# TarDiff: Target-Oriented Diffusion Guidance for Synthetic Electronic Health Record Time Series Generation

**Bowen Deng[1]***, **Chang Xu[2]†, Hao Li[3], Yuhao Huang[4], Min Hou[5], Jiang Bian[2]**
[1] Peking University   [2] Microsoft Research Asia   [3] University of Manchester
[4] Nanjing University   [5] Hefei University of Technology
devin@stu.pku.edu.cn   {chanx, jiang.bian}@microsoft.com
hao.li-2@manchester.ac.uk   huangyh@smail.nju.edu.cn
hmhoumin@gmail.com

## Abstract

Synthetic Electronic Health Record (EHR) time-series generation is crucial for advancing clinical machine learning models, as it helps address data scarcity by providing more training data. However, most existing approaches focus primarily on replicating statistical distributions and temporal dependencies of real-world data. We argue that fidelity to observed data alone does not guarantee better model performance, as common patterns may dominate, limiting the representation of rare but important conditions. This highlights the need for generate synthetic samples to improve performance of specific clinical models to fulfill their target outcomes. To address this, we propose *TarDiff*, a novel target-oriented diffusion framework that integrates task-specific influence guidance into the synthetic data generation process. Unlike conventional approaches that mimic training data distributions, TarDiff optimizes synthetic samples by quantifying their expected contribution to improving downstream model performance through influence functions. Specifically, we measure the reduction in task-specific loss induced by synthetic samples and embed this influence gradient into the reverse diffusion process, thereby steering the generation towards utility-optimized data. Evaluated on **six** publicly available EHR datasets, TarDiff achieves state-of-the-art performance, outperforming existing methods by up to 20.4% in AUPRC and 18.4% in AUROC. Our results demonstrate that TarDiff not only preserves temporal fidelity but also enhances downstream model performance, offering a robust solution to data scarcity and class imbalance in healthcare analytics.

## 1 Introduction

Healthcare is a cornerstone of societal well-being, especially as the world faces an aging population and the rising burden of chronic diseases, placing increasing pressure on healthcare systems worldwide [Liang et al., 2024]. Traditionally, medical diagnoses have relied on human expertise, but with advancements in machine learning, Electronic Health Records (EHRs)—which digitally store a patient's medical history, including demographic attributes [Maweu et al., 2021, Li et al., 2023], vital signs [Tseng et al., 2022], and lab measurements—have become invaluable for clinical research [Kaushik et al., 2020, Schlegel et al., 2023]. EHR time series data, such as Electroencephalography (EEG) for neurological analysis and Electrocardiography (ECG) for heart condition diagnosis, provide

---

*The work was conducted during the internship of Bowen Deng and Hao Li at Microsoft Research.
†Corresponding author.

critical insights for medical decision-making. Leveraging these rich data sources, machine learning models trained in a data-driven manner are then applied to various downstream tasks, including disease diagnosis, prognosis prediction, and treatment planning [Shickel et al., 2017, Goldstein et al., 2016, Li et al., 2024b, Nagar et al., 2024].

However, obtaining and utilizing EHR data remains a significant challenge due to medical-related factors, such as strict privacy regulations, data incompleteness resulting from sensor failures, and difficulties in accurate labeling. As a result, synthesizing EHR data has gained increasing attention. Existing work has explored a wide range of methods including rule-based approaches and generative models. Rule-based techniques—such as time warping, jittering, and interpolation[Wang et al., 2024, Wen et al., 2020]—are favored for their simplicity and efficiency, yet they often fail to capture the intricate temporal dependencies and pathological patterns present in clinical data. GANs [Schön et al., 2023] and VAEs [Kingma, 2013] have been employed to generate high-fidelity signals, with notable examples including TimeGANs [Yoon et al., 2019], TimeVAE [Desai et al., 2021], and TimeVQ-VAE [Lee et al., 2023], which have demonstrated improved performance in various biohealthcare applications. More recently, diffusion models [Ho et al., 2020, Song et al., 2020, Fan et al., 2024, Li et al., 2025] have emerged as a promising alternative; methods such as TimeDiff [Tian et al., 2024], DiffusionTS [Yuan and Qiao, 2024] and BioDiffusion [Li et al., 2024a] similarly seek to approximate the underlying data distribution through iterative refinement, emphasizing the generation of realistic time series data.

While these efforts have significantly advanced the field, they primarily focus on generating synthetic data by mimicking the empirical distribution of training samples. We argue that this approach inherently overlooks the effectiveness of the generated data for downstream medical tasks. For instance, in rare disease diagnosis, where positive samples are scarce in real-world datasets, generating synthetic data purely based on observed distributions may exacerbate biases in the downstream model [Gupta et al., 2021], making it more inclined to diagnose common conditions while failing to recognize rare diseases effectively. Consequently, such data generation strategies do not necessarily enhance model performance in rare disease diagnosis and may even worsen the imbalance in predictive accuracy [Huo et al., 2022]. Therefore, we propose that EHR data generation should not merely replicate statistical patterns of training data but should instead be guided by its utility in training more effective models to fulfill the targets of downstream tasks. This calls for an adaptive generation strategy that actively optimizes synthetic data to enhance model learning for specific medical applications.

In this work, we investigate how to develop a model for generating EHR data that is specifically tailored to enhance downstream model performance for target tasks. Recent advancements in diffusion models have shown promising results in time series generation [Tian et al., 2024, Yuan and Qiao, 2024, Li et al., 2024a], with conditional diffusion models [Huang et al., 2025, Dhariwal and Nichol, 2021, Ho and Salimans, 2022] offering significant insights for this work. These models have the ability to guide the generation process towards a predefined goal, which provides a compelling direction for EHR generation. Therefore, an intuitive idea is to guide the diffusion process towards generating data that benefits downstream tasks. However, one of the key challenges is how to represent and quantify the impact of the generated data on the performance of downstream models for specific tasks. To address this, we draw inspiration from influence functions[Koh and Liang, 2017, Anand et al., 2023, Cook, 1977], a robust statistics technique that measures the impact of a single data point on an estimator, revealing how observations affect model parameters and predictions. In machine learning, influence functions help understand model behavior, debug predictions, and identify influential training points by tracing a model's predictions back to its training data [Hou et al., 2024].

Building upon this foundation, we propose a diffusion framework that integrates the influence of synthetic samples as a form of guidance, with the goal of generating samples that yield the most positive impact on specific clinical tasks. In our approach, the influence of a generated sample is defined as the reduction in the task-specific loss on a guidance set drawn from the same distribution as the downstream task when the sample is incorporated into the training data. By estimating the influence of intermediate samples during the diffusion process, we leverage the gradient information associated with the influence to steer the generation process toward producing data that are optimally beneficial for downstream tasks. By incorporating this influence gradient into the reverse diffusion process, our method actively guides the generation toward producing samples that are more likely to improve performance of healthcare prediction tasks. Ultimately, the influence mechanism bridges the

gap between data authenticity and clinical utility, ensuring that the generated time series are not only realistic but also tailored to enhance specific healthcare tasks.

To summarize, this paper makes the following contributions: First, we introduce a novel influence mechanism that quantifies the clinical utility of synthetic samples by measuring the expected reduction in task-specific loss. Second, we propose TarDiff, an influence guided diffusion model that integrates this mechanism into the reverse diffusion process, effectively steering generation toward samples that enhance downstream model performance. Third, we validate our approach on multiple clinical tasks, demonstrating that the synthetic medical time series produced by our method are not only realistic but also yield significant improvements in clinical outcomes.

## 2 Preliminary

### 2.1 Diffusion Models

**Denoising Diffusion Probabilistic Models**    The core idea behind DDPM[Ho et al., 2020] is the modeling of a forward diffusion process, which progressively adds noise to the data, and a reverse diffusion process, which denoises the data to ultimately generate realistic samples from noise. Let $x_0$ represent a data sample drawn from the real data distribution $p_{\text{data}}(x_0)$. The forward diffusion process introduces Gaussian noise to the data over $T$ discrete time steps, transforming $x_0$ into pure noise $x_T$, which is assumed to follow a standard normal distribution, i.e., $x_T \sim \mathcal{N}(0, I)$. At each time step $t$, the data $x_t$ is modeled as a noisy version of $x_{t-1}$ via the conditional Gaussian distribution:

$$\begin{aligned}
q(x_t|x_{t-1}) &= \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I), \\
q(x_t|x_0) &= \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I).
\end{aligned} \tag{1}$$

$\beta_t$ denotes the noise schedule controlling the amount of noise added at each step, and $\bar{\alpha}_t = \prod_{s=1}^{t}(1 - \beta_s)$ represents the cumulative noise factor. The forward diffusion process is Markovian, and the joint distribution of the noisy data can be expressed as $q(x_{1:T}|x_0) = \prod_{t=1}^{T} q(x_t|x_{t-1})$, where $x_{1:T} = (x_1, x_2, ..., x_T)$ represents the sequence of noisy variables from $x_0$ to $x_T$.

The reverse diffusion process seeks to recover the original data $x_0$ from the noise $x_T$ by learning a generative model. Specifically, at each time step, the model predicts the mean of the reverse distribution $p_\theta(x_{t-1}|x_t)$, which is also assumed to be Gaussian:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(t)), \tag{2}$$

where $\mu_\theta(x_t, t)$ and $\Sigma_\theta(t)$ represent the model-predicted mean and covariance, respectively, parameterized by $\theta$.

The model is trained by minimizing the following loss function, leveraging the Markov property of the diffusion process:

$$L(\theta) = \mathbb{E}_q\left[\|\hat{\epsilon}_\theta(x_t, t) - \epsilon_t\|^2\right], \tag{3}$$

$\hat{\epsilon}_\theta(x_t, t)$ is the model's prediction of the noise added at time step $t$, and $\epsilon_t$ is the actual noise introduced during the forward diffusion process. This loss function encourages the model to accurately predict the added noise, enabling it to reverse the diffusion process and recover the original data $x_0$.

**Conditional Diffusion**    In conditional diffusion[Ho and Salimans, 2022], the reverse process is explicitly guided by an additional condition $y$(i.e.,class label). Specifically, the probability is modeled as:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, y, t), \Sigma_\theta(t)), \tag{4}$$

where the mean function $\mu_\theta(x_t, y, t)$ incorporates not only the current state $x_t$ and the time step $t$ but also the condition $y$, and $\Sigma_\theta(t)$ denotes the covariance at time $t$. This formulation enables the model to generate samples that adhere to both the learned data distribution and the desired attributes specified by $y$, thereby achieving controlled synthesis.

In the following, we adapt the conditional diffusion framework to time series generation. To make this concrete, we first define the general time series generation problem.

3

**Classifier-Guided Diffusion.** This approach[Dhariwal and Nichol, 2021] augments reverse diffusion with signals from an classifier that estimates the conditional likelihood $p(y \mid x_t)$ at every timestep $t$. The gradient $\nabla_{x_t} \log p(y \mid x_t)$ indicates the direction in sample space that most increases the probability of label $y$; adding this vector therefore nudges the denoising trajectory toward regions consistent with the desired condition.

Formally, the mean of the Gaussian transition in Eq. (4) is replaced by

$$\tilde{\mu}_\theta(x_t, y, t) \; = \; \mu_\theta(x_t, y, t) \; + \; \alpha \, \nabla_{x_t} \log p(y \mid x_t), \tag{5}$$

where the scalar $\alpha$ controls guidance strength. The resulting reverse kernel becomes

$$p_\theta(x_{t-1} \mid x_t) = \mathcal{N}\big(x_{t-1}; \, \tilde{\mu}_\theta(x_t, y, t), \, \Sigma_\theta(t)\big). \tag{6}$$

By favouring states that yield larger $\log p(y \mid x_t)$, classifier guidance systematically steers the generative process toward samples that satisfy the target label, compensating for mismatches in the diffusion model's original conditional distribution.

## 2.2 Task Formulation

Let $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T\}$ be a continuous time series of length $T$, where each observation $\mathbf{x}_t \in \mathbb{R}^d$. Suppose we have a dataset $\mathcal{D}_0 = \{\mathbf{X}^{(i)}\}_{i=1}^N$, composed of $N$ independent sequences drawn from an unknown underlying distribution

$$p_{\text{data}}(\mathbf{X}) \; = \; p_{\text{data}}\big(\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T\big). \tag{7}$$

Our goal is to learn a generative model $p_G(\mathbf{X})$ that approximates $p_{\text{data}}(\mathbf{X})$, enabling us to sample new sequences

$$\hat{\mathbf{X}} = \{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \ldots, \hat{\mathbf{x}}_T\} \quad \text{with} \quad \hat{\mathbf{X}} \sim p_G(\mathbf{X}). \tag{8}$$

A common strategy is to minimize some divergence measure between $p_{\text{data}}$ and $p_G$,

$$\min_{p_G} \; D\big(p_{\text{data}}(\mathbf{X}) \, \big\| \, p_G(\mathbf{X})\big), \tag{9}$$

where $D(\cdot \| \cdot)$ could be the KL divergence. In practice, $p_G$ must capture both local dependencies (e.g., between adjacent time steps) and global trends. Let $S(\mathbf{X})$ denote relevant statistics (e.g., autocorrelation or cross-correlation). Then a suitable generative model should satisfy

$$\mathbb{E}_{\hat{\mathbf{X}} \sim p_G}\big[S(\hat{\mathbf{X}})\big] \; \approx \; \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}}\big[S(\mathbf{X})\big], \tag{10}$$

so that synthetic sequences reflect the essential temporal structures of real data.

For conditional diffusion models, the generative process can be guided by incorporating conditional information $y$, such as class labels, as formulated in Equation (4). This conditioning mechanism allows the model to generate time series that align with specific contexts or constraints while preserving both local dependencies and global trends.

# 3 Methodology

In this section, we describe our TarDiff framework in detail. As can be seen from Figure 1, our method builds upon a conditional diffusion model to generate useful synthetic data through explicitly incorporating task-specific influence signals into the reverse diffusion process.

## 3.1 Influence Formulation

We consider a training dataset $\mathcal{D}_{\text{train}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and a collection of downstream tasks $\mathcal{T} = \{T_1, T_2, \ldots\}$. For any task $T \in \mathcal{T}$, the task-specific parameters are obtained by minimising the empirical loss on $\mathcal{D}_{\text{train}}$:

$$\phi_T^* = \arg\min_\phi \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{train}}} \ell_T\big(\mathbf{x}_i, y_i; \phi\big). \tag{11}$$
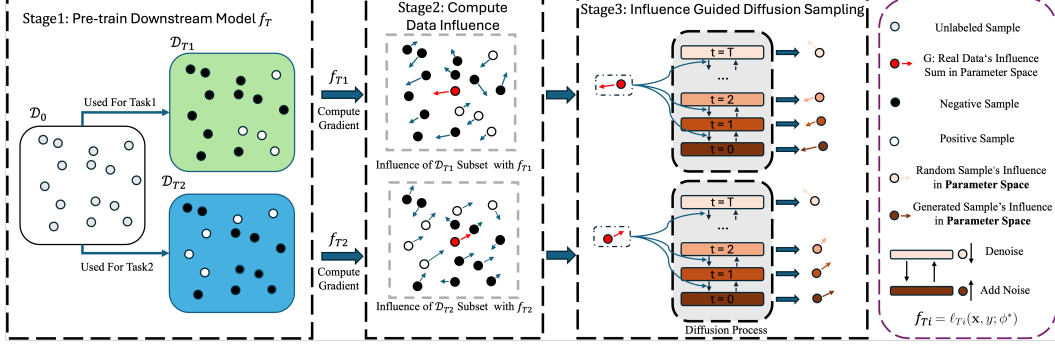
4

Figure 1: **Overview of the Influence Guidance Diffusion framework.** In **Stage 1**, we construct task-specific datasets from the original dataset $\mathcal{D}_{train}$ and train downstream models $f_{T_i}$. In **Stage 2**, we compute each sample's gradient-based influence for total influence $\mathbf{G}$ based on $\mathcal{D}_{T_i}$ and $f_{T_i}$. In **Stage 3**, we leverage influence signals guide the reverse diffusion process with computing $\Delta\mathcal{L}_T(\hat{z}) = \nabla_\phi \ell_T(\mathbf{x}_t, y_t; \phi) \cdot G$. All symbols are detailed in the legend on the right.

To mitigate limitations arising from insufficient or noisy training data, we generate a synthetic sample $\hat{z} = (\mathbf{x}, y)$ and augment the original data, yielding $\mathcal{D}_0 \cup \{\hat{z}\}$, where $\mathcal{D}_0 := \mathcal{D}_{\text{train}}$. Retraining on the augmented set gives

$$\phi_T^{\hat{z}} = \arg\min_\phi \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}_0 \cup \{\hat{z}\}} \ell_T\big(\mathbf{x}_i, y_i; \phi\big). \tag{12}$$

For a samples $(\mathbf{x}', y')$, we define

$$H_\phi\big(\mathbf{x}, y, \mathbf{x}', y'\big) \;=\; \ell_T\big(\mathbf{x}', y'; \phi^{\hat{z}}\big) - \ell_T\big(\mathbf{x}', y'; \phi_T^*\big), \tag{13}$$

the change in downstream loss on $(\mathbf{x}', y')$ caused by adding $(\mathbf{x}, y)$ to the training set.[3]

Let $\mathcal{P}$ denote the underlying data-generating distribution that produces unseen i.i.d. samples at evaluation time. We define the *influence* of $\hat{z}$ on task $T$ as the expected reduction in loss over this distribution:

$$\Delta\mathcal{L}_T(\hat{z}) \;\triangleq\; -\mathbb{E}_{(\mathbf{x}', y') \sim \mathcal{P}}\big[H_\phi\big(\mathbf{x}, y, \mathbf{x}', y'\big)\big] = -\mathbb{E}_{(\mathbf{x}', y') \sim \mathcal{P}}\Big[\ell_T\big(\mathbf{x}', y'; \phi^{\hat{z}}\big) - \ell_T\big(\mathbf{x}', y'; \phi_T^*\big)\Big]. \tag{14}$$

Our goal is to synthesise the sample that maximises this influence:

$$\hat{z}^* = \arg\max_{\hat{z}} \; \Delta\mathcal{L}_T(\hat{z}), \tag{15}$$

thereby ensuring that the generated data yields the greatest expected performance gain on unseen, i.i.d. instances of the target task.

### 3.2 Influence Guidance Diffusion

While our pre-trained conditional diffusion model with parameters $\theta^*$ is adept at generating realistic medical time series dataset $\mathcal{D}_0$, its generation is driven solely by the learned data distribution conditioned on $y$ (e.g., the label), with the reverse diffusion process described by Equation (4). In this standard setting, to generate a time series sample with a specific label, one simply samples from Gaussian noise and iteratively denoises according to (4)—each step conditioned only on the time step $t$ and the label $y$—to yield data consistent with the designated label.

Section 2.1 reviewed how an auxiliary classifier steers the reverse kernel via the gradient $\nabla_{x_t} \log p(y \mid x_t)$, yielding the modified mean $\tilde{\mu}_\theta(x_t, y, t) = \mu_\theta(x_t, y, t) + \alpha \nabla_{x_t} \log p(y \mid x_t)$.

Inspired by classifier guidance, we incorporate an additional control signal that measures the impact of a synthetic sample on downstream performance. Let $\mathcal{D}_{guide} = \{(\mathbf{x}'_j, y'_j)\}_{j=1}^{N_g}$ be a **guidance set** drawn i.i.d. from the same data-generating distribution $\mathcal{P}$ introduced in Section 3.1.

---

[3]Throughout, $\phi^{\hat{z}}$ depends on $(\mathbf{x}, y)$ as in Eq. (2), and $\phi_T^*$ is the original optimum in Eq. (1).

Recalling the point-wise loss change $H_\phi(\cdot)$ defined in Eq. (13), the influence of a sample $\hat{z} = (\mathbf{x}, y)$ on task $T$ is estimated by

$$\Delta\mathcal{L}_T(\hat{z}) \;=\; \sum_{(\mathbf{x}',y')\in\mathcal{D}_{guide}} H_\phi\big(\mathbf{x}, y, \mathbf{x}', y'\big). \tag{16}$$

At each diffusion step $t$ we treat $\hat{z}_t = (x_t, y)$ and replace the classifier-guidance term $\nabla_{x_t} \log p(y \mid x_t)$ with the gradient of the influence estimate, $\nabla_{x_t}\Delta\mathcal{L}_T(\hat{z}_t)$. Note that the constant factor in Eq. (16) is omitted, as it does not affect the direction of the gradient used for guidance.

Consequently, the reverse update becomes

$$\tilde{\mu}_\theta(x_t, y, t) = \mu_\theta(x_t, y, t) \;+\; \alpha\,\nabla_{x_t}\Delta\mathcal{L}_T(\hat{z}_t), \tag{17}$$

where $\alpha$ controls guidance strength.

By steering the denoising trajectory toward samples that exhibit higher $\Delta\mathcal{L}_T(\hat{z})$, we ensure that the generated medical time series data not only conform to the specified condition label but also actively enhance downstream task performance across the guidance set. This approach mitigates spurious correlations that can arise from purely label-conditioned diffusion, as the model explicitly seeks synthetic samples that yield beneficial effects on a broad set of guidance set examples.

## 3.3   Estimates of Influence

In 3.1 and 3.2, we defined the concept of influence for synthesized samples and discussed how to generate samples with the highest possible influence during the generation process. Building on this foundation, the generation of task-specific medical time series requires evaluating the influence of a generated sample on downstream tasks. The primary challenge is to efficiently quantify how synthetic samples impact the model's parameter updates, which are designed to minimize the task-specific loss function $\ell_T$.

A straightforward but computationally expensive approach would involve iteratively adding each candidate synthetic sample $\hat{z}$ to the training set, retraining the model from scratch, and measuring the resultant performance change. Since this process requires $\mathcal{O}(n)$ retraining steps, it leads to prohibitive computational costs, particularly for large-scale datasets and complex models. Thus, optimizing for both efficiency and accuracy in influence estimation is crucial to enable practical applications of this method.

To circumvent the need for exhaustive retraining, a framework was proposed in [Charpiat et al., 2019, Anand et al., 2023] that approximates the parameter shift $\delta\phi$ induced by a synthetic sample $\hat{z}$ via gradient-based analysis, thereby significantly reducing computational overhead while maintaining accuracy in influence estimation.

When estimating the impact of a single sample on model parameters, we build upon the classic conclusion presented in [Charpiat et al., 2019]. As derived in [Charpiat et al., 2019], to change the value of the model's prediction on a sample $\mathbf{x}$, denoted as $f_\phi(\mathbf{x})$, by a small quantity $\varepsilon$, The parameters $\phi$ can be updated according to $\delta\phi = \frac{\varepsilon\nabla_\phi f_\phi(\mathbf{x})}{\|\nabla_\phi f_\phi(\mathbf{x})\|^2}$. After this parameter update, the new prediction at $\mathbf{x}$, $f_{\phi+\delta\phi}(\mathbf{x})$ is given by:

$$\begin{aligned} f_{\phi+\delta\phi}(\mathbf{x}) &= f_\phi(\mathbf{x}) + \nabla_\phi f_\phi(\mathbf{x}) \cdot \delta\phi + \mathcal{O}(\|\delta\phi\|^2) \\ &= f_\phi(\mathbf{x}) + \varepsilon + \mathcal{O}(\varepsilon^2) \end{aligned} \tag{18}$$

where $\mathcal{O}(\|\delta\phi\|^2)$ denotes higher-order terms that are negligible for sufficiently small updates. This foundational perspective allows us to analyze the impact of individual samples on model behavior without necessitating full retraining. Furthermore, we can assess how such a parameter change affects the model's prediction on another sample $\mathbf{x}'$:

$$\begin{aligned} f_{\phi+\delta\phi}(\mathbf{x}') &= f_\phi(\mathbf{x}') + \nabla_\phi f_\phi(\mathbf{x}') \cdot \delta\phi + \mathcal{O}(\|\delta\phi\|^2) \\ &= f_\phi(\mathbf{x}') + \varepsilon \frac{\nabla_\phi f_\phi(\mathbf{x}') \cdot \nabla_\phi f_\phi(\mathbf{x})}{\|\nabla_\phi f_\phi(\mathbf{x})\|^2} + \mathcal{O}(\|\delta\phi\|^2). \end{aligned} \tag{19}$$

The term $\varepsilon \frac{\nabla_\phi f_\phi(\mathbf{x}') \cdot \nabla_\phi f_\phi(\mathbf{x})}{\|\nabla_\phi f_\phi(\mathbf{x})\|^2}$ quantifies the influence of the parameter update caused by $\mathbf{x}$ on the model's prediction at $\mathbf{x}'$. Specifically, the numerator $\nabla_\phi f_\phi(\mathbf{x}') \cdot \nabla_\phi f_\phi(\mathbf{x})$ represents the alignment

between the gradients of $f_\phi$ at $\mathbf{x}'$ and $\mathbf{x}$. A higher alignment indicates that changes in $\phi$ due to $\mathbf{x}$ will have a more pronounced effect on the prediction at $\mathbf{x}'$.

**Influence on Performance.**    To extend this analysis to our setting, we replace the model's prediction $f_\phi(\mathbf{x})$ with a target loss function $\ell_\mathrm{T}(\mathbf{x}, y; \phi)$, transitioning from analyzing the influence on predictions to quantifying the influence on the overall optimization objective. Specifically, the influence of $\mathbf{x}$ on the loss at another sample $\mathbf{x}'$ can be expressed as:

$$
\begin{aligned}
\ell_\mathrm{T}(\mathbf{x}', y'; \phi + \delta\phi) =& \ell_\mathrm{T}(\mathbf{x}', y'; \phi) + \mathcal{O}(\varepsilon^2) \\
& + \varepsilon \frac{\nabla_\phi \ell_\mathrm{T}(\mathbf{x}', y'; \phi) \cdot \nabla_\phi \ell_\mathrm{T}(\mathbf{x}, y; \phi)}{\|\nabla_\phi \ell_\mathrm{T}(\mathbf{x}, y; \phi)\|^2}.
\end{aligned}
\tag{20}
$$

Replacing $\mathbf{x}$ with synthetic data $\hat{z} = (\mathbf{x}, y)$ and aggregating its influence across the guidance set $\mathcal{D}_{guide}$ leads to:

$$
\begin{aligned}
\Delta\mathcal{L}_T(\hat{z}) &= - \sum_{(\mathbf{x}', y') \in \mathcal{D}_{guide}} H_\phi(\mathbf{x}, y, \mathbf{x}', y') \\
&= \sum_{(\mathbf{x}', y') \in \mathcal{D}_{guide}} [\ell_T(\mathbf{x}', y'; \phi) - \ell_T(\mathbf{x}', y'; \phi + \delta\phi)] \\
&= - \sum_{(\mathbf{x}', y') \in \mathcal{D}_{guide}} \varepsilon \frac{\nabla_\phi \ell(\mathbf{x}', y'; \phi) \cdot \nabla_\phi \ell(\hat{z}; \phi)}{\|\nabla_\phi \ell(\hat{z}; \phi)\|^2}.
\end{aligned}
\tag{21}
$$

Notice that we can denote the gradient accumulation of the guidance set $\mathcal{D}_{guide}$ by $\mathbf{G}$, we can end up with the following equation:

$$
\mathbf{G} = - \sum_{(\mathbf{x}', y') \in \mathcal{D}_{guide}} \varepsilon \frac{\nabla_\phi \ell(\mathbf{x}', y'; \phi)}{\|\nabla_\phi \ell(\hat{z}; \phi)\|^2}
\tag{22}
$$

$$
\Delta\mathcal{L}_T(\hat{z}) = \nabla_\phi \ell(\hat{z}; \phi) \cdot G
\tag{23}
$$

Where $\Delta\mathcal{L}_T(\hat{z})$ measures the positive impact of the synthetic sample $\hat{z}$ on the target loss. Specifically, a greater increase in $\Delta\mathcal{L}_T(\hat{z})$ indicates that the model has better optimized its objective function, implying improved generalization to unseen data. This improvement manifests as enhanced performance metrics (e.g., accuracy, AUC, or F1-score) on downstream tasks. By generating synthetic samples that maximize this influence—i.e., by maximizing $\Delta\mathcal{L}_T(\hat{z})$—we aim to augment the dataset with medical time series data that is beneficial for downstream tasks. This process addresses challenges such as data sparsity, class imbalance, and noisy measurements often encountered in medical datasets. By introducing synthetic samples that capture task-relevant patterns (e.g., subtle physiological changes or rare but critical events), the model can learn more representative and clinically significant features, ultimately improving its robustness and reliability in tasks such as mortality prediction and disease diagnosis.

### 3.4   TarDiff Pipeline

In this section, we provide an end-to-end overview of the Influence Guided Diffusion pipeline for generating high-quality medical time-series data customized to a target task $\mathcal{T}$. The complete procedure is illustrated in Algorithm 1, which consists of three main steps: (1) pre-training the downstream model, (2) computing influence gradients, and (3) performing guided diffusion sampling.

**Step 1: Pre-train Downstream Model.** We first optimize the downstream task model $f_\phi$ on the original dataset $\mathcal{D}_{train}$, aiming to find parameters $\phi^*$ that minimize the target loss function $\ell(\cdot; \phi)$ over $(\mathbf{x}, y) \in \mathcal{D}_{train}$. This yields a well-trained model capable of capturing task-specific knowledge relevant to the subsequent generation process.

**Step 2: Compute Data Influence.** Using the trained parameters $\phi^*$, we then compute per-sample gradients $\nabla_\phi \ell(\mathbf{x}_i, y_i; \phi^*)$ for each sample in guidance dataset $\mathcal{D}_{guide}$. Accumulating and normalizing these gradients produces a single vector $\mathbf{G}$ that reflects the aggregated influence of the dataset on the

downstream model. This gradient cache $\mathbf{G}$ is leveraged to guide the diffusion model, ensuring that generated samples maximize their impact on the target task.

**Step 3: Influence-Guided Diffusion Sampling.** With $\mathbf{G}$ fixed, we initialize $\mathbf{x}_T$ from a standard Gaussian. At each reverse diffusion step $t$, the model outputs $\mu_t = \mu_\theta(\mathbf{x}_t, y, t)$. We then compute an influence-driven guidance term, $\mathbf{J} \leftarrow \nabla_{\mathbf{x}_t}\big(\mathbf{G} \cdot \nabla_\phi \ell(\mathbf{x}_t, y; \phi^*)\big)$, where $\ell$ is the downstream loss, and $\mathbf{G}$ encapsulates the influence from the guidance set. We update the mean via $\tilde{\mu}_t \leftarrow \mu_t + w \cdot \mathbf{J}$, where $w$ controls the strength of task-oriented guidance. Finally, we sample $\mathbf{x}_{t-1} \sim \mathcal{N}(\tilde{\mu}_t, \Sigma_\theta(t))$. Iterating this procedure through all diffusion steps yields the final synthetic sample $\hat{z} = \mathbf{x}_0$, biased toward improving downstream performance.

By combining a pre-trained downstream task model, influence gradient aggregation, and guidance-based diffusion sampling, TarDiff provides a unified pipeline for generating task-specific synthetic data. This pipeline ensures the fidelity of generated medical time-series data while aligning it closely with the optimization objective of the target task, making it particularly suitable in scenarios such as diagnostic improvement or risk prediction in healthcare applications.

---

**Algorithm 1** TarDiff Pipeline

---

**Require:** Original dataset $\mathcal{D}_{train}$, guidance subset $\mathcal{D}_{guide}$ Pretrained conditional diffusion model $\mu_\theta(\mathbf{x}_t, y, t), \Sigma_\theta(t)$. Downstream task model $f_\phi$ with random initialization, Loss function $\ell(\cdot; \phi)$, Total diffusion steps $T$, Influence scaling factor $w$.

**Ensure:** Synthetic sample $\hat{z}$ optimizing task $\mathcal{T}$.

    **Step 1: Pre-train downstream model**
    Optimize $\phi^* \leftarrow \arg\min_\phi \sum_{(\mathbf{x}, y) \in \mathcal{D}_{train}} \ell(\mathbf{x}, y; \phi)$
    **Step 2: Compute Data Influence.**
    Initialize $\mathbf{G} \leftarrow 0$
    **for** $i = 1$ to $|\mathcal{D}_{guide}|$ **do**
        Get sample $(\mathbf{x}_i, y_i) \in \mathcal{D}_{guide}$
        Compute per-sample gradient:
        $\mathbf{g}_i \leftarrow \nabla_\phi \ell(\mathbf{x}_i, y_i; \phi^*)$
        Accumulate gradients: $\mathbf{G} \leftarrow \mathbf{G} + \mathbf{g}_i$
    **end for**
    Normalize: $\mathbf{G} \leftarrow \frac{1}{|\mathcal{D}_0|}\mathbf{G}$
    **Step 3: Influence Guided diffusion sampling**
    Initialize $\mathbf{x}_T \sim \mathcal{N}(0, I)$
    **for** $t = T$ to $1$ **do**
        (a) $\mu_t \leftarrow \mu_\theta(\mathbf{x}_t, y, t)$.
        (b) $\mathbf{J} \leftarrow \nabla_{\mathbf{x}_t}\big[\mathbf{G} \cdot \nabla_\phi \ell(\mathbf{x}_t, y_t; \phi^*)\big]$
        (c) $\tilde{\mu}_t \leftarrow \mu_t + w \cdot \mathbf{J}$.
        (d) $\mathbf{x}_{t-1} \sim \mathcal{N}(\tilde{\mu}_t, \Sigma_\theta(t))$.
    **end for**
    Return $\hat{z} \leftarrow (\mathbf{x}_0, y_0)$

---

## 4 Experiment Setup and Result Analysis

Our experiments are designed to systematically evaluate the proposed method in terms of data quality, augmentation effectiveness, influence guidance mechanism, and computational efficiency. Specifically, we investigate: (1) the feasibility of entirely replacing real data with synthetic data (Section 4.2), (2) the effectiveness of synthetic data as augmentation (Section 4.3), (3) the specific impact of the influence guidance mechanism (Section 4.6), and (4) the computational efficiency of our proposed method (Section 4.5). Throughout all experiments, we adopt the standard validation split as the **i.i.d. guidance set** introduced in Section 3.1.

### 4.1 Datasets, Baselines, and Evaluation Metrics

**Datasets.** We evaluate our method across multiple datasets. **MIMIC-III** [Johnson et al., 2016] includes multivariate ICU data (7 features, 24 steps) from 20,920 samples. **eICU** [Pollard et al., 2018] provides ICU data (3 features, 288 steps after preprocessed) from over 200,000 admissions.

Table 1: Performance Comparison of Synthetic Data Generation Methods on MIMICIII and eICU

| Method | MIMIC-III | | | | eICU | | | |
| | Mortality | | ICU Stay | | Mortality | | ICU Stay | |
| | AUPRC | AUROC | AUPRC | AUROC | AUPRC | AUROC | AUPRC | AUROC |
|---|---|---|---|---|---|---|---|---|
| TimeGAN | 0.1402 | 0.5144 | 0.3645 | 0.5431 | 0.1592 | 0.6238 | 0.4213 | 0.4620 |
| TimeVAE | 0.0957 | 0.5392 | 0.3939 | 0.5656 | 0.1046 | 0.5233 | 0.4753 | 0.5321 |
| TimeVQVAE | 0.0874 | 0.5182 | 0.3790 | 0.5389 | 0.1216 | 0.5721 | 0.4520 | 0.5287 |
| DiffusionTS | 0.0865 | 0.5330 | 0.3451 | 0.4946 | 0.1292 | 0.5594 | 0.4692 | 0.5261 |
| BioDiffusion | 0.0964 | 0.5335 | 0.3370 | 0.4905 | 0.1435 | 0.5872 | 0.4625 | 0.5323 |
| Real Data | 0.1736 | 0.6350 | 0.4618 | 0.6282 | 0.2072 | 0.6869 | 0.6004 | 0.6615 |
| TarDiff | 0.1799 | 0.6373 | 0.4183 | 0.5800 | 0.1698 | 0.6308 | 0.5583 | 0.6184 |

Additionally, we test generalizability on four physiological signal datasets: *APAVA* [Escudero et al., 2006], *PTB* [Goldberger et al., 2000], *TDBRAIN* [Van Dijk et al., 2022], and *ADFD* [Miltiadous et al., 2023], covering diverse ECG and EEG signals. All datasets use an 80%-10%-10% split for training, validation, and test sets. Detailed descriptions and preprocessing steps are provided in Appendix A.

**Baselines.** We compare our approach with several state-of-the-art generative methods: *TimeGAN* [Yoon et al., 2019], a GAN-based model; *TimeVAE* [Desai et al., 2021], a variational autoencoder model; *Diffusion-TS* [Yuan and Qiao, 2024], an unconditional diffusion model; *TimeVQ-VAE* [Lee et al., 2023] and *BioDiffusion* [Li et al., 2024a], both conditional generative models capable of directly generating label-conditioned time series. For unconditional models (TimeGAN, TimeVAE, Diffusion-TS), we train separate class-specific models. To evaluate the downstream utility of generated data, we use *TimesNet* [Wu et al., 2022], a state-of-the-art time-series classification architecture, to measure classification performance.

**Evaluation Metrics.** Given the inherent class imbalance in clinical datasets, we use Area Under the Receiver Operating Characteristic Curve (AUROC) and Area Under the Precision-Recall Curve (AUPRC) as primary evaluation metrics. Detailed metric definitions and formulas are provided in Appendix B.

## 4.2 Train on Synthetic, Test on Real (TSTR)

We evaluate the potential of the generated data to serve as a substitute for original data in training high-performance models for clinical tasks. The downstream classifier (`TimesNet` [Wu et al., 2022]) is trained exclusively on synthetic data produced by each time-series generation method, then evaluated on real test data to assess its generalization capability.

Table 1 summarizes the TSTR performance on MIMIC-III and eICU datasets, while the results on high-frequency EEG (*APAVA*, *ADFTD*, *TDBrain*) and ECG (*PTB*) datasets are presented separately in Table 2. Consistently, `TarDiff` achieves state-of-the-art performance in terms of AUROC and AUPRC. Notably, in the TSTR setting, models trained solely on `TarDiff`-generated samples outperform those trained on synthetic data from other baselines across all tasks, including both standard EHR classification and EEG/ECG-based diagnoses. Despite having no access to original data during training, the downstream classifiers attain high scores on real test sets, confirming that `TarDiff`-generated samples effectively capture the essential clinical or physiological features needed for robust predictive modeling.

These findings emphasize that our approach not only produces realistic time-series data but also generates samples that are carefully optimized to improve downstream clinical predictions in a variety of scenarios, thereby validating the robustness and utility of `TarDiff` in critical healthcare applications ranging from general EHR to high-frequency EEG/ECG analytics.

## 4.3 Train on Synthetic and Real, Test on Real (TSRTR)

**Setups.** In this experiment, we investigate the effect of incorporating synthetic EHR data—generated by different time series generation framework—into the training set alongside real data, and subsequently testing the trained model on a real test set. We examine five synthetic-to-real mixing ratios:

Table 2: Performance Comparison of Synthetic Data Generation Methods on ECG and EEG Datasets

| Method | APAVA | | ADFD | | PTB | | TDBrain | |
|---|---|---|---|---|---|---|---|---|
| | AUPRC | AUROC | AUPRC | AUROC | AUPRC | AUROC | AUPRC | AUROC |
| TimeGAN | 0.61229 | 0.51033 | 0.35839 | 0.54113 | 0.86766 | 0.78173 | 0.60932 | 0.62780 |
| TimeVAE | 0.74266 | 0.68569 | 0.42466 | 0.60053 | 0.95092 | 0.89445 | 0.58464 | 0.58565 |
| TimeVQVAE | 0.63722 | 0.55500 | 0.33884 | 0.50578 | 0.94862 | 0.89843 | 0.51336 | 0.54659 |
| DiffusionTS | 0.57083 | 0.47062 | 0.33286 | 0.50039 | 0.87372 | 0.75979 | 0.49900 | 0.49319 |
| BioDiffusion | 0.63294 | 0.54325 | 0.46042 | 0.63753 | 0.85887 | 0.74805 | 0.55134 | 0.56796 |
| Real Data | 0.76692 | 0.72063 | 0.43349 | 0.62390 | 0.96768 | 0.93306 | 0.96424 | 0.96153 |
| TarDiff | 0.76519 | 0.77097 | 0.47950 | 0.64429 | 0.95435 | 0.90532 | 0.65420 | 0.64444 |

0.2, 0.4, 0.6, 0.8, and 1.0, where the ratio indicates the size of the synthetic dataset relative to the real training set. For each ratio $\alpha$, the combined training set is defined as

$$\mathcal{D}_{\text{train}} = \mathcal{D}_{\text{real}} \cup \mathcal{D}_{\text{synthetic}}(\alpha),$$

with $\mathcal{D}_{\text{synthetic}}(\alpha)$ representing the synthetic data sampled to match the specified proportion $\alpha$.

We evaluate the downstream performance on six datasets using key metrics AUROC. Figure 2 illustrates the AUROC performance of various methods across two datasets under different synthetic-to-real data mix ratios. Our method demonstrates superior performance, especially at higher mix ratios.
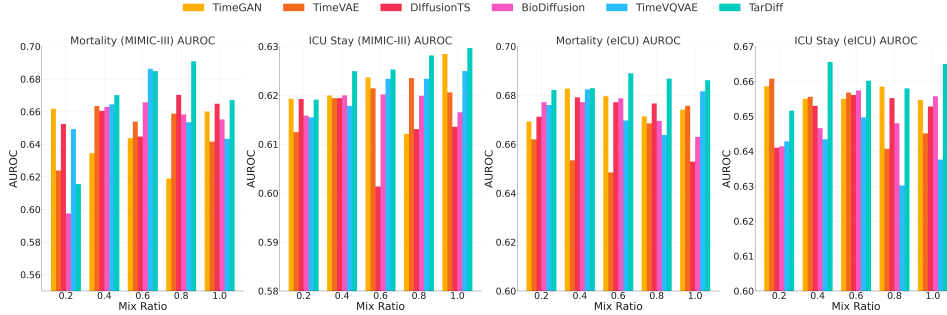


Figure 2: Comparison of AUROC values for the Mortality and ICU Stay task on the MIMIC III and eICU dataset with synthetic-to-real data mix ratios from 0.2 to 1.0.

**Results.** The results indicate that TarDiff outperforms other time series generation baselines across the majority of tasks and mixing ratios, demonstrating a generally upward trend as the synthetic proportion increases from 0.2 to 1.0. Although minor fluctuations occur, TarDiff-derived samples exhibit a sustained positive impact on model performance by leveraging real data guidance throughout the generation process. This synergy enables the synthetic data to capture clinically relevant features more effectively, thereby improving outcomes on these datasets. In contrast, while some baselines occasionally show improvements, their performance tends to degrade or become inconsistent at higher mixing ratios, suggesting a limited capacity to maintain useful signal in purely synthetic settings. Overall, these findings highlight TarDiff's robust ability to generate synthetic samples that enhance downstream performance and reinforce clinically meaningful patterns when integrated with real data.

Overall, these findings underscore the potential of leveraging synthetic EHR time series data to supplement limited real-world datasets, thereby boosting performance on critical clinical tasks.

## 4.4 Influence Guidance under Class Imbalance

This subsection investigates whether influence-guided diffusion mitigates the inherent label imbalance in clinical prediction tasks on MIMIC-III and eICU.

**Gradient Analysis.** We first measure $\ell_2$-norms of gradients obtained from a pretrained TimesNet mortality classifier. Table 3 indicates that minority samples ($\approx$9–11% of instances) exhibit substan-

tially larger gradient magnitudes than majority samples, suggesting these cases are intrinsically harder to classify and thus receive stronger guidance signals.

Table 3: Mean gradient norms ($\pm$ std) for majority vs. minority samples.

| Dataset | Majority | Minority |
|---|---|---|
| MIMIC-III | $1.06 \pm 1.23$ | $\mathbf{16.85 \pm 2.48}$ |
| eICU | $5.41 \pm 5.05$ | $\mathbf{37.86 \pm 8.73}$ |

**Minority-Class Performance.** Table 4 reports minority-class $F_1$ scores when real data are augmented with synthetic data generated by TarDiff (TSRTR protocol). Compared with a real-only baseline, TarDiff more than doubles the minority $F_1$ on MIMIC-III (+93%) and improves eICU by 44%, demonstrating that unified influence guidance already alleviates imbalance without explicit class weighting.

Table 4: Minority-class $F_1$ under different generation strategies.

| Method | MIMIC-III | eICU |
|---|---|---|
| TRTR (Real-only) | 0.056 | 0.013 |
| TarDiff | **0.108** | **0.018** |

**Class-Specific Guidance.** To further isolate the effect of guidance signals, gradients are recomputed on (i) all guidance samples, (ii) majority-only samples, and (iii) minority-only samples. Table 5 shows that minority-only guidance attains the highest minority $F_1$ (0.163 on MIMIC-III; 0.025 on eICU), while majority-only guidance degrades performance.

Table 5: Minority-class $F_1$ with class-specific guidance.

| Guidance Source | MIMIC-III | eICU |
|---|---|---|
| All Samples | 0.108 | 0.018 |
| Majority-only | 0.066 | 0.012 |
| **Minority-only** | **0.163** | **0.025** |

**Discussion.** The pronounced gradient disparity (Table 3) indicates that influence guidance naturally focuses on under-represented events. Consequently, TarDiff improves rare-class metrics without additional hyper-parameters and retains flexibility to apply targeted gradients for further gains, providing a principled mechanism for alleviating class imbalance in medical time-series generation.

## 4.5 Complexity Analysis.

We analyze the additional computational costs introduced by our gradient-guided diffusion framework relative to a standard diffusion-based generation pipeline. First, consider the one-time overhead of training a downstream task network (e.g., a classifier or regressor) and caching its gradients over a target set. Specifically, the downstream model must be trained until convergence, followed by forward-backward passes on the target set to obtain gradient norms with respect to the downstream loss. If $N_t$ is the size of the target set, $f_T(L, D)$ represents the complexity of a single forward-backward pass for the downstream model on time-series data of length $L$ and feature dimensionality $D$, and if these gradients are stored once for reuse, the overall cost of this stage can be approximated by

$$O\big(N_t \cdot f_T(L, D)\big). \tag{24}$$

In practice, on eight real-world datasets, we measure this one-time overhead to range from 10s to 167s; please see the Appendix C for the detailed statistics.

Next, during sampling, we incorporate gradient guidance by projecting intermediate samples $\mathbf{x}_t$ onto directions derived from the cached gradient norms. Because the downstream network is not

11

re-invoked at each diffusion step, the extra per-step overhead is limited to dot products and final gradient computations, denoted by $g(L, D)$. Hence, for $T$ diffusion steps and a batch size $B_{\text{sample}}$, the cost of gradient-guided sampling is

$$O\big(T \cdot B_{\text{sample}} \cdot g(L, D)\big). \tag{25}$$

In comparison, a standard diffusion framework (e.g., DDPM) typically incurs a sampling cost of

$$O\big(T \cdot B_{\text{sample}} \cdot h(L, D)\big), \tag{26}$$

where $h(L, D)$ is the computational complexity for each step without gradient guidance. Consequently, the additional overhead ratio can be approximated by

$$\frac{T \cdot B_{\text{sample}} \cdot g(L, D)}{T \cdot B_{\text{sample}} \cdot h(L, D)} = \frac{g(L, D)}{h(L, D)}, \tag{27}$$

which is typically small because $g(\cdot)$ involves only lightweight vector or matrix operations. Moreover, we compare our sampling speed with other baselines on a diffusion-based model and observe faster sampling for our gradient-guided approach. Full results can be found in the Appendix C. Overall, our gradient guidance requires a one-time downstream network training plus negligible extra work at each sampling step, yet yields a substantial improvement in aligning generated samples with downstream tasks—an important advantage in applications such as medical time-series data.

### 4.6   Sample Influence with Performance

To rigorously assess the effectiveness of our proposed TarDiff, we conduct experiments targeting two objectives: (1) investigate whether TarDiff can successfully modulate the influence $\Delta \mathcal{L}_T(\hat{z})$ of generated samples; and (2) evaluate how these influence adjustments affect model performance in downstream classification tasks.

Specifically, we partition the original validation set into two subsets: (i) **Guidance-Val subset**, used exclusively for generating samples during the guidance diffusion process, and (ii) **Evaluation-Val subset**, used for selecting the optimal guidance scale by assessing downstream task performance. After selecting the optimal guidance scale, we report the final model performance on the **entire validation set**, thus ensuring an unbiased assessment of TarDiff's generalization capabilities.

Figure 3 illustrates experimental outcomes across various influence scales using samples generated from the *Guidance-Val subset*, with performance metrics evaluated on the *Evaluation-Val subset*. The figure comprises two panels: (i) the left panel depicts changes in sample influence values for both Mortality and ICU Stay tasks, and (ii) the right panel reports the corresponding AUROC performance for these tasks.

As shown in the left panel of Figure 3, varying the influence scale from $-1000$ to $+1000$ markedly affects sample influence values. For the Mortality task, the sample influence consistently declines as the scale increases; conversely, the ICU Stay task exhibits the opposite trend. Correspondingly, in the right panel, these influence adjustments yield improvements in AUROC performance, with both tasks reaching optimal performance at moderate-to-high influence scales (around $+1000$).

By separately tuning the guidance scale using clearly defined subsets of the validation data and subsequently evaluating the final performance on the entire validation set, we ensure the reported results accurately reflect TarDiff's effectiveness and generalization capabilities.

## 5   Related work

Recent advances in time series generation leverage deep generative models such as GANs, VAEs, and diffusion-based approaches [Yoon et al., 2019, Desai et al., 2021, Lee et al., 2023, Huang et al., 2025]. TimeGAN [Yoon et al., 2019] combines adversarial training with supervised embedding, thereby aligning the generated sequences with the real data's temporal structure. Meanwhile, TimeVAE [Desai et al., 2021] and TimeVQVAE [Lee et al., 2023] adopt latent representations to capture salient patterns, which helps improve both the reconstruction quality and overall fidelity of the synthetic time series. Diffusion-based models like TimeDP [Huang et al., 2025] further enhance realism by incorporating domain prompts in their denoising process, enabling more accurate generation of
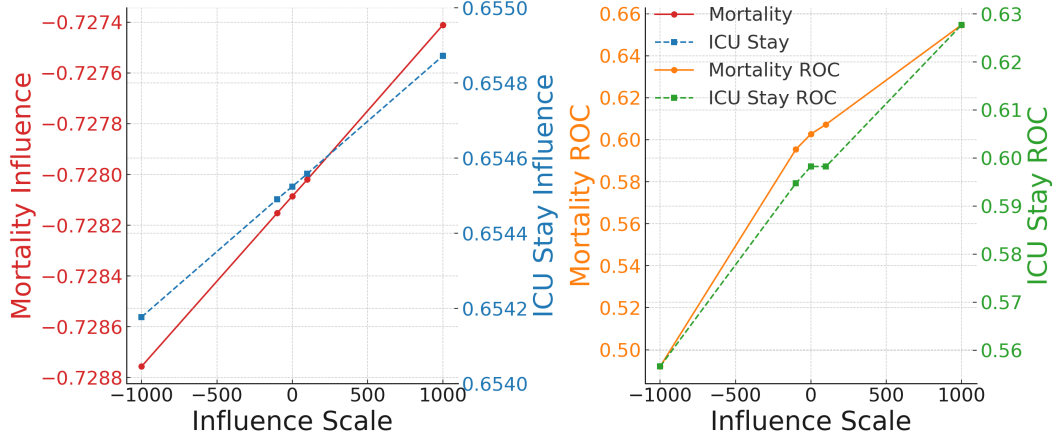
Figure 3: Influence scale analysis conducted by generating samples from the *Guidance-Val subset* and assessing AUROC performance on the *Guidance-Val subset* for Mortality and ICU Stay tasks, with scales ranging from -1000 to 1000. The left panel illustrates sample influence value changes, while the right panel shows AUROC performance across different scales.

complex temporal signals. Empirically, these approaches have demonstrated credible temporal fidelity in diverse domains such as finance [Huang et al., 2024] and medicine [Chen et al., 2024], leading to promising applications in areas like data augmentation, anomaly detection, and privacy-preserving analytics.

In healthcare applications, generating synthetic patient time series presents unique challenges, including privacy protection and the need for clinically relevant patterns. GAN-based methods have been widely used for realistic EHR synthesis [Choi et al., 2017], and diffusion models have also emerged as a promising approach, demonstrating efficacy in producing high-fidelity synthetic EHR time-series data [Tian et al., 2024, Karami et al., 2024]. Conditioning on clinical variables, as seen in methods like MEGAN [Chen et al., 2022], enables high-fidelity, multi-perspective ECG generation, thereby facilitating applications like data augmentation and simulation-based scenario testing in clinical research.

# 6 Conclusion.

In this paper, we present a task-based framework for electronic medical record time series generation that guides diffusion in generating synthetic data by estimating the impact of synthetic samples on specific downstream task models, maximising the impact of synthetic data on clinical tasks. Comprehensive experiments on six datasets demonstrate that our framework not only improves the influence of the generated data on the target task. but also significantly enhances downstream model performance. These results underscore the potential of TarDiff to mitigate data scarcity and privacy concerns in healthcare.

# References

Nikhil Anand, Joshua Tan, and Maria Minakova. Influence scores at scale for efficient language data sampling. *arXiv preprint arXiv:2311.16298*, 2023.

Guillaume Charpiat, Nicolas Girard, Loris Felardos, and Yuliya Tarabalka. Input similarity from the neural network perspective. *Advances in Neural Information Processing Systems*, 32, 2019.

Jiabo Chen, Yongfan Lai, Deyun Zhang, Yue Wang, Shijia Geng, Hongyan Li, and Shenda Hong. Diffusets: 12-lead ecg generation conditioned on clinical text reports and patient-specific information. In *Artificial Intelligence and Data Science for Healthcare: Bridging Data-Centric AI and People-Centric Healthcare*, 2024.

Jintai Chen, Kuanlun Liao, Kun Wei, Haochao Ying, Danny Z Chen, and Jian Wu. Me-gan: Learning panoptic electrocardio representations for multi-view ecg synthesis conditioned on heart diseases. In *International Conference on Machine Learning*, pages 3360–3370. PMLR, 2022.

Edward Choi, Siddharth Biswal, Bradley A. Malin, Jon Duke, Walter F. Stewart, and Jimeng Sun. Generating multi-label discrete electronic health records using generative adversarial networks. *CoRR*, abs/1703.06490, 2017.

R Dennis Cook. Detection of influential observation in linear regression. *Technometrics*, 19(1):15–18, 1977.

Abhyuday Desai, Cynthia Freeman, Zuhui Wang, and Ian Beaver. Timevae: A variational auto-encoder for multivariate time series generation. *arXiv preprint arXiv:2111.08095*, 2021.

Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.

J Escudero, Daniel Abásolo, Roberto Hornero, Pedro Espino, and Miguel López. Analysis of electroencephalograms in alzheimer's disease patients with multiscale entropy. *Physiological measurement*, 27(11):1091, 2006.

Xinyao Fan, Yueying Wu, Chang Xu, Yuhao Huang, Weiqing Liu, and Jiang Bian. MG-TSD: multi-granularity time series diffusion models with guided learning process. In *ICLR*. OpenReview.net, 2024.

A Goldberger, L Amaral, L Glass, J Hausdorff, PC Ivanov, R Mark, HE Stanley, and PhysioToolkit PhysioBank. Physionet: Components of a new research resource for complex physiologic signals components of a new research resource for complex physiologic signals. *Circulation*, 101:e215–e220, 2000.

Benjamin A Goldstein, Ann Marie Navar, Michael J Pencina, and John PA Ioannidis. Opportunities and challenges in developing risk prediction models with electronic health records data: a systematic review. *Journal of the American Medical Informatics Association: JAMIA*, 24(1):198, 2016.

Aman Gupta, Deepak Bhatt, and Anubha Pandey. Transitioning from real to synthetic data: Quantifying the bias in model. *arXiv preprint arXiv:2105.04144*, 2021.

Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

Min Hou, Yueying Wu, Chang Xu, Yu-Hao Huang, Chenxi Bai, Le Wu, and Jiang Bian. Invdiff: Invariant guidance for bias mitigation in diffusion models. *CoRR*, abs/2412.08480, 2024.

Yu-Hao Huang, Chang Xu, Yang Liu, Weiqing Liu, Wu-Jun Li, and Jiang Bian. Controllable financial market generation with diffusion guided meta agent. *arXiv preprint arXiv:2408.12991*, 2024.

Yu-Hao Huang, Chang Xu, Yueying Wu, Wu-Jun Li, and Jiang Bian. Timedp: Learning to generate multi-domain time series with domain prompts. *arXiv preprint arXiv:2501.05403*, 2025.

Zepeng Huo, Xiaoning Qian, Shuai Huang, Zhangyang Wang, and Bobak J Mortazavi. Density-aware personalized training for risk prediction in imbalanced medical data. In *Machine Learning for Healthcare Conference*, pages 101–122. PMLR, 2022.

Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.

Hojjat Karami, Mary-Anne Hartley, David Atienza, and Anisoara Ionescu. Timehr: Image-based time series generation for electronic health records. *CoRR*, abs/2402.06318, 2024.

Shruti Kaushik, Abhinav Choudhury, Pankaj Kumar Sheron, Nataraj Dasgupta, Sayee Natarajan, Larry A. Pickett, and Varun Dutt. AI in healthcare: Time-series forecasting using statistical, neural, and ensemble architectures. *Frontiers Big Data*, 3:4, 2020.

Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 1885–1894. PMLR, 2017.

Daesoo Lee, Sara Malacarne, and Erlend Aune. Vector quantized time series generation with a bidirectional prior model. *arXiv preprint arXiv:2303.04743*, 2023.

Hao Li, Yuping Wu, Viktor Schlegel, Riza Batista-Navarro, Thanh-Tung Nguyen, Abhinav Ramesh Kashyap, Xiao-Jun Zeng, Daniel Beck, Stefan Winkler, and Goran Nenadic. Team: PULSAR at probsum 2023: PULSAR: pre-training with extracted healthcare terms for summarising patients' problems and data augmentation with black-box large language models. In *BioNLP@ACL*, pages 503–509. Association for Computational Linguistics, 2023.

Hao Li, Yu-Hao Huang, Chang Xu, Viktor Schlegel, Ren-He Jiang, Riza Batista-Navarro, Goran Nenadic, and Jiang Bian. Bridge: Bootstrapping text to control time-series generation via multi-agent iterative optimization and diffusion modelling. *arXiv preprint arXiv:2503.02445*, 2025.

Xiaomin Li, Mykhailo Sakevych, Gentry Atkinson, and Vangelis Metsis. Biodiffusion: A versatile diffusion model for biomedical signal synthesis. *CoRR*, abs/2401.10282, 2024a.

Yizhi Li, Ge Zhang, Xingwei Qu, Jiali Li, Zhaoqun Li, Noah Wang, Hao Li, Ruibin Yuan, Yinghao Ma, Kai Zhang, Wangchunshu Zhou, Yiming Liang, Lei Zhang, Lei Ma, Jiajun Zhang, Zuowen Li, Wenhao Huang, Chenghua Lin, and Jie Fu. Cif-bench: A chinese instruction-following benchmark for evaluating the generalizability of large language models. In *ACL (Findings)*, pages 12431–12446. Association for Computational Linguistics, 2024b.

Yuanyuan Liang, Yanbing Ju, Xiao-Jun Zeng, Hao Li, Peiwu Dong, and Tian Ju. A user-generated content-based social network large-scale group decision-making approach in healthcare service: Case study of general practitioners selection in uk. *Expert Systems with Applications*, page 125542, 2024.

Barbara Mukami Maweu, Rittika Shamsuddin, Sagnik Dakshit, and Balakrishnan Prabhakaran. Generating healthcare time series data for improving diagnostic accuracy of deep neural networks. *IEEE Trans. Instrum. Meas.*, 70:1–15, 2021.

Andreas Miltiadous, Katerina D Tzimourta, Theodora Afrantou, Panagiotis Ioannidis, Nikolaos Grigoriadis, Dimitrios G Tsalikakis, Pantelis Angelidis, Markos G Tsipouras, Euripidis Glavas, Nikolaos Giannakeas, et al. A dataset of scalp eeg recordings of alzheimer's disease, frontotemporal dementia and healthy subjects from routine eeg. *Data*, 8(6):95, 2023.

Aishik Nagar, Viktor Schlegel, Thanh-Tung Nguyen, Hao Li, Yuping Wu, Kuluhan Binici, and Stefan Winkler. Llms are not zero-shot reasoners for biomedical information extraction. *CoRR*, abs/2408.12249, 2024.

Tom J Pollard, Alistair EW Johnson, Jesse D Raffa, Leo A Celi, Roger G Mark, and Omar Badawi. The eicu collaborative research database, a freely available multi-center database for critical care research. *Scientific data*, 5(1):1–13, 2018.

Viktor Schlegel, Hao Li, Yuping Wu, Anand Subramanian, Thanh-Tung Nguyen, Abhinav Ramesh Kashyap, Daniel Beck, Xiao-Jun Zeng, Riza Theresa Batista-Navarro, Stefan Winkler, and Goran Nenadic. PULSAR at mediqa-sum 2023: Large language models augmented by synthetic dialogue convert patient dialogues to medical records. In *CLEF (Working Notes)*, volume 3497 of *CEUR Workshop Proceedings*, pages 1668–1679. CEUR-WS.org, 2023.

Julian Schön, Raghavendra Selvan, Lotte Nygård, Ivan Richter Vogelius, and Jens Petersen. Explicit temporal embedding in deep generative latent models for longitudinal medical image synthesis. *CoRR*, abs/2301.05465, 2023.

Benjamin Shickel, Patrick James Tighe, Azra Bihorac, and Parisa Rashidi. Deep ehr: a survey of recent advances in deep learning techniques for electronic health record (ehr) analysis. *IEEE journal of biomedical and health informatics*, 22(5):1589–1604, 2017.

Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.

Muhang Tian, Bernie Chen, Allan Guo, Shiyi Jiang, and Anru R Zhang. Reliable generation of privacy-preserving synthetic electronic health record time series via diffusion models. *Journal of the American Medical Informatics Association*, 31(11):2529–2539, 2024.

Tzu-Wei Tseng, Chang-Fu Su, and Feipei Lai. Fast healthcare interoperability resources for inpatient deterioration detection with time-series vital signs: Design and implementation study. *JMIR Medical Informatics*, 10(10):e42429, 2022.

Hanneke Van Dijk, Guido Van Wingen, Damiaan Denys, Sebastian Olbrich, Rosalinde Van Ruth, and Martijn Arns. The two decades brainclinics research archive for insights in neurophysiology (tdbrain) database. *Scientific data*, 9(1):333, 2022.

Yihe Wang, Nan Huang, Taida Li, Yujun Yan, and Xiang Zhang. Medformer: A multi-granularity patching transformer for medical time-series classification. *arXiv preprint arXiv:2405.19363*, 2024.

Qingsong Wen, Liang Sun, Fan Yang, Xiaomin Song, Jingkun Gao, Xue Wang, and Huan Xu. Time series data augmentation for deep learning: A survey. *arXiv preprint arXiv:2002.12478*, 2020.

Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. *arXiv preprint arXiv:2210.02186*, 2022.

Jinsung Yoon, Daniel Jarrett, and Mihaela Van der Schaar. Time-series generative adversarial networks. *Advances in neural information processing systems*, 32, 2019.

Xinyu Yuan and Yan Qiao. Diffusion-ts: Interpretable diffusion for general time series generation. *arXiv preprint arXiv:2403.01742*, 2024.

# A  Dataset Details

## A.1  Critical Care EHR Datasets (MIMIC-III and eICU)

**MIMIC-III**[Johnson et al., 2016] is a large, publicly available database comprising de-identified health-related data associated with over 40,000 patients who stayed in critical care units of the Beth Israel Deaconess Medical Center between 2001 and 2012 . For our analysis, we focus on the first 24 hours of hospitalization for each patient, resulting in 20,920 samples. Each sample is a 24-step multivariate time series with 7 features: *heart rate*, *systolic blood pressure*, *diastolic blood pressure*, *mean blood pressure*, *respiration rate*, *temperature*, and *oxygen saturation*. We split this dataset into training, validation, and test sets with proportions of 80%, 10%, and 10%, respectively.

**eICU**[Pollard et al., 2018] is a multi-center critical care dataset comprising de-identified health data from over 200,000 admissions to intensive care units (ICUs) across the United States between 2014 and 2015. For our analysis, we extract time-series measurements (*heart rate*, *respiratory rate*, and *oxygen saturation*) from the initial 24-hour window of ICU admission. Data are sampled every 5 minutes, resulting in 288 time steps. Each time step includes 3 features, providing a granular view of patient status. The dataset is partitioned into training, validation, and test sets using an 80%, 10%, 10% ratio.

**Tasks.** We evaluate our framework on two tasks: *(i) Mortality Prediction:* Determine whether a patient will die during the hospital stay. In the MIMIC-III dataset, the positive-to-negative ratio is 1,680 : 19,240, while in the eICU dataset, it is 3,173 : 27,892. *(ii) ICU Length-of-Stay Prediction:* Predict whether a patient's ICU stay exceeds three days. For the MIMIC-III dataset, the positive-to-negative ratio is 2,869 : 18,051, and for the eICU dataset, it is 13,206 : 17,859.

Table 6 provides detailed information on training label distributions within each dataset. Given the inherent label imbalance in both MIMIC-III and eICU, generating clinically meaningful synthetic data remains a significant challenge.

Table 6: MIMIC-III and eICU Dataset Overview.

| Dataset | Task | All Samples | Negative Samples | Positive Samples | Features | Seq Length |
|---------|------|-------------|------------------|------------------|----------|------------|
| MIMIC | Mortality | **20,920** | 19,240 | 1,680 | 7 | 24 |
|  | ICU Stay | **20,920** | 18,051 | 2,869 | 7 | 24 |
| eICU | Mortality | **31,065** | 27,892 | 3,173 | 3 | 288 |
|  | ICU Stay | **31,065** | 17,859 | 13,206 | 3 | 288 |

## A.2 Specialized Physiological Signal Datasets (EEG and ECG)

For these four datasets, we followed the preprocessing and data partitioning settings (training, validation, and test splits) described in [Wang et al., 2024]. Below, we provide a brief introduction to each dataset and their associated tasks.

**APAVA** dataset is a public EEG time series dataset comprising recordings from 23 subjects, including 12 Alzheimer's disease (AD) patients and 11 healthy controls. Each subject underwent approximately 30 trials, with each trial consisting of a 5-second EEG recording sampled at 256Hz across 16 channels. The task associated with APAVA is binary classification, distinguishing AD patients from healthy individuals.

**ADFTD** (Alzheimer's Disease and Frontotemporal Dementia) dataset is an EEG dataset specifically curated to study Alzheimer's disease (AD) and frontotemporal dementia (FTD). It includes EEG recordings from subjects diagnosed with AD, subjects diagnosed with FTD, and healthy controls. EEG signals are recorded across multiple channels at approximately 500Hz and undergo standard preprocessing steps such as filtering and downsampling. The task for ADFTD is a three-class classification to differentiate among AD, FTD, and healthy subjects.

**PTB** Diagnostic ECG Database is a publicly available collection of 549 high-resolution 15-lead ECG recordings from 290 subjects, aged between 17 and 87 years. Each recording includes the standard 12 leads along with 3 Frank leads (Vx, Vy, Vz), digitized at 1000Hz with 16-bit resolution over a ±16.384 mV range. The database encompasses a variety of cardiac conditions, including myocardial infarction, cardiomyopathy, and bundle branch block, as well as recordings from healthy controls. The task on the PTB dataset is binary classification between patients diagnosed with myocardial infarction and healthy controls.

**TDBrain** dataset comprises EEG recordings from multiple channels collected from subjects performing eye-closed tasks. The dataset includes EEG data from subjects diagnosed with Parkinson's disease and healthy controls. The associated task is binary classification to distinguish Parkinson's disease patients from healthy individuals.

## A.3 Scalability Across Different Dataset Sizes

We conducted experiments across multiple datasets varying significantly in sample size, channel numbers, and sequence lengths. This diversity enables us to evaluate the scalability and robustness of our model comprehensively. The datasets cover diverse medical signals, including EHR (MIMIC-III, eICU), ECG (PTB), and EEG (ADFD, APAVA, TDBRAIN), with varying complexity and dimensionality.

As summarized in Table 7, our method consistently demonstrates solid performance improvements across datasets of different sizes and modalities.

## A.4 Controllability of Guidance Set Scale

To further demonstrate the scalability and controllability of our guidance mechanism, we conducted an experiment analyzing the impact of varying guidance set sizes on model performance using the PTBrain dataset. As illustrated in Figure 4, despite fluctuations, the model's performance remains consistently stable across different guidance set sizes. This indicates that our proposed method

Table 7: Dataset scales summary

| Dataset | Samples | Channels | Length |
|---|---|---|---|
| eICU | 31,065 | 3 | 288 |
| MIMIC-III | 26,150 | 7 | 24 |
| ADFD | 69,752 | 19 | 256 |
| PTB | 64,356 | 15 | 288 |
| TDBRAIN | 6,240 | 33 | 256 |
| APAVA | 5,967 | 16 | 256 |

provides effective control over the guidance scale, allowing users to flexibly adjust it according to practical dataset constraints or computational resources. Consequently, this mitigates potential concerns regarding the scalability and practical applicability of our approach.



Figure 4: Model performance with varying guidance set sizes on PTBrain dataset.

# B    Evaluation Metrics Details

**Area Under the Receiver Operating Characteristic Curve (AUROC)** measures the ability of a classifier to distinguish between classes. It is calculated as follows:

$$\text{AUROC} = \int_0^1 \text{TPR}(x) \, d(\text{FPR}(x))$$

where TPR is the True Positive Rate and FPR is the False Positive Rate across different decision thresholds.

**Area Under the Precision-Recall Curve (AUPRC)** evaluates the classifier's performance in imbalanced classification problems by considering precision and recall:

$$\text{AUPRC} = \int_0^1 \text{Precision}(x) \, d(\text{Recall}(x))$$

where precision is the fraction of relevant instances among retrieved instances, and recall is the fraction of relevant instances retrieved over the total relevant instances.

These metrics provide a comprehensive evaluation of classification models, particularly useful in healthcare scenarios with significant class imbalance.

## C  Runtime and Overhead Comparison

As briefly discussed in Section 4.5, our approach introduces a one-time overhead for training and gradient caching, as well as a per-step overhead during sampling. This appendix details the runtime measurements for both the one-time overhead and the sampling process.

### C.1  One-time Overhead

Table 8 shows the one-time cost for gradient caching on eight different datasets, which ranges from 10s to 167s.

Table 8: One-time cost for gradient caching.

| Dataset/Task | Overhead (s) |
|---|---|
| TDBRAIN | 10.51 |
| APAVA | 15.49 |
| ADFD | 167.81 |
| PTB | 144.65 |
| eICU_mortality | 34.01 |
| eICU_ICUStay | 33.93 |
| MIMIC_mortality | 27.51 |
| MIMIC_ICUStay | 27.01 |

### C.2  Sampling Runtime Comparison

To evaluate the sampling efficiency, we compare our TarDiff against representative GAN-based, VAE-based, and Diffusion-based approaches. As seen in Table 9, our method achieves competitive or faster sampling compared to other diffusion-based methods.

Table 9: Sampling runtime comparison across different generation methods.

| Backbone Type | Method | Sampling Time (s/sample) |
|---|---|---|
| GAN-based | TimeGAN | 0.0005 |
| VAE-based | TimeVQVAE | 0.0047 |
| | TimeVQE | 0.0006 |
| Diffusion-based | BioDiffusion | 0.3008 |
| | DiffusionTS | 0.1340 |
| | **TarDiff** | **0.0259** |

## D  Model Structure and Implementation Details

The denoising network in our diffusion model employs a one-dimensional U-Net architecture specifically designed for multi-channel time-series data. The model initializes with 64 channels and features multiple resolution levels, each comprising three residual blocks. We apply progressive channel multipliers of [1, 2, 4, 4] to enhance feature representation at coarser resolutions. To effectively capture long-range temporal dependencies, attention mechanisms with eight heads are incorporated at resolutions of 1, 2, and 4. Additionally, the model integrates scale-shift normalization and residual connections for up-sampling and down-sampling to stabilize training. Contextual embeddings are projected into a 32-dimensional latent space. Furthermore, the architecture supports classifier-free guidance and spatial transformations, optimizing its performance for classification tasks.

Training was conducted using a batch size of 256 for 20,000 iterations with a fixed learning rate of 0.0001. During sampling, we consistently employed a guidance scale of 100 for generated samples. All experiments were carried out on a single NVIDIA A100 GPU with 80GB of memory.

# E    Additional Results on Fidelity & Privacy

Table 10: Distribution Similarity (DS) ↓

| | MIMIC | | eICU | |
| --- | --- | --- | --- | --- |
| **Method** | Mortality | ICU Stay | ICU Stay | Mortality |
| TarDiff | 0.000201 | 0.000000 | 0.1488 | 0.1810 |
| TimeGAN | 0.000201 | 0.000000 | 0.3781 | 0.2471 |
| TimeVAE | 0.000000 | 0.0304 | 0.3668 | 0.1709 |
| TimeVQ-VAE | 0.0325 | 0.0015 | 0.0000 | 0.3663 |
| DiffusionTS | 0.000101 | 0.4451 | 0.5000 | 0.4931 |
| BioDiffusion | 0.000000 | 0.4989 | 0.4968 | 0.5000 |

Table 11: Membership Inference Risk (MIR) ↓

| | MIMIC | | eICU | |
| --- | --- | --- | --- | --- |
| **Method** | Mortality | ICU Stay | ICU Stay | Mortality |
| TarDiff | 0.6761 | 0.6787 | 0.6667 | 0.6667 |
| BioDiffusion | 0.7316 | 0.8114 | 0.7736 | 0.8199 |
| TimeVQ-VAE | 0.6792 | 0.6818 | 0.6667 | 0.6668 |
| TimeGAN | 0.6949 | 0.6762 | 0.6668 | 0.6668 |
| TimeVAE | 0.6788 | 0.9169 | 0.6667 | 0.6668 |
| DiffusionTS | 0.9683 | 0.6811 | 0.9349 | 0.9976 |

Table 12: Privacy Score (PS) ↓

| | MIMIC | | eICU | |
| --- | --- | --- | --- | --- |
| **Method** | Mortality | ICU Stay | ICU Stay | Mortality |
| TarDiff | 0.5819 | 0.5669 | 0.5795 | 0.5770 |
| BioDiffusion | 0.6226 | 0.6656 | 0.5954 | 0.5744 |
| TimeVQ-VAE | 1.2898 | 1.2115 | 0.5006 | 0.5671 |
| TimeGAN | 0.7035 | 0.8837 | 0.5198 | 0.6420 |
| TimeVAE | 0.9856 | 0.9616 | 0.5071 | 0.5028 |
| DiffusionTS | 0.8671 | 0.9097 | 0.6558 | 0.6821 |

To demonstrate TarDiff's fidelity and privacy preservation, we provide Discriminative Score (DS, measures how easily a classifier can distinguish synthetic samples from real ones.), Predictive Score (PS, evaluates how accurately models trained on synthetic data perform when predicting outcomes on real test data.), and Membership Inference Risk (MIR, assesses the privacy risk by quantifying vulnerability to membership inference attacks.)

# F    Visualization of Synthetic Data

We visualize the positive and negative samples generated by different methods. The results indicate that while TimeGAN performs relatively well for positive samples, its negative samples exhibit significant fluctuations and lack trend variations. Additionally, TimeVAE and TimeVQVAE produce overly smoothed sequences. In contrast, our approach more closely aligns with the distribution of real data.
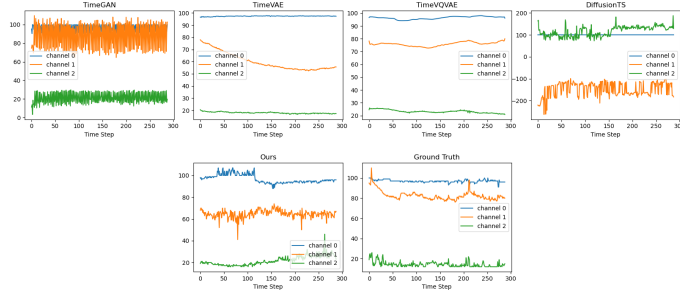
## F.1    More TSRTS Results

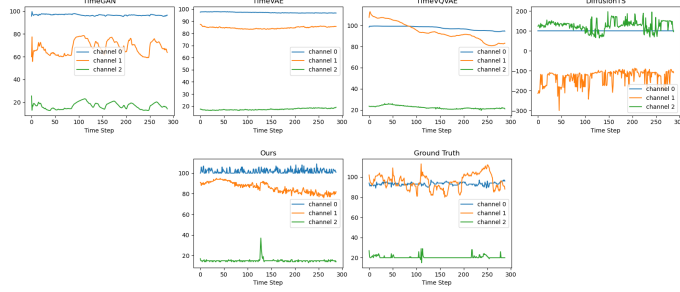Figure 5: Visualization of negative samples generated by different methods for ICU-Stay on eICU



Figure 6: Visualization of positive samples generated by different methods for ICU-Stay on eICU
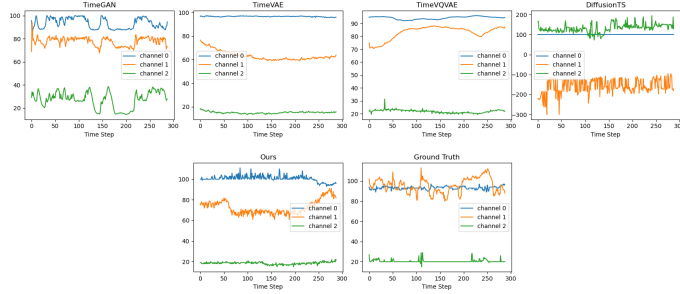


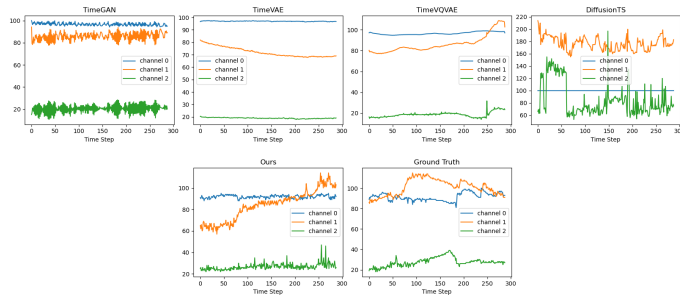Figure 7: Visualization of negative samples generated by different methods for Mortality on eICU



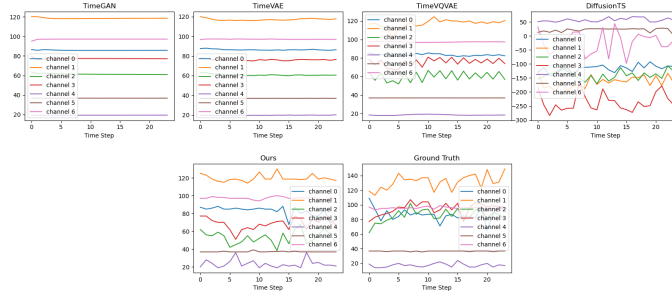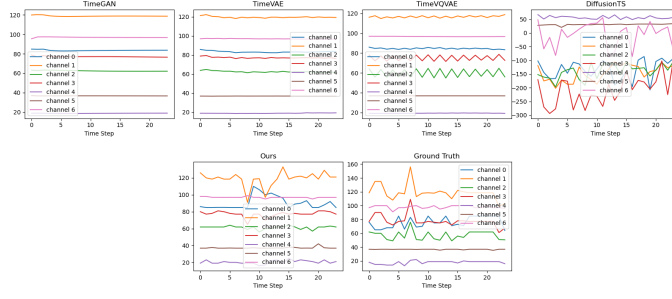Figure 8: Visualization of positive samples generated by different methods for Mortality on eICU

Figure 9: Visualization of negative samples generated by different methods for ICU Stay on MIMIC-III



Figure 10: Visualization of positive samples generated by different methods for ICU Stay on MIMIC-III
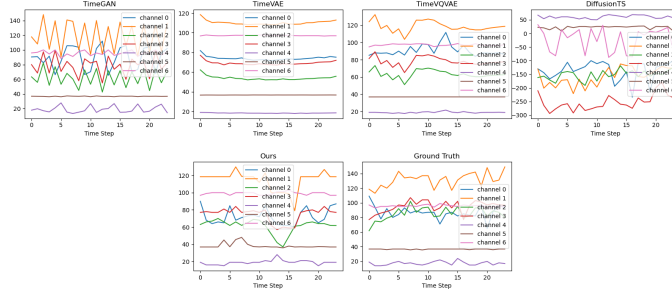


Figure 11: Visualization of negative samples generated by different methods for Mortality on MIMIC-III
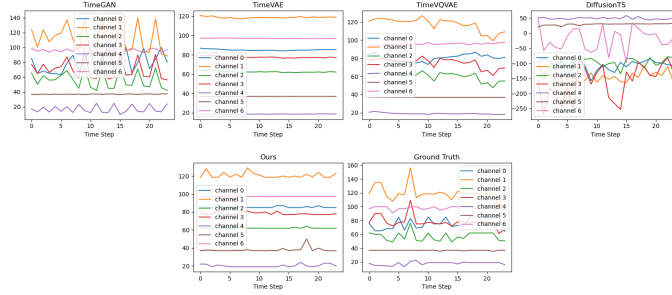


Figure 12: Visualization of positive samples generated by different methods for Mortality on MIMIC-III