof Stamm, Thomas Hofmann, Volkmar Falk, and Carsten Eickhoff. Machine learning for real-time prediction of complications in critical care: a retrospective study. *The Lancet Respiratory Medicine*, 6(12):905–914, 2018.
18. Sanjay Purushotham, Chuizheng Meng, Zhengping Che, and Yan Liu. Benchmarking deep learning models on large healthcare datasets. *Journal of biomedical informatics*, 83:112–134, 2018.
19. Emma Rocheteau, Pietro Liò, and Stephanie Hyland. Temporal pointwise convolutional networks for length of stay prediction in the intensive care unit. In *Proceedings of the Conference on Health, Inference, and Learning*, pages 58–68, 2021.
20. Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Rezende, and Daan Wierstra. Draw: A recurrent neural network for image generation. In *International conference on machine learning*,



pages 1462–1471. PMLR, 2015.

21. Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016.