

Generating Accurate Synthetic Survival Data by Conditioning on Outcomes

Mohammad Ashhad

MOHAMMAD.ASHHAD@KAUST.EDU.SA

BESE

King Abdullah University of Science and Technology (KAUST)

Thuwal, KSA

Ricardo Henao

RICARDO.HENAO@DUKE.EDU

Department of Bioinformatics & Biotatistics

Duke University

Durham, USA

Abstract

Synthetically generated data can improve privacy, fairness, and data accessibility; however, it can be challenging in specialized scenarios such as survival analysis. One key challenge in this setting is censoring, *i.e.*, the timing of an event is unknown in some cases. Existing methods struggle to accurately reproduce the distributions of both observed and censored event times when generating synthetic data. We propose a conceptually simple approach that generates covariates conditioned on event times and censoring indicators by leveraging existing tabular data generation models without making assumptions about the mechanism underlying censoring. Experiments on real-world datasets demonstrate that our method consistently outperforms baselines and improves downstream survival model performance.

1. Introduction

Synthetic data generation is the process of creating artificial data that mimic the statistical properties and patterns of real-world data. This technique has gained significant attention in various machine learning settings, including data privacy and data augmentation [Jordon et al. \(2022\)](#). The primary motivation behind the generation of synthetic data is to address challenges associated with limited availability, privacy concerns, or population imbalance that is often prevalent in real-world data ([Zhang et al., 2017](#); [Wang et al., 2021](#)). For example, researchers and organizations could train and evaluate models using synthetic data without compromising sensitive or proprietary information. Moreover, synthetic data can augment existing datasets, enabling more robust model performance. Alternatively, it can protect data privacy by providing a means to share and exchange data without revealing sensitive information, facilitating collaboration across different domains ([de Benedetti et al., 2020](#)).

Survival analysis, also known as time-to-event analysis, is a family of statistical methods that are used to analyze and model the time until the occurrence of a specific event (or outcome) of interest. These methods are widely applied in various fields, including biomedical research, operations research, engineering, economics, and social sciences ([Kaso et al., 2022](#); [Lillelund et al., 2023](#); [Danacica and Babucea, 2010](#); [Gross et al., 2014](#)). For example, assessing the effectiveness of medical treatments ([Singh and Mukhopadhyay, 2011](#)), predicting

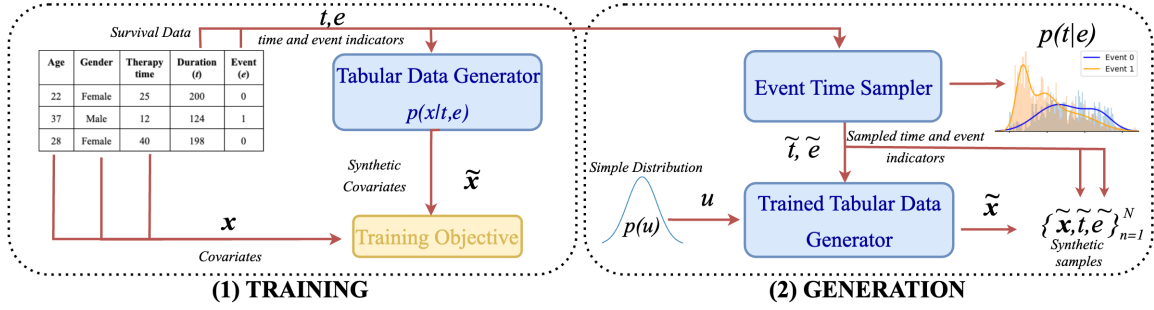


Figure 1: Block diagram of the proposed method. First, a conditional tabular data generator (we consider variational autoencoders, adversarial generators, diffusion-based models and large language models) is trained to learn to sample covariates from $p(\mathbf{x}|\mathbf{t}, e)$. After training, times \tilde{t} and event indicators \tilde{e} are sampled from a one-dimensional generator for $p(t, e = 1)$ and $p(t, e = 0)$ (we use a Dirichlet process mixture model, DPMM). These are then fed into the trained generator along with $\mathbf{u} \sim p(\mathbf{u})$, where $p(\mathbf{u})$ is a simple distribution. The generator then repeatedly generates synthetic covariates to completing the synthetic dataset $\mathcal{D} = \{(\tilde{\mathbf{x}}, \tilde{t}, \tilde{e})\}_{n=1}^N$.

equipment failure rates (de Cos Juez et al., 2010), or analyzing customer churn in the business domain (Danacica and Babucea, 2010). The primary goal of survival analysis is to estimate the probability (distribution) of an event occurring over time, given a set of covariates or risk factors. One of the distinctive challenges in survival analysis involves dealing with censored data, which occurs when the event of interest is not observed for some individuals within the study period. This can happen for various reasons, such as loss of follow-up, measurement failure, study termination, or the occurrence of competing risks (Salerno and Li, 2023). The handling of censored data requires tailored statistical methods to avoid biased survival estimates. Another challenge is that, oftentimes, sample sizes in survival data are relatively small, or the proportion of observed events relative to those with censoring is small, thus causing overfitting issues which negatively impact generalization ability.

In most domains, such as clinical trials or engineering studies, collecting large amounts of survival data can be challenging, time-consuming, and costly. Synthetic data generation allows researchers to create large datasets with desired characteristics, enabling more robust model prototyping, development, and evaluation. Synthetic survival data, which are predominantly tabular (or structured), can be generated using generative models that are specifically developed for tabular data, *e.g.*, autoencoders (Xu et al., 2019), adversarial generators (Yoon et al., 2020), diffusion generators (Kotelnikov et al., 2023), and even large language models (LLMs) (Borisov et al., 2022). However, in addition to the well-known challenges associated with generating tabular data, such as the appropriate handling of categorical and continuous data, mixed data types, and their joint distributions (Xu et al., 2019), survival data generation, especially in the medical domain, faces some unique challenges. These are due to mainly unavoidable differences in the distributions for observed and censored events, and their (unknown) underlying generation mechanism given the covariates. In practice, this challenge causes mismatches between these distributions when comparing real-world and synthetic data generated from it (Norcliffe et al., 2023). Consequently, such mismatches are likely to cause survival models trained on such synthetic data to underperform relative to

the real-world data in terms of discrimination and calibration. So motivated, our work offers the following contributions:

- We propose a simple method (see Figure 1) to generate survival data by conditioning the generation of covariates on event times and censoring indicators after sampling these from a model for the distribution of event times.
- We show that our *generator-agnostic* methodology can be easily extended to use LLM-based data generators to obtain high-quality synthetic survival data, an application that to the best of our knowledge has not been explored.
- Experiments on five real-world survival analysis datasets demonstrate the capabilities of the proposed method in terms of the quality of the generated observed, censored, and covariate distributions, as well as the discrimination and calibration performance of survival analysis models trained on synthetic data and evaluated on real-world data.

Generalizable Insights about Machine Learning in the Context of Healthcare

- **Synthetic data is crucial for advancing healthcare research while preserving privacy.** Healthcare datasets contain sensitive patient information, making them difficult to share across institutions. Synthetic data generation creates artificial data that mimics statistical properties of real clinical data without exposing patient identities, enabling broader collaboration while addressing limited data availability in specialized clinical studies.
- **Survival analysis poses unique challenges that require specialized modeling approaches.** The distinctive feature of survival data is censoring, where the timing of an event is unknown for some patients. Conventional data generation approaches struggle with correctly reproducing distributions of both observed and censored event times, which are critical for accurate prognostic assessment in healthcare.
- **Conditioning on event times and censoring indicators rather than covariates yields more accurate synthetic survival data.** Our approach of generating covariates conditioned on event times and censoring indicators preserves the critical time-to-event distributions by construction. This improves performance across multiple real-world clinical datasets while maintaining simplicity, making it practical for healthcare researchers.

2. Related Work

Generative models have emerged as powerful tools for synthesizing realistic data in various domains, including images, text, and tabular data. These models aim to learn the underlying probability distributions of the training data and generate new samples that exhibit similar characteristics. Three prominent classes of generative models have gained significant traction: *generative adversarial networks* (GANs), *variational autoencoders* (VAEs), and *diffusion-based models*. GANs employ an adversarial training paradigm, where a generator network learns to produce synthetic data samples, while a discriminator network aims to distinguish between real and generated samples (Goodfellow et al., 2020). This adversarial process drives the generator to produce increasingly realistic samples. VAEs leverage variational inference techniques to learn a latent representation of the data, allowing the generation of new samples

by sampling from the learned latent space (Kingma and Welling, 2013). Diffusion-based models, such as *denoising diffusion probabilistic model* (DDPM) (Ho et al., 2020), gradually add noise to the data and then learn to reverse this process, generating new samples by denoising random points. These generative models have shown remarkable success in various applications, including image synthesis (Kang et al., 2023), text generation (Su et al., 2022), and video generation (Jiang et al., 2023). In survival analysis, generative models have been applied to estimate event-time distributions and hazard functions (Chapfuwa et al., 2018; Zhou et al., 2022).

Tabular data stands out as a prevalent data format in machine learning (ML), with more than 65% of datasets found on the Google Dataset Search platform (datasetsearch.research.google.com) comprising tabular files, typically in comma separated or spreadsheet formats (Benjelloun et al., 2020). Although conventional generative methods are not optimally tailored for tabular data due to the mixture of continuous and categorical variables (Xu et al., 2019), modified versions have been developed for this domain. These include the *conditional tabular generative adversarial network* (CTGAN) (Xu et al., 2019), which leverages the GAN framework to generate synthetic data preserving multivariate distributions and relationships, the *tabular variational autoencoder* (TVAE) (Xu et al., 2019), and the *anonymization through data synthesis using generative adversarial network* (ADS-GAN) (Yoon et al., 2020). The *tabular denoising diffusion probabilistic model* (TabDDPM) is a recent approach that leverages denoising diffusion probabilistic models to generate high-fidelity synthetic tabular data (Kotelnikov et al., 2023). Large language models (LLMs) have also shown potential for tabular data generation, using fine-tuning on tabular data represented in token form (Borisov et al., 2022).

In the generation of synthetic survival data, early statistical models (Bender et al., 2005; Austin, 2012) transformed uniform samples into survival times but did not generate covariates. More recent techniques have incorporated deep learning into the generative process. Ranganath et al. (2016) proposed using deep exponential families to generate survival data, but this approach has limited flexibility on the learned distributions. Miscouridou et al. (2018) and Zhou et al. (2022) relaxed this assumption but still focused on generating survival times and censoring statuses conditioned on the covariates, rather than generating the covariates themselves. Recently, SurvivalGAN (Norcliffe et al., 2023) was developed, generating synthetic data in three steps: *i*) a conditional GAN (ADS-GAN) generates covariates (\mathbf{x}) and samples the event indicator (e) from the empirical distribution; *ii*) a survival function model (DeepHit (Lee et al., 2018)) predicts survival functions for the generated covariates; and *iii*) these outputs are used by a regression model (XGboost (Chen and Guestrin, 2016)) to predict the event-time (t), generating the complete triplet (\mathbf{x}, t, e) . Although effective, this method is complex with multiple models, each having their own limitations.

Our work explores as an alternative a much simpler method that leverages existing generators for one-dimensional (event time) distributions and tabular (covariates) data and repurposes them for survival data without the need for dedicated networks for the prediction of the survival function or event/censoring distributions.

3. Methods

Problem Definition Instances (or subjects) of survival data can be represented in general as a triplet $z = (\mathbf{x}, t, e)$. Here, $\mathbf{x} \in \mathcal{X}$ denotes m -dimensional tabular covariates that describe an instance’s state at an initial (or index) time, encompassing both continuous and categorical covariates. Then, $t_i \in \mathcal{T}$ represents the time of a specific event relative to the initial time, thus $t \geq 0$ and $\mathcal{T} \equiv \mathbb{R}_+$. Lastly, $e_i \in \mathcal{E}$ stands for the event indicator, commonly $\mathcal{E} = \{0, 1\}$, where $e = 1$ indicates the event of interest occurs at time t , while $e = 0$ indicates that the event of interest has not occurred up to time t . In this work, we only consider *right censoring* as it is the predominant form in real-world datasets, however, the proposed method can be extended to left or interval censoring Klein and Moeschberger (2006).

Background Survival analysis is a statistical framework used to analyze and model the time until the occurrence of the event of interest, also known as the survival time or time-to-event. Survival analysis involves modeling the conditional probability density function $p(t|\mathbf{x})$, to estimate the likelihood of the event of interest occurs at time t given the covariates \mathbf{x} . From this, the survival function is derived, representing the probability that the event has not taken place by time t , *i.e.*,

$$S(t | \mathbf{x}) = \int_t^\infty p(t' | \mathbf{x}) dt', \quad (1)$$

where $S(t | \mathbf{x})$ is an estimate of the proportion of instances (subjects) with covariates \mathbf{x} who have survived up to time t . When the initial time is zero and given that events cannot occur at $t \leq 0$, thus $S(0|\mathbf{x}) = 1$. Additionally, since $p(t|\mathbf{x})$ is a valid probability distribution (nonnegative), then $S(t|\mathbf{x})$ is a monotonically decreasing function. Time-to-event approximation involves estimating the expected lifetime for any given covariate value, denoted as $\mu(\mathbf{x})$. Specifically, this is obtained as $\mu(\mathbf{x}) = \int_0^\infty t' p(t'|\mathbf{x}) dt'$, which, through integration by parts, simplifies to the area under the survival curve: $\mu(\mathbf{x}) = \int_0^\infty S(t|\mathbf{x}) dt$.

Survival models typically fall into one of two categories: *i*) parametric such as the accelerated failure time (Weibull, 1951), and log-logistic (Prentice, 1976) models; or *ii*) non-parametric such as the Kaplan-Meier estimator (Kaplan and Meier, 1958) and Cox proportional hazards model (Cox, 1972). Moreover, deep-learning versions of these have been proposed, *e.g.*, DeepSurv (Katzman et al., 2018), DeepHit (Lee et al., 2018), DATE (Chapfuwa et al., 2018), *etc.*

Conditioning on Event Time and Type Synthetic survival data generation involves the generation of samples from the complete joint distribution $p(\mathbf{x}, t, e)$. In practice, one can either sample from it directly (and unconditionally) using generative models for tabular data, or via conditioning using for instance $p(t|\mathbf{x}, e)p(\mathbf{x})p(e)$ or $p(\mathbf{x}|t, e)p(t|e)p(e)$. The former is the approach used in Norcliffe et al. (2023), in which the samples $\tilde{\mathbf{x}}$ and \tilde{e} , from the marginals $p(\mathbf{x})$ and $p(e)$, are obtained using a conditional GAN (ADS-GAN) generator and the empirical distribution for the event indicators, respectively, and subsequently the samples \tilde{t} from the conditional $p(t|\mathbf{x}, e)$ are generated (deterministically) using a regression model. One important drawback of this approach is that the quality of the samples for event times \tilde{t} from $p(t|\mathbf{x}, e)$ is both dependent on the quality of the approximation $\tilde{t} \sim p_\phi(t|\mathbf{x}, e)$ (with parameters ϕ) and that of $p(\mathbf{x})$ via $\tilde{\mathbf{x}} \sim p_\psi(\mathbf{x}|\mathbf{u})$ parameterized by ψ , and \mathbf{u} being sampled from a simple distribution, *e.g.*, uniform or Gaussian. As a result, approximation error in

covariates \mathbf{x} is compounded with that of t , resulting in event and censoring distributions that do not necessarily match the real data. Consequently, Norcliffe et al. (2023) also proposed metrics to quantify the quality of these distributions (see Section 4).

In an effort to alleviate these key issues, we reverse the conditioning and instead sample *both* event times and type from their joint distribution via $p(t|e)$ and $p(e)$, using a simple one-dimensional generator. In our experiments, we consider a Dirichlet process mixture model (DPMM) (Blei and Jordan, 2006) to separately fit models for $p(t|e = 0)$ and $p(t|e = 1)$. Note that this is possible by assuming without loss of generality that the observed and censoring times are conditionally independent given the covariates, which also aligns with the common assumption of censoring at random in survival analysis, which posits that the censoring mechanism is independent of the unobserved survival times, conditional on the covariates. Then, we sample the covariates from $p(\mathbf{x}|t, e)$ using a conditional generator as follows:

$$\tilde{e} \sim p(e), \tilde{t} \sim p(t|\tilde{e}), \mathbf{u} \sim p(\mathbf{u}), \tilde{\mathbf{x}} \sim p_{\theta}(\mathbf{x}|\tilde{t}, \tilde{e}, \mathbf{u}), \quad (2)$$

where $p_{\theta}(\mathbf{x}|\tilde{t}, \tilde{e}, \mathbf{u})$ is a conditional generator parameterized by θ , while $p(\mathbf{u})$ is a simple distribution. Repeated sampling from the mechanism in (2) allows one to obtain a synthetic dataset $\mathcal{D} = \{(\mathbf{x}_n, t_n, e_n)\}_{n=1}^N$ whose empirical conditionals for event and censoring times readily match the ground-truth distributions, $p(t|e = 1)$ and $p(t|e = 0)$, respectively, and synthetic covariates that acknowledge their association with the event of interest while accounting for censoring.

Importantly, using (2): *i*) eliminates the need for a supervised model to generate event times (XGboost in Norcliffe et al. (2023)); *ii*) eliminates the need for a separate model to generate survival distributions (DeepHit in Norcliffe et al. (2023)), and *iii*) guarantees the quality of the observed and censored event distributions. Moreover, from a practical perspective, (2) offers flexibility, since $p_{\theta}(\mathbf{x}|\tilde{t}, \tilde{e}, \mathbf{u})$ can be modeled, in principle, with any conditional generator. In the experiments, we consider TVAE, CTGAN, ADS-GAN, TabDDPM and LLMs.

Note that in (2) we are not required to sample from $p(t|e = 0)$ and $p(t|e = 1)$ using a DPMM. Alternatively, one may use univariate (kernel) density estimators and then draw \tilde{t} and \tilde{e} accordingly, especially, if the dataset is small and the number of unique values of t in \mathcal{D} is small. Moreover, in situations where privacy is a less sensitive concern, one could simply draw $p(t|e = 0)$ and $p(t|e = 1)$ from their empirical distributions. We consider this setting in our experiments as a means to quantifying the impact on performance of using the DPMM to sample the event times.

3.1. Adapting Conditional Generators to Survival Data

Existing tabular generators (see Section 2) use distinct strategies to implement conditioning. Below, we briefly describe how they are adapted for survival data generation.

CTGAN This model being a conditional adversarial generator, synthesizes data using $G(\mathbf{u}, \mathbf{c})$, where $G(\cdot)$ is the generator specified as a neural network, \mathbf{u} is a vector sampled from a simple distribution, *e.g.*, a standard Gaussian distribution, *e.g.*, $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and \mathbf{c} is a one-hot vector representing a discrete conditioning covariate. See Xu et al. (2019) for additional details. In order to use $G(\mathbf{u}, \mathbf{c})$ as a sampling mechanism for $p_{\theta}(\mathbf{x}|\tilde{t}, \tilde{e}, \mathbf{u})$ in (2) we

simply set $\mathbf{c} = E_t(\tilde{t}) \oplus \tilde{e}$, where $E_t(\cdot)$ is an m -dimensional sinusoidal time embedding (Wang and Chen, 2020) and \oplus is the concatenation operator. In our experiments, we set $m = 4$.

TVAE The autoencoding formulation in Xu et al. (2019) does not specify explicitly how to perform conditional generation for the tabular VAE. However, the simplest strategy involves setting the encoder and decoder pair as $\mathbf{u} \sim \mathcal{N}(\mu(\mathbf{x}), \sigma^2(\mathbf{x}))$ and $\tilde{\mathbf{x}} \sim p_\theta(\mathbf{x}|\mathbf{c}, \mathbf{u})$, respectively, where here \mathbf{u} is the latent representation for covariates \mathbf{x} , $\mu(\mathbf{x})$ and $\sigma^2(\mathbf{x})$ are two neural networks for the mean and variance functions of the latent representation \mathbf{u} , $p_\theta(\mathbf{x}|\mathbf{c}, \mathbf{u})$ is a probabilistic decoder specified using neural networks (see Xu et al. (2019) for details), \mathbf{c} is a one-hot vector as above for CTGAN, and the input to the decoder conveniently implemented by concatenating \mathbf{z} and \mathbf{c} . Similarly to CTGAN, we make $\mathbf{c} = E_t(\tilde{t}) \oplus \tilde{e}$ in our implementation to sample from $p_\theta(\mathbf{x}|\tilde{t}, \tilde{e}, \mathbf{u})$ in (2) via $p_\theta(\mathbf{x}|\mathbf{c} = E_t(\tilde{t}) \oplus \tilde{e}, \mathbf{u})$.

ADS-GAN This alternative adversarial model specification encourages de-identifiability by letting the generator be $\tilde{\mathbf{x}} = G(\mathbf{u}, \mathbf{x}, \mathbf{c})$, *i.e.*, covariates \mathbf{x} are also used as input to the generation function $G(\cdot)$, to encourage the model to generate samples $\tilde{\mathbf{x}}$ that are distinct from \mathbf{x} to preserve privacy. See Yoon et al. (2020) for additional details. Consistent with CTGAN and TVAE above, we simply set $\mathbf{c} = E_t(\tilde{t}) \oplus \tilde{e}$.

TabDDPM This model designed specifically for tabular data uses a combination of Gaussian and multinomial diffusion processes to handle numerical and categorical features, respectively. Notably, each covariate uses a separate forward diffusion processes. The reverse diffusion function in Kotelnikov et al. (2023) is set as $\mathbf{x}_{is} = g_i(\mathbf{x}_i, \mathbf{x}_{i0}, s)$, where $g_i(\cdot)$ is modeled using neural networks with identity and softmax activations for continuous and discrete covariates, respectively, $\mathbf{x}_{is} = h_x(\mathbf{x}_i) + h_s(E_t(s)) + E_c(\mathbf{c})$ is the representation of the i -th covariate in \mathbf{x} at diffusion step s , $h_x(\mathbf{x}_i)$ is a fully connected layer with linear activation, $h_s(\cdot)$ is composed of two fully connected layers with sigmoid linear activations, $E_c(\cdot)$ is a standard (trainable) categorical embedding, and $s = 0, \dots, S$, is such that $\mathbf{x}_{iS} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ or $\mathbf{x}_{iS} \sim \text{Cat}(\mathbf{1}/K_i)$, for K_i categories (distinct values), for continuous or discrete covariates, respectively. Note that, effectively, $g_i(\cdot)$ models the residuals of \mathbf{x}_{is} at diffusion step s rather than \mathbf{x}_{is} itself (Nichol and Dhariwal, 2021). For additional details of the formulation and components of the model architecture, see Kotelnikov et al. (2023). For our implementation, we set $\mathbf{c} = E_t(\tilde{t}) + E_s(\tilde{e})$ and set $m = 128$ as the embedding dimension.

4. Experiments

Baselines and setup We compare our methodology against the following baselines: generative adversarial networks for anonymization (ADS-GAN) (Yoon et al., 2020); conditional generative adversarial networks for tabular data (CTGAN) (Xu et al., 2019); variational autoencoder for tabular data (TVAE) (Xu et al., 2019); tabular denoising diffusion probabilistic models (TabDDPM) (Kotelnikov et al., 2023); and SurvivalGAN (Norcliffe et al., 2023). Note that only the latter is specific to survival data, whereas all the others generate tabular data *unconditionally*, *i.e.*, from the joint $p(\mathbf{x}, t, e)$. For CTGAN, TVAE, ADS-GAN, and TabDDPM models, we report metrics both directly using the models *without* conditioning (Unconditional) for survival data generation, as well as our methodology (see Section 3.1), *i.e.*, using them as conditional generators given event times and censoring indicators sampled using DPMMs.

To evaluate downstream performance, survival models are trained on synthetic data and tested on real data using the Train on Synthetic Test on Real (TSTR) paradigm (Esteban et al., 2017). Specifically, the original dataset is divided into three folds, and the synthetic data generator is trained on two folds while the third is reserved for testing. Synthetic data equivalent (in size) to the training data is then generated, and downstream models are trained on this synthetic dataset and evaluated on the held-out *real test set*. This process is repeated for all three fold combinations. We consider various survival models: linear (CoxPH) (Cox, 1972), gradient boosting (SurvivalXGBoost) (Barnwal et al., 2022), and neural networks (DeepHit) (Lee et al., 2018), and report metrics for the best-performing model. For each dataset, benchmark, and experimental setting, we report mean and standard deviation of performance metrics using 5 random seeds. To streamline the benchmarking, we utilized the Synthcity library (Qian et al., 2024), which provides implementations of a variety of synthetic tabular data generation models and benchmarking utilities. Detailed experimental settings and hyperparameters are in Appendix B.3. The source code for reproducing experiments is available at https://github.com/ashhadm/synthetic_survival_data.

Datasets We benchmark our methodology on a variety of real-world medical datasets. Specifically: *i)* *Study to understand prognoses preferences outcomes and risks of treatment* (SUPPORT) (Knaus et al., 1995); *ii)* *Molecular taxonomy of breast cancer international consortium* (METABRIC) (Curtis et al., 2012); *iii)* *ACTG 320 clinical trial dataset* (AIDS) (Hammer et al., 1997); *iv)* *Rotterdam & German breast cancer study group* (GBSG) (Schumacher et al., 1994); and *v)* *Assay of serum free light chain* (FLCHAIN) (Dispenzieri et al., 2012). See Appendix B.2 for additional details.

Metrics To evaluate the quality of the generated synthetic survival data, various metrics are employed, which can be categorized into three groups: *covariates quality*, *event-time distribution quality*, and *downstream performance*. For assessing the quality of the generated covariates $\tilde{\mathbf{x}}$, the Jensen-Shannon (JS) distance and Wasserstein distance (WS) are used to measure the divergence between the generated and original covariate distributions. For the quality of the event-time distributions we quantify the alignment between ground-truth and generated temporal marginals, namely, $p(t, e)$ is evaluated using the Kaplan-Meier (KM) divergence, optimism, and short-sightedness metrics as previously described in Norcliffe et al. (2023). The KM divergence compares the mean absolute difference between the synthetic and real survival function estimates, while optimism and short-sightedness are a proxy for their bias and variance, respectively. These three metrics capture the accuracy of the generated censoring and event distributions. Finally, to assess downstream performance, survival models are trained on the synthetic data and evaluated on real dataset. Specifically, we consider the concordance index (C-index) (Harrell et al., 1982) and the Brier score (Brier, 1950). The former measures the discriminative ability of the survival model, while the latter quantifies the calibration of the probabilistic predictions.

4.1. Synthetic Survival Data Generation Benchmark

Covariate quality metrics: Results in Table 1 compare the similarity between the distribution of synthetic samples and the original data. First, we assess the overall (covariance) structure of the synthetic covariates relative to the original data via the JS and WS distances.

Table 1: Quality, downstream and event-time metrics. Performance reported for best model between ADS-GAN, TVAE, CTGAN and TabDDPM. Original is for the survival model trained on the real (training) data. Subscripts are standard deviations for 5 repetitions.

Metric	Method	AIDS	METABRIC	SUPPORT	GBSG	FLCHAIN
JS distance (\downarrow)	SurvivalGAN	0.013 _{0.005}	0.009 _{0.000}	0.008 _{0.004}	0.008 _{0.001}	0.009 _{0.005}
	Unconditional	0.006 _{0.001}	0.007 _{0.000}	0.005 _{0.004}	0.005 _{0.002}	0.002 _{0.000}
	Ours	0.006 _{0.005}	0.006 _{0.002}	0.004 _{0.005}	0.004 _{0.002}	0.001 _{0.003}
WS distance (\downarrow)	SurvivalGAN	0.112 _{0.015}	0.039 _{0.002}	0.043 _{0.004}	0.019 _{0.005}	0.052 _{0.000}
	Unconditional	0.065 _{0.005}	0.031 _{0.005}	0.036 _{0.005}	0.013 _{0.004}	0.016 _{0.005}
	Ours	0.063 _{0.004}	0.030 _{0.002}	0.032 _{0.002}	0.011 _{0.004}	0.016 _{0.002}
C-Index (\uparrow)	Original	0.760 _{0.005}	0.636 _{0.004}	0.616 _{0.002}	0.695 _{0.006}	0.870 _{0.004}
	SurvivalGAN	0.735 _{0.005}	0.625 _{0.000}	0.602 _{0.004}	0.668 _{0.005}	0.870 _{0.002}
	Unconditional	0.779 _{0.002}	0.649 _{0.004}	0.625 _{0.002}	0.679 _{0.002}	0.879 _{0.004}
	Ours	0.785 _{0.025}	0.652 _{0.005}	0.626 _{0.004}	0.682 _{0.002}	0.880 _{0.005}
Brier Score (\downarrow)	Original	0.062 _{0.005}	0.200 _{0.004}	0.195 _{0.002}	0.205 _{0.005}	0.095 _{0.004}
	SurvivalGAN	0.068 _{0.005}	0.205 _{0.004}	0.202 _{0.002}	0.212 _{0.005}	0.096 _{0.004}
	Unconditional	0.060 _{0.005}	0.200 _{0.004}	0.199 _{0.002}	0.207 _{0.005}	0.086 _{0.005}
	Ours	0.060 _{0.004}	0.197 _{0.005}	0.198 _{0.002}	0.210 _{0.004}	0.085 _{0.004}
Optimism ($\rightarrow 0$)	SurvivalGAN	0.021 _{0.005}	0.011 _{0.002}	0.016 _{0.004}	0.006 _{0.003}	0.134 _{0.005}
	Unconditional	0.002 _{0.002}	0.001 _{0.005}	0.001 _{0.003}	0.004 _{0.005}	0.005 _{0.004}
	Ours	0.001 _{0.001}	0.001 _{0.001}	-0.001 _{0.001}	-0.004 _{0.000}	-0.003 _{0.001}
Shortsightedness ($\rightarrow 0$)	SurvivalGAN	0.007 _{0.003}	0.124 _{0.004}	0.020 _{0.002}	0.019 _{0.005}	0.005 _{0.003}
	Unconditional	0.002 _{0.005}	0.000 _{0.002}	0.002 _{0.003}	0.014 _{0.002}	0.003 _{0.005}
	Ours	0.001 _{0.001}	0.000 _{0.001}	0.000 _{0.002}	-0.012 _{0.001}	-0.002 _{0.002}
KM Divergence (\downarrow)	SurvivalGAN	0.021 _{0.004}	0.082 _{0.005}	0.064 _{0.002}	0.049 _{0.003}	0.134 _{0.004}
	Unconditional	0.015 _{0.003}	0.019 _{0.003}	0.011 _{0.004}	0.026 _{0.005}	0.007 _{0.002}
	Ours	0.003 _{0.001}	0.013 _{0.002}	0.006 _{0.001}	0.007 _{0.001}	0.005 _{0.001}

Our models outperformed or matched baselines in all 5 datasets for JS distance, and WS distance. Full results are shown in Appendix C.

Downstream Performance We conduct a comparative analysis of survival models trained with synthetic data generated by our methodology against models trained with data from baseline methods. A favorable outcome is achieved when a model trained with synthetic data performs comparably to or occasionally even better than a model trained with real data, while also outperforming models trained with alternative synthetic data sources. For reference, we also report the C-Index and Brier Score for survival models trained on the original data. C-index and Brier score serve as the most widely used indicators of performance, as they encapsulate the entire conditional distribution of covariates, event/censoring times, and event indicators $p(t, e|x)$. Results in Table 1 demonstrate that in C-index, we outperform the baselines across all datasets while in Brier Score, we outperform the baselines in 4 of 5 datasets. Further, in most cases, we were also able to achieve better performance than survival models trained on the original data.

Event-time distribution metrics The proposed method demonstrates superior performance in preserving the underlying event-time distribution characteristics across all datasets as shown in Table 1. Our approach achieves notably lower KM divergence scores compared to both SurvivalGAN and Unconditional baselines. The method also exhibits minimal optimism and shortsightedness, with values consistently close to the ideal score of zero across all datasets.

Table 2: Quality and downstream performance metrics for synthetic survival data generation using LLMs. *Ours* refer to our best model between ADS-GAN, TVAE, CTGAN and TabDDPM from Table 1. Subscripts are standard deviations for 5 repetitions.

Metric	Method	AIDS	METABRIC	SUPPORT	GBSG	FLCHAIN
JS Distance (\downarrow)	SurvivalGAN	0.013 _{0.005}	0.009 _{0.000}	0.008 _{0.004}	0.008 _{0.001}	0.009 _{0.005}
	Ours	0.006 _{0.005}	0.006 _{0.002}	0.004 _{0.005}	0.004 _{0.002}	0.001 _{0.003}
	LLM	0.004 _{0.001}	0.006 _{0.001}	0.003 _{0.000}	0.007 _{0.001}	0.001 _{0.000}
	LLM (<i>Ours</i>)	0.003 _{0.001}	0.005 _{0.001}	0.002 _{0.000}	0.006 _{0.001}	0.001 _{0.000}
WS Distance (\downarrow)	SurvivalGAN	0.112 _{0.015}	0.039 _{0.002}	0.043 _{0.004}	0.019 _{0.005}	0.052 _{0.000}
	Ours	0.063 _{0.004}	0.030 _{0.002}	0.032 _{0.002}	0.011 _{0.004}	0.016 _{0.002}
	LLM	0.046 _{0.001}	0.000 _{0.000}	0.020 _{0.002}	0.011 _{0.001}	0.020 _{0.002}
	LLM (<i>Ours</i>)	0.040 _{0.003}	0.000 _{0.000}	0.025 _{0.001}	0.010 _{0.001}	0.015 _{0.001}
C-Index (\uparrow)	SurvivalGAN	0.735 _{0.005}	0.625 _{0.000}	0.602 _{0.004}	0.668 _{0.005}	0.870 _{0.002}
	Ours	0.785 _{0.025}	0.652 _{0.005}	0.626 _{0.004}	0.682 _{0.002}	0.880 _{0.005}
	LLM	0.725 _{0.010}	0.623 _{0.002}	0.627 _{0.000}	0.672 _{0.002}	0.878 _{0.001}
	LLM (<i>Ours</i>)	0.787 _{0.002}	0.655 _{0.002}	0.628 _{0.001}	0.684 _{0.001}	0.880 _{0.000}
Brier Score (\downarrow)	SurvivalGAN	0.068 _{0.005}	0.205 _{0.004}	0.202 _{0.002}	0.212 _{0.005}	0.096 _{0.004}
	Ours	0.060 _{0.004}	0.197 _{0.005}	0.198 _{0.002}	0.210 _{0.004}	0.085 _{0.004}
	LLM	0.063 _{0.001}	0.201 _{0.000}	0.200 _{0.001}	0.207 _{0.001}	0.090 _{0.001}
	LLM (<i>Ours</i>)	0.060 _{0.001}	0.196 _{0.001}	0.198 _{0.001}	0.207 _{0.002}	0.083 _{0.001}

Ablation In our ablation study, we investigated the impact of an alternate sampling strategy for time and event variables during the generation of synthetic data. Although our base method uses DPMM to sample these variables, we also explore directly sampling *both* event times and type from their joint distribution via $p(t|e)$ and $p(e)$, using their empirical distributions. This modification led to modest but consistent improvements across all evaluation metrics while being remarkably simple and eliminating the need for a separate model to generate event times (See Table 13 in the Appendix). Importantly, this performance gain did not come at the cost of privacy, as evidenced by the median and minimum Distance to Closest Record (DCR) values (see Table 14 and 15 in the Appendix). In fact, the direct sampling approach achieved a higher median DCR in 3 out of 5 datasets and a higher minimum DCR in 4 out of 5 datasets when compared to *unconditional* models, albeit by small margins. These results suggest that directly sampling t and e from the underlying distributions preserves the utility of the generated synthetic data while simplifying the generation process. Full results are shown in Appendix C. Based on the comprehensive evaluation metrics shown in the table, ADS-GAN consistently demonstrates superior performance in preserving the overall survival distribution, achieving the highest C-Index scores across most datasets while maintaining competitive Brier Scores. TVAE excels in terms of data fidelity metrics, showing particularly strong performance in JS and WS distances, making it particularly suitable when the priority is modeling the marginals of the original dataset. For overall utility and general-purpose synthetic survival data generation, ADS-GAN emerges as the most well-rounded model, offering the best balance between distribution preservation (high C-Index) and data fidelity (competitive JS/WS distances) while maintaining comparable performance in terms of calibration.

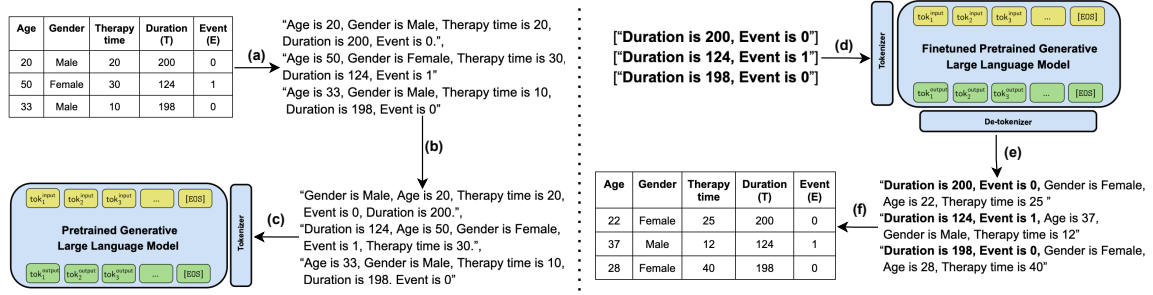


Figure 2: Training and sampling procedure for survival data generation using LLMs.

4.2. Fine Tuning an LLM for Survival Data Generation

Generation of realistic tabular data (GReaT) is a recently proposed approach to generating high-quality synthetic tabular data using LLMs [Borisov et al. \(2022\)](#). This is achieved by representing the tabular data as a sequence of text and training the language model to generate new sequences that correspond to valid and plausible tabular data instances. We adapt GReaT to generate synthetic survival data by conditioning the generation on time-to-event and event-type. The fine-tuning of a pre-trained auto-regressive LLM on the encoded tabular data for data generation as proposed in [Borisov et al. \(2022\)](#) involves the following steps.

Textual encoding and feature permutation: The tabular data with M column names $\{f_m\}_{m=1}^M$ and thus, M -dimensional rows $\{\mathbf{x}_n\}_{n=1}^N$ are converted into textual representation. Each row (sample) \mathbf{x}_n is encoded as a sentence with elements $\mathbf{t}_n = \{t_{nm}\}_{m=1}^M$, where $t_{nm} = [f_m, \text{"is"}, x_{nm}, \text{" "}]$ contains the column name f_m and its value x_{nm} . **Model training:** The LLM is trained using DistilGPT2 ([Li et al., 2021](#)) on the textually encoded dataset $\{\mathbf{t}_n\}_{n=1}^N$, with elements of \mathbf{t}_n permuted at random to remove pseudo-positional information as column order in a tabular dataset is in principle non-informative. **Sampling:** Feature permutations during training enable the model to start generation with any combination of features and values. To generate synthetic data conditionally, we prompt the trained model with conditioning sequences sampled from the DPMM model for $p(t, e)$, where survival times \tilde{t} are drawn from the fitted DPMM for each event type $e \in \{0, 1\}$, and let it generate the remaining tokens to complete the textual feature vector. Unconditional generation follows [Borisov et al. \(2022\)](#). The training and sampling procedure is shown in Figure 2. Table 2 compares the performance of GReaT with and without conditional generation (LLM (*Ours*) and LLM respectively in Table 2), against our best generator from Table 1. We observe that conditional generation consistently enhances LLM’s performance over the unconditional variant and baseline generators. Note however that GReaT is much more costly compared to other models as shown in Appendix B.1.

4.3. Sub-population Level Evaluation of Synthetic Data

In this experiment, we evaluate the performance of the proposed methodology at the sub-population level using the AIDS dataset, using race (White, Black and Hispanic) to define the sub-populations. Performance evaluation is carried out via race-stratified K -fold cross-validation. We consider survival models in three scenarios: *i*) trained on the real data; *ii*)

Table 3: Downstream performance metrics for survival models trained on Real Data, *Synthetic*, and *Synthetic (Balanced)*.

Method	Race	Synthetic		Synthetic (Balanced)	
		C-index	Brier Score	C-index	Brier Score
ADS-GAN (Ours)	All	0.718 _{0.004}	0.070 _{0.001}	0.741 _{0.002}	0.066 _{0.002}
	Race 1	0.718 _{0.002}	0.067 _{0.002}	0.723 _{0.001}	0.062 _{0.002}
	Race 2	0.718 _{0.002}	0.071 _{0.003}	0.723 _{0.001}	0.062 _{0.002}
	Race 3	0.760 _{0.005}	0.071 _{0.003}	0.761 _{0.010}	0.062 _{0.001}
SurvivalGAN	All	0.663 _{0.004}	0.100 _{0.020}	0.683 _{0.010}	0.076 _{0.010}
	Race 1	0.663 _{0.003}	0.092 _{0.005}	0.676 _{0.002}	0.072 _{0.010}
	Race 2	0.663 _{0.003}	0.095 _{0.010}	0.676 _{0.010}	0.073 _{0.020}
	Race 3	0.668 _{0.010}	0.095 _{0.010}	0.698 _{0.010}	0.073 _{0.010}
Method	Race	Real Data			
		C-index		Brier Score	
Original	All	0.735 _{0.010}		0.075 _{0.010}	
	Race 1	0.724 _{0.001}		0.069 _{0.002}	
	Race 2	0.724 _{0.001}		0.072 _{0.001}	
	Race 3	0.778 _{0.020}		0.072 _{0.001}	

trained on synthetic data with the same race proportion as the original data (*Synthetic*); and *iii*) trained on synthetic data with balanced race samples while preserving the distribution of observed and censored events for each race (*Synthetic (Balanced)*).

Based on the results shown in Table 3, for the survival models trained on the original AIDS dataset, the C-index differs across races, with the model performing better on Hispanic (0.778) when compared to White (0.724) and Black (0.724), with a $0.778/0.724 \approx 1.07$ ratio. When training using our synthetic data (ADS-GAN conditioned on time and event) with the same distribution as the original data, the C-index values reflect a similar performance ratio of 1.06 between races. For the balanced distribution scenario, all performance metrics improve at the expense of reducing the performance ratio between Hispanic and White/Black observed in the original data to 1.05. Furthermore, the proposed model consistently outperforms SurvivalGAN, which is less able to capture the race performance difference with ratios 1.01 and 1.03 for Synthetic and Synthetic (Balanced), respectively.

4.4. Robustness to Limited Training Data

To further verify the advantage of our reverse conditioning approach, we conducted an experiment examining how different methods perform with reduced training data. Using the AIDS dataset, we trained models on progressively smaller subsets (100%, 75%, and 50%) of the original training data and evaluated their performance across different performance metrics. As shown in Table 4, our reverse conditioning approach consistently outperforms both unconditional models and SurvivalGAN in all metrics, the performance gap widening as training data become more limited. With only 50% of the training data, our method still maintains strong performance, while other approaches show a more substantial degradation.

Table 4: Performance comparison with reduced training data on the AIDS dataset.

Method	JS Distance	WS Distance	C-Index	Brier Score	KM Divergence
100% Training Data					
Ours	0.006 _{0.005}	0.063 _{0.004}	0.785 _{0.025}	0.060 _{0.004}	0.003 _{0.001}
Unconditional	0.006 _{0.002}	0.065 _{0.005}	0.779 _{0.002}	0.060 _{0.005}	0.015 _{0.003}
SurvivalGAN	0.013 _{0.005}	0.112 _{0.015}	0.735 _{0.005}	0.068 _{0.005}	0.021 _{0.004}
75% Training Data					
Ours	0.007 _{0.004}	0.069 _{0.005}	0.780 _{0.023}	0.062 _{0.005}	0.004 _{0.001}
Unconditional	0.008 _{0.003}	0.072 _{0.006}	0.770 _{0.015}	0.063 _{0.005}	0.018 _{0.004}
SurvivalGAN	0.016 _{0.006}	0.120 _{0.016}	0.720 _{0.012}	0.072 _{0.006}	0.027 _{0.005}
50% Training Data					
Ours	0.009 _{0.005}	0.079 _{0.004}	0.762 _{0.025}	0.065 _{0.004}	0.007 _{0.002}
Unconditional	0.011 _{0.004}	0.085 _{0.005}	0.755 _{0.018}	0.068 _{0.005}	0.025 _{0.005}
SurvivalGAN	0.022 _{0.007}	0.135 _{0.020}	0.695 _{0.015}	0.078 _{0.008}	0.038 _{0.006}

Importantly, KM Divergence, which directly measures how well the models capture the survival distribution, remains consistently better with our approach, showing a 3.6x improvement over unconditional models and 5.4x improvement over SurvivalGAN when using only 50% of the training data. Similarly, our method’s C-Index decreases by only 2.9% when reducing training data by half, compared to a 5.4% decrease for SurvivalGAN.

These results demonstrate that directly conditioning the generation of covariates on event-time and censoring indicators provides greater robustness when training data is limited, a crucial advantage in healthcare settings where large annotated datasets are often difficult to obtain. By separating the generation of event times from covariates, our approach utilizes the information available in small datasets more efficiently, making it particularly suitable for practical clinical applications.

5. Discussion

This work proposed a simple yet effective methodology for generating high-quality synthetic survival data by conditioning the generation of covariates on event times and censoring indicators sampled from a one-dimensional distribution approximated with a DPMM. Through extensive experiments on multiple real-world datasets, we demonstrated that our approach outperforms several competitive baselines across various evaluation metrics that assess the quality of the generated covariate distributions, alignment with the ground-truth event-time distributions, and the downstream performance of survival models trained on the synthetic data. Moreover, we showcased the applicability of LLMs for survival data generation by fine-tuning them in a conditional manner on the textual representations of tabular data and how the proposed method preserves the sub-population-level performance characteristics of real-world data. Finally, we also demonstrated that our method performs better than baselines when there is less data available for training, as is common in survival studies.

Limitations Despite promising results, our work has limitations. First, the quality of generated data depends on the representativeness and diversity of the original dataset used for training the generative models. If the training data exhibit biases or lack variability,

these will propagate to the synthetic data. Second, while our approach ensures accurate reproduction of event-time and censoring distributions, it does not consider time-varying covariates, which may be relevant in certain applications. Finally, further research is needed to address bias and equity in survival data. Though we examine survival models’ behavior trained on synthetic data at a sub-population level, we acknowledge that bias and equity are multifaceted challenges extending beyond this study. These are exciting avenues for further research.

6. Broader Impact

The ability to generate realistic synthetic survival datasets can have far-reaching impacts across various domains, especially in privacy-sensitive applications like healthcare and clinical research. Synthetic data can enable model development, benchmarking, and collaboration while preserving patient confidentiality and complying with data protection regulations. Furthermore, our methodology can potentially address the common challenge of limited data availability in survival analysis by augmenting existing datasets or creating entirely new synthetic datasets tailored to specific requirements. While synthetic survival data is specific to the domain to which it is applied, limiting the potential for misuse, it is important to acknowledge the possibility of reinforcing biases present in the training data, as is the case with any generative model. Though we aim to understand the behavior of survival models trained on synthetic data across sub-populations, we recognize that addressing bias and ensuring equity are complex challenges that extend beyond the scope of this study. Thus, it is crucial to exercise caution and implement appropriate safeguards to mitigate potential biases and promote fairness in the development and deployment of such models. We have no conflicts of interest to declare. All datasets used are publicly available and de-identified. Our work complies with relevant data protection regulations. We encourage users of our method to carefully consider the ethical implications and potential biases when applying it to sensitive healthcare data.

References

- Peter C. Austin. Generating survival times to simulate cox proportional hazards models with time-varying covariates. *Statistics in medicine*, 31(29):3946–3958, 2012.
- Avinash Barnwal, Hyunsu Cho, and Toby Hocking. Survival regression with accelerated failure time model in xgboost. *Journal of Computational and Graphical Statistics*, 31(4):1292–1302, 2022.
- Ralf Bender, Thomas Augustin, and Maria Blettner. Generating survival times to simulate cox proportional hazards models. *Statistics in medicine*, 24(11):1713–1723, 2005.
- Omar Benjelloun, Shiyu Chen, and Natasha Noy. Google dataset search by the numbers. In *International Semantic Web Conference*, pages 667–682. Springer, 2020.
- David M. Blei and Michael I. Jordan. Variational inference for dirichlet process mixtures. 2006.

- Vadim Borisov, Kathrin Sekler, Tobias Leemann, Martin Pawelczyk, and Gjergji Kasneci. Language models are realistic tabular data generators. *arXiv preprint arXiv:2210.06280*, 2022.
- Glenn W. Brier. Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3, 1950.
- Paidamoyo Chapfuwa, Chenyang Tao, Chunyuan Li, Courtney Page, Benjamin Goldstein, Lawrence Carin Duke, and Ricardo Henao. Adversarial time-to-event modeling. In *International Conference on Machine Learning*, pages 735–744. PMLR, 2018.
- Nitesh V. Chawla, Kevin W. Bowyer, Lawrence O. Hall, and W. Philip Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16: 321–357, 2002.
- Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- David R. Cox. Regression models and life-tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 34(2):187–202, 1972.
- Christina Curtis, Sohrab P. Shah, Suet-Feung Chin, Gulisa Turashvili, Oscar M. Rueda, Mark J. Dunning, Doug Speed, Andy G. Lynch, Shamith Samarajiwa, and Yinyin Yuan. The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature*, 486(7403):346–352, 2012.
- Daniela-Emanuela Danacica and Ana-Gabriela Babucea. Using survival analysis in economics. *survival*, 11:15, 2010.
- Juan de Benedetti, Namir Oues, Zhenchen Wang, Puja Myles, and Allan Tucker. Practical lessons from generating synthetic healthcare data with bayesian networks. In *ECML PKDD 2020 Workshops: Workshops of the European Conference on Machine Learning and Knowledge Discovery in Databases (ECML PKDD 2020): SoGood 2020, PDFL 2020, MLCS 2020, NFMCP 2020, DINA 2020, EDML 2020, XKDD 2020 and INRA 2020, Ghent, Belgium, September 14–18, 2020, Proceedings*, pages 38–47. Springer, 2020.
- Francisco Javier de Cos Juez, P.J. García Nieto, J. Martínez Torres, and J. Taboada Castro. Analysis of lead times of metallic components in the aerospace industry through a supported vector machine model. *Mathematical and computer modelling*, 52(7-8):1177–1184, 2010.
- Angela Dispenzieri, Jerry A. Katzmman, Robert A. Kyle, Dirk R. Larson, Terry M. Therneau, Colin L. Colby, Raynell J. Clark, Graham P. Mead, Shaji Kumar, and L. Joseph Melton III. Use of nonclonal serum immunoglobulin free light chains to predict overall survival in the general population. In *Mayo Clinic Proceedings*, volume 87, pages 517–523. Elsevier, 2012.
- Cristóbal Esteban, Stephanie L. Hyland, and Gunnar Rätsch. Real-valued (medical) time series generation with recurrent conditional gans. *arXiv preprint arXiv:1706.02633*, 2017.

- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- Samuel R. Gross, Barbara O’Brien, Chen Hu, and Edward H. Kennedy. Rate of false conviction of criminal defendants who are sentenced to death. *Proceedings of the National Academy of Sciences*, 111(20):7230–7235, 2014.
- Scott M. Hammer, Kathleen E. Squires, Michael D. Hughes, Janet M. Grimes, Lisa M. Demeter, Judith S. Currier, Joseph J. Eron Jr, Judith E. Feinberg, Henry H. Balfour Jr, and Lawrence R. Deyton. A controlled trial of two nucleoside analogues plus indinavir in persons with human immunodeficiency virus infection and cd4 cell counts of 200 per cubic millimeter or less. *New England Journal of Medicine*, 337(11):725–733, 1997.
- Frank E. Harrell, Robert M. Califf, David B. Pryor, Kerry L. Lee, and Robert A. Rosati. Evaluating the yield of medical tests. *Jama*, 247(18):2543–2546, 1982.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Yuming Jiang, Shuai Yang, Tong Liang Koh, Wayne Wu, Chen Change Loy, and Ziwei Liu. Text2performer: Text-driven human video generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22747–22757, 2023.
- James Jordon, Lukasz Szpruch, Florimond Houssiau, Mirko Bottarelli, Giovanni Cherubin, Carsten Maple, Samuel N. Cohen, and Adrian Weller. Synthetic data—what, why and how? *arXiv preprint arXiv:2205.03257*, 2022.
- Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung Park. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10124–10134, 2023.
- Edward L. Kaplan and Paul Meier. Nonparametric estimation from incomplete observations. *Journal of the American statistical association*, 53(282):457–481, 1958.
- Abdene Weya Kaso, Gebi Agero, Zewdu Hurissa, Taha Kaso, Helen Ali Ewune, Habtamu Endashaw Hareru, and Alemayehu Hailu. Survival analysis of covid-19 patients in ethiopia: a hospital-based study. *PLoS One*, 17(5):e0268280, 2022.
- Jared L. Katzman, Uri Shaham, Alexander Cloninger, Jonathan Bates, Tingting Jiang, and Yuval Kluger. DeepSurv: personalized treatment recommender system using a cox proportional hazards deep neural network. *BMC medical research methodology*, 18:1–12, 2018.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- John P Klein and Melvin L Moeschberger. *Survival analysis: techniques for censored and truncated data*. Springer Science & Business Media, 2006.

- William A. Knaus, Frank E. Harrell, Joanne Lynn, Lee Goldman, Russell S. Phillips, Alfred F. Connors, Neal V. Dawson, William J. Fulkerson, Robert M. Califf, and Norman Desbiens. The support prognostic model: Objective estimates of survival for seriously ill hospitalized adults. *Annals of internal medicine*, 122(3):191–203, 1995.
- Akim Kotelnikov, Dmitry Baranchuk, Ivan Rubachev, and Artem Babenko. Tabddpm: Modelling tabular data with diffusion models. In *International Conference on Machine Learning*, pages 17564–17579. PMLR, 2023.
- Changhee Lee, William Zame, Jinsung Yoon, and Mihaela Van Der Schaar. Deephit: A deep learning approach to survival analysis with competing risks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Tianda Li, Yassir El Mesbahi, Ivan Kobzyev, Ahmad Rashid, Atif Mahmud, Nithin Anchuri, Habib Hajimolahoseini, Yang Liu, and Mehdi Rezagholizadeh. A short study on compressing decoder-based language models. *arXiv preprint arXiv:2110.08460*, 2021.
- Christian Marius Lillelund, Fernando Pannullo, Morten Opprud Jakobsen, and Christian Fischer Pedersen. Predicting survival time of ball bearings in the presence of censoring. *arXiv preprint arXiv:2309.07188*, 2023.
- Xenia Miscouridou, Adler Perotte, Noémie Elhadad, and Rajesh Ranganath. Deep survival analysis: Nonparametrics and missingness. In *Machine Learning for Healthcare Conference*, pages 244–256. PMLR, 2018.
- Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pages 8162–8171. PMLR, 2021.
- Alexander Norcliffe, Bogdan Cebere, Fergus Imrie, Pietro Lio, and Mihaela van der Schaar. Survivalgan: Generating time-to-event data for survival analysis. In *International Conference on Artificial Intelligence and Statistics*, pages 10279–10304. PMLR, 2023.
- Ross L. Prentice. A generalization of the probit and logit methods for dose response curves. *Biometrics*, pages 761–768, 1976.
- Zhaozhi Qian, Rob Davis, and Mihaela van der Schaar. Synthcity: a benchmark framework for diverse use cases of tabular synthetic data. *Advances in Neural Information Processing Systems*, 36, 2024.
- Rajesh Ranganath, Adler Perotte, Noémie Elhadad, and David Blei. Deep survival analysis. In *Machine Learning for Healthcare Conference*, pages 101–114. PMLR, 2016.
- Stephen Salerno and Yi Li. High-dimensional survival analysis: Methods and applications. *Annual review of statistics and its application*, 10:25–49, 2023.
- M. Schumacher, G. Bastert, H. Bojar, K. Hübner, M. Olschewski, W. Sauerbrei, C. Schmoor, C. Beyerle, R. L. Neumann, and H. F. Rauschecker. Randomized 2 x 2 trial evaluating hormonal treatment and the duration of chemotherapy in node-positive breast cancer patients. german breast cancer study group. *Journal of Clinical Oncology*, 12(10):2086–2093, 1994.

- Ritesh Singh and Keshab Mukhopadhyay. Survival analysis in clinical trials: Basics and must know areas. *Perspectives in clinical research*, 2(4):145–148, 2011.
- Yixuan Su, Tian Lan, Yan Wang, Dani Yogatama, Lingpeng Kong, and Nigel Collier. A contrastive framework for neural text generation. *Advances in Neural Information Processing Systems*, 35:21548–21561, 2022.
- Yu-An Wang and Yun-Nung Chen. What do position embeddings learn? an empirical study of pre-trained language model positional encoding. *arXiv preprint arXiv:2010.04903*, 2020.
- Zhenchen Wang, Puja Myles, and Allan Tucker. Generating and evaluating cross-sectional synthetic electronic healthcare data: Preserving data utility and patient privacy. *Computational Intelligence*, 37(2):819–851, 2021.
- Waloddi Weibull. A statistical distribution function of wide applicability. *Journal of applied mechanics*, 1951.
- Lei Xu, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Modeling tabular data using conditional gan. *Advances in neural information processing systems*, 32, 2019.
- Jinsung Yoon, Lydia N. Drumright, and Mihaela Van Der Schaar. Anonymization through data synthesis using generative adversarial networks (ads-gan). *IEEE journal of biomedical and health informatics*, 24(8):2378–2388, 2020.
- Jun Zhang, Graham Cormode, Cecilia M. Procopiuc, Divesh Srivastava, and Xiaokui Xiao. Privbayes: Private data release via bayesian networks. *ACM Transactions on Database Systems (TODS)*, 42(4):1–41, 2017.
- Zilong Zhao, Aditya Kunar, Robert Birke, and Lydia Y. Chen. Ctab-gan: Effective table data synthesizing. In *Asian Conference on Machine Learning*, pages 97–112. PMLR, 2021.
- Xingyu Zhou, Wen Su, Changyu Liu, Yuling Jiao, Xingqiu Zhao, and Jian Huang. Deep generative survival analysis: Nonparametric estimation of conditional survival function. *arXiv preprint arXiv:2205.09633*, 2022.

Table 5: Training time per-iteration (TTPI) and generation time (GT) for synthetic survival data generation (in seconds).

Metric	Method	AIDS	METABRIC	SUPPORT	GBSG	FLCHAIN
TTPI (\downarrow)	SurvivalGAN	0.178 _{0.004}	0.260 _{0.003}	1.234 _{0.013}	0.283 _{0.005}	1.024 _{0.012}
	TVAE	0.081 _{0.003}	0.126 _{0.002}	0.725 _{0.005}	0.134 _{0.003}	0.527 _{0.004}
	TabDDPM	0.056 _{0.001}	0.048 _{0.002}	0.201 _{0.001}	0.061 _{0.003}	0.183 _{0.002}
	CTGAN	0.174 _{0.002}	0.242 _{0.004}	1.231 _{0.012}	0.226 _{0.003}	0.891 _{0.011}
	ADS-GAN	0.141 _{0.001}	0.235 _{0.003}	1.139 _{0.011}	0.252 _{0.004}	0.827 _{0.013}
	TVAE (u)	0.136 _{0.004}	0.186 _{0.003}	1.023 _{0.014}	0.187 _{0.002}	0.735 _{0.005}
	TabDDPM (u)	0.046 _{0.003}	0.050 _{0.002}	0.215 _{0.003}	0.070 _{0.004}	0.183 _{0.002}
	CTGAN (u)	0.193 _{0.004}	0.282 _{0.003}	1.312 _{0.012}	0.287 _{0.003}	1.028 _{0.011}
	ADS-GAN (u)	0.214 _{0.003}	0.291 _{0.004}	1.404 _{0.015}	0.306 _{0.002}	1.061 _{0.013}
	SurvivalGAN	0.396 _{0.014}	0.421 _{0.023}	0.896 _{0.085}	0.407 _{0.054}	0.715 _{0.053}
GT (\downarrow)	TVAE	0.084 _{0.002}	0.109 _{0.003}	0.365 _{0.031}	0.114 _{0.013}	0.243 _{0.011}
	TabDDPM	11.870 _{0.135}	9.430 _{0.227}	49.611 _{0.284}	17.121 _{0.107}	35.116 _{0.195}
	CTGAN	0.085 _{0.013}	0.089 _{0.014}	0.156 _{0.003}	0.068 _{0.009}	0.103 _{0.009}
	ADS-GAN	0.082 _{0.012}	0.085 _{0.003}	0.149 _{0.004}	0.064 _{0.012}	0.101 _{0.014}
	TVAE (u)	0.128 _{0.023}	0.135 _{0.004}	0.468 _{0.095}	0.124 _{0.003}	0.281 _{0.117}
	TabDDPM (u)	11.785 _{0.367}	9.466 _{0.337}	50.017 _{0.817}	18.085 _{0.567}	34.937 _{0.857}
	CTGAN (u)	0.079 _{0.003}	0.087 _{0.004}	0.192 _{0.035}	0.073 _{0.004}	0.124 _{0.117}
	ADS-GAN (u)	0.089 _{0.004}	0.098 _{0.013}	0.212 _{0.045}	0.085 _{0.013}	0.111 _{0.003}

Appendix A. Ethics Statement

This study focuses on synthetic survival data generation, which has important ethical implications. While our method aims to preserve patient privacy by generating synthetic data, we acknowledge the potential risks of reinforcing biases present in the original datasets. We have made efforts to evaluate our approach across different sub-populations to assess fairness, but further work is needed to fully address bias and equity concerns in survival analysis. The synthetic data generated should not be used for real clinical decision-making without extensive validation. We have no conflicts of interest to declare. All datasets used are publicly available and de-identified. Our work complies with relevant data protection regulations. We encourage users of our method to carefully consider the ethical implications and potential biases when applying it to sensitive healthcare data.

Appendix B. Experimental Details

B.1. Computational Cost

All experiments, except for the LLM fine-tuning (see Section 4), were conducted on Google Colab Pro using a T4 GPU. For the LLM fine-tuning experiments, an NVIDIA A100 GPU was utilized on Colab. In Table 5 we report the training time per iteration (TTPI) along with the time taken for synthetic data generation (GT) for all models used in Section 4.1, while the training and generation time for Section 4.2 are reported in Table 6.

Table 6: Training and generation time for synthetic survival data generation using LLMs (in seconds).

Metric	Method	AIDS	METABRIC	SUPPORT	GBSG	FLCHAIN
TTPI	LLM (Ours)	5.182 _{0.113}	9.584 _{0.195}	49.756 _{0.004}	6.698 _{0.214}	23.367 _{0.015}
GT	LLM	14.237 _{0.153}	121.451 _{0.207}	270.798 _{0.994}	23.516 _{0.054}	77.154 _{0.187}
	LLM (Ours)	623.092 _{2.005}	912.083 _{1.767}	5519.845 _{5.574}	812.298 _{0.256}	1140.475 _{2.697}

Table 7: Summary statistics of the datasets used in the study.

Dataset	No. instances	No. censored instances	No. features
AIDS	1151	96	11
METABRIC	1904	801	9
FLCHAIN	7874	5705	9
GBSG	2232	965	7
SUPPORT	8873	2837	14

B.2. Datasets

We benchmark our methodology on a variety of medical datasets summarized in Table 7. Specifically: *i*) Study to understand prognoses preferences outcomes and risks of treatment (SUPPORT) (Knaus et al., 1995); *ii*) Molecular taxonomy of breast cancer international consortium (METABRIC) (Curtis et al., 2012); *iii*) ACTG 320 clinical trial dataset (AIDS) (Hammer et al., 1997); *iv*) Rotterdam & German breast cancer study group (GBSG) (Schumacher et al., 1994); and *v*) Assay of serum free light chain (FLCHAIN) (Dispenzieri et al., 2012). Pre-processed versions of METABRIC, SUPPORT, and GBSG can be found at: <https://github.com/havakv/pycox>. AIDS and FLCHAIN datasets can be downloaded from <https://github.com/sebp/scikit-survival/tree/master/sksurv/datasets/data>. For the FLCHAIN dataset, missing values in continuous covariates were imputed to the mean, while in discrete covariates they were imputed to the mode. All of these datasets are publicly available hence the experiments can be readily reproduced. In parts of our code (see Section 3.1 and 4), we utilize and modify the Synthcity library (<https://github.com/vanderschaarlab/synthcity>) which is protected under the *Apache-2.0* license. All rights to Synthcity are reserved by the original authors (Qian et al., 2024).

B.3. Hyperparameters

For reproducibility purposes, all hyperparameters are specified below. Table 8 lists the hyperparameters for the downstream survival models used in the benchmarks. Further, Tables 9 and 10 provide the hyperparameters for all generative models employed in the study.

Table 8: Hyperparameters for the survival models used in Section 4.

Method	Parameter	Parameter Value
CoxPH	Estimation Method	Breslow
	Penalizer	0.0
	L^1 Ratio	0.0
SurvivalXGBoost	Objective	Survival: AFT
	Evaluation Metric	AFT Negative Log Likelihood
	AFT Loss Distribution	Normal
	AFT Loss Distribution Scale	1.0
	No. Estimators	100
	Column Subsample Ratio (by node)	0.5
	Maximum Depth	5
	Subsample Ratio	0.5
	Learning Rate	5×10^{-2}
	Minimum Child Weight	50
	Tree Method	Histogram
	Booster	Dart
Deephit	No. Durations:	1000
	Batch Size	100
	Epochs	2000
	Learning Rate	1×10^{-2}
	Hidden Width	300
	α	0.28
	σ	0.38
	Dropout Rate	0.2
	Patience	20

Table 9: Hyperparameters used for the LLM in Section 4.2.

Method	Parameter	Parameter Value
GReaT (DistilGPT2)	Batch Size	32
	No. Iterations	1000
	Learning Rate	5×10^{-5}
	Optimizer	AdamW
	Sampling Temperature	0.7
	Sampling Batch Size	100

Table 10: Hyperparameters of the generative models used in synthetic benchmarks in Section 4.1.

Model	Parameter	Parameter Value
ADS-GAN	No. Iterations	10000
	Generator no. Hidden Layers	2
	Generator Hidden Units	500
	Generator Non-linearity	ReLU
	Generator Dropout Rate	0.1
	Discriminator No. Hidden Layers	2
	Discriminator Hidden Units	500
	Discriminator Non-linearity	Leaky ReLU
	Discriminator Dropout Rate	0.1
	Learning Rate	1×10^{-3}
	Weight Decay	1×10^{-3}
	Batch Size	200
	Gradient Penalty (λ)	10
	Identifiability Penalty	0.1
	Encoder Max Clusters	5
	Early Stopping Patience	5
CTGAN	No. Iterations	2000
	Generator No. Hidden Layers	2
	Generator Hidden Units	500
	Generator Non-linearity	ReLU
	Learning Rate	1×10^{-3}
	Weight Decay	1×10^{-3}
	Discriminator No. Hidden Layers	2
	Discriminator Hidden Units	500
	Discriminator Non-linearity	Leaky ReLU
	Gradient Penalty (λ)	10
	Batch Size	200
	Early Stopping Patience	5
SurvivalGAN	Uncensoring Model	Survival Function Regression
	Time-to-event strategy	Survival Function
	Censoring Strategy	Random
	Dataloader Sampling Strategy	Imbalance Time Censoring
TVAE	No. Iterations	1000
	Batch Size	200
	Learning Rate	1×10^{-3}
	Weight Decay	1×10^{-5}
	Encoder No. Hidden Layers	3
	Encoder Hidden Units	500
	Encoder Non-linearity	Leaky ReLU
	Encoder Dropout Rate	0.1
	Decoder No. Hidden Layers	3
	Decoder Hidden Units	500
	Decoder Non-linearity	Leaky ReLU
	Decoder Dropout Rate	0
	Early Stopping Patience	5
	Data Encoder Max Clusters	10
	Embedding Width	500
TabDDPM	No. Iterations	1000
	Batch Size	1024
	Learning Rate	2×10^{-3}
	Weight Decay	1×10^{-4}
	No. of Time-Steps	1000
	Scheduler	Cosine
	Gaussian Loss Type	MSE

Appendix C. Additional Performance Metrics

Full quality metrics: Below we provide the comprehensive scores of all models evaluated in the paper. Table 11 presents the *covariate quality* and *downstream performance* metrics for all models assessed in Section 4.1. In Table 12, we report the *event-time distribution quality* metrics, including optimism, short-sightedness, and KM Divergence, for both conditional and unconditional models. Although our base method uses DPMM to sample t and e , we also explore directly sampling *both* event times and type from their joint distribution via $p(t|e)$ and $p(e)$, using their empirical distributions. This comparison is shown in Table 13.

Privacy experiment: To explore the acceptability of bootstrapping t and e when generating synthetic data, we employed the Distance to Closest Record (DCR) metric to evaluate the privacy preservation capabilities of various synthetic data generation methods (Zhao et al., 2021). The DCR quantifies the Euclidean distance between each synthetic record and its nearest real counterpart. A higher DCR value indicates a lower risk of privacy breach. We report the median and minimum DCR for all synthetic survival data generators used in our study, with the addition of a Synthetic Minority Oversampling Technique (SMOTE) (Chawla et al., 2002) baseline. SMOTE, originally proposed for minority class oversampling, is a simple interpolation-based method that generates synthetic points as convex combinations of real data points and their k -th nearest neighbors. In this study, we generalized and applied SMOTE to synthetic data generation to bootstrap the entire data point $(\tilde{x}, \tilde{t}, \tilde{e})$, for comparison purposes. The results, presented in Table 14, demonstrate that the median DCR for the methods where t and e were bootstrapped (denoted by a dagger †) was higher in 3 of 5 data sets, although by a small margin when compared to *unconditional* models (denoted by (u)). A similar observation can be made in Table 15 where minimum DCR was higher for our methods in 4 of 5 datasets compared to the *unconditional* models. In general, the median and minimum DCR values were largely similar between the methods when sampling from empirical t and e , sampling t and e from fitted DPMM or generating them along with the covariates as a joint distribution (unconditional) suggesting that sampling them is not likely to impact privacy. However, SMOTE consistently exhibited the lowest DCR in all datasets, indicating potential privacy concerns. These findings provide empirical evidence that bootstrapping t and e is generally acceptable from a privacy perspective. However, we note that even the most stringent minimum DCR does not provide privacy guarantees, so it needs to be interpreted with care.

Table 11: Quality and downstream performance metrics. Models conditioned on empirical t and e are highlighted (\dagger) and u refers to unconditional models. Subscripts are standard deviations for 5 repetitions.

Metric	Method	AIDS	METABRIC	SUPPORT	GBSG	FLCHAIN
JS distance (\downarrow)	SurvivalGAN	0.013 _{0.005}	0.009 _{0.000}	0.008 _{0.004}	0.008 _{0.001}	0.009 _{0.005}
	TVAE \dagger	0.007 _{0.004}	0.008 _{0.005}	0.004 _{0.000}	0.005 _{0.001}	0.002 _{0.005}
	TabDDPM \dagger	0.007 _{0.000}	0.007 _{0.001}	0.013 _{0.005}	0.005 _{0.002}	0.001 _{0.004}
	CTGAN \dagger	0.013 _{0.001}	0.020 _{0.015}	0.005 _{0.003}	0.003 _{0.002}	0.004 _{0.000}
	ADS-GAN \dagger	0.006 _{0.002}	0.009 _{0.005}	0.005 _{0.001}	0.004 _{0.004}	0.010 _{0.015}
	TVAE	0.008 _{0.002}	0.009 _{0.000}	0.004 _{0.005}	0.005 _{0.002}	0.003 _{0.004}
	TabDDPM	0.006 _{0.005}	0.009 _{0.004}	0.014 _{0.002}	0.004 _{0.000}	0.001 _{0.003}
	CTGAN	0.013 _{0.004}	0.007 _{0.002}	0.005 _{0.005}	0.004 _{0.002}	0.003 _{0.001}
	ADS-GAN	0.006 _{0.003}	0.006 _{0.002}	0.005 _{0.004}	0.005 _{0.005}	0.005 _{0.000}
	TVAE (u)	0.011 _{0.002}	0.009 _{0.005}	0.007 _{0.001}	0.007 _{0.004}	0.003 _{0.000}
	TabDDPM (u)	0.006 _{0.001}	0.007 _{0.004}	0.006 _{0.002}	0.005 _{0.005}	0.002 _{0.000}
	CTGAN (u)	0.007 _{0.005}	0.012 _{0.002}	0.005 _{0.004}	0.008 _{0.001}	0.005 _{0.003}
	ADS-GAN (u)	0.006 _{0.002}	0.007 _{0.000}	0.007 _{0.005}	0.005 _{0.002}	0.005 _{0.004}
	SurvivalGAN	0.112 _{0.015}	0.039 _{0.002}	0.043 _{0.004}	0.019 _{0.005}	0.052 _{0.000}
WS distance (\downarrow)	TVAE \dagger	0.061 _{0.004}	0.028 _{0.005}	0.032 _{0.002}	0.013 _{0.000}	0.016 _{0.005}
	TabDDPM \dagger	0.159 _{0.025}	0.089 _{0.004}	0.308 _{0.025}	0.056 _{0.002}	0.028 _{0.005}
	CTGAN \dagger	0.095 _{0.002}	0.133 _{0.015}	0.034 _{0.005}	0.013 _{0.004}	0.019 _{0.000}
	ADS-GAN \dagger	0.082 _{0.005}	0.037 _{0.002}	0.036 _{0.004}	0.011 _{0.005}	0.018 _{0.004}
	TVAE	0.063 _{0.004}	0.037 _{0.005}	0.032 _{0.002}	0.011 _{0.004}	0.016 _{0.002}
	TabDDPM	0.160 _{0.005}	0.104 _{0.002}	0.309 _{0.004}	0.037 _{0.005}	0.028 _{0.000}
	CTGAN	0.097 _{0.002}	0.037 _{0.004}	0.035 _{0.005}	0.014 _{0.002}	0.018 _{0.004}
	ADS-GAN	0.079 _{0.005}	0.030 _{0.002}	0.034 _{0.000}	0.011 _{0.005}	0.017 _{0.002}
	TVAE (u)	0.075 _{0.004}	0.031 _{0.005}	0.037 _{0.002}	0.013 _{0.004}	0.017 _{0.005}
	TabDDPM (u)	0.079 _{0.005}	0.031 _{0.002}	0.049 _{0.004}	0.015 _{0.000}	0.016 _{0.005}
	CTGAN (u)	0.069 _{0.000}	0.041 _{0.004}	0.036 _{0.005}	0.017 _{0.002}	0.021 _{0.004}
	ADS-GAN (u)	0.065 _{0.005}	0.035 _{0.000}	0.038 _{0.004}	0.013 _{0.005}	0.017 _{0.002}
	SurvivalGAN	0.735 _{0.005}	0.625 _{0.000}	0.602 _{0.004}	0.668 _{0.005}	0.870 _{0.002}
	TVAE \dagger	0.737 _{0.004}	0.612 _{0.005}	0.583 _{0.000}	0.672 _{0.001}	0.872 _{0.005}
C-Index (\uparrow)	TabDDPM \dagger	0.660 _{0.075}	0.589 _{0.015}	0.536 _{0.005}	0.663 _{0.002}	0.876 _{0.004}
	CTGAN \dagger	0.746 _{0.005}	0.628 _{0.015}	0.577 _{0.004}	0.665 _{0.015}	0.874 _{0.002}
	ADS-GAN \dagger	0.797 _{0.015}	0.655 _{0.005}	0.623 _{0.002}	0.684 _{0.004}	0.880 _{0.005}
	TVAE	0.783 _{0.025}	0.630 _{0.005}	0.602 _{0.004}	0.672 _{0.002}	0.868 _{0.005}
	TabDDPM	0.670 _{0.015}	0.603 _{0.004}	0.530 _{0.005}	0.659 _{0.002}	0.875 _{0.004}
	CTGAN	0.760 _{0.005}	0.629 _{0.004}	0.602 _{0.002}	0.668 _{0.005}	0.874 _{0.000}
	ADS-GAN	0.785 _{0.025}	0.652 _{0.005}	0.626 _{0.004}	0.682 _{0.002}	0.880 _{0.005}
	TVAE (u)	0.735 _{0.004}	0.646 _{0.002}	0.604 _{0.005}	0.671 _{0.004}	0.878 _{0.005}
	TabDDPM (u)	0.759 _{0.005}	0.649 _{0.004}	0.625 _{0.002}	0.679 _{0.005}	0.879 _{0.004}
	CTGAN (u)	0.779 _{0.002}	0.647 _{0.005}	0.606 _{0.004}	0.679 _{0.002}	0.878 _{0.005}
	ADS-GAN (u)	0.776 _{0.004}	0.636 _{0.005}	0.601 _{0.002}	0.663 _{0.004}	0.878 _{0.002}
	SurvivalGAN	0.068 _{0.005}	0.205 _{0.004}	0.202 _{0.002}	0.212 _{0.005}	0.096 _{0.004}
	TVAE \dagger	0.059 _{0.004}	0.199 _{0.005}	0.207 _{0.004}	0.214 _{0.002}	0.095 _{0.005}
	TabDDPM \dagger	0.063 _{0.005}	0.212 _{0.004}	0.217 _{0.002}	0.215 _{0.005}	0.096 _{0.004}
Brier Score (\downarrow)	CTGAN \dagger	0.061 _{0.004}	0.199 _{0.005}	0.205 _{0.004}	0.215 _{0.015}	0.089 _{0.005}
	ADS-GAN \dagger	0.059 _{0.005}	0.197 _{0.004}	0.198 _{0.002}	0.213 _{0.005}	0.084 _{0.004}
	TVAE	0.060 _{0.004}	0.197 _{0.005}	0.208 _{0.004}	0.211 _{0.005}	0.091 _{0.002}
	TabDDPM	0.061 _{0.005}	0.209 _{0.004}	0.218 _{0.005}	0.210 _{0.004}	0.091 _{0.002}
	CTGAN	0.060 _{0.004}	0.199 _{0.005}	0.204 _{0.002}	0.211 _{0.004}	0.089 _{0.005}
	ADS-GAN	0.061 _{0.005}	0.195 _{0.004}	0.198 _{0.002}	0.214 _{0.005}	0.085 _{0.004}
	TVAE (u)	0.061 _{0.004}	0.204 _{0.005}	0.206 _{0.004}	0.210 _{0.002}	0.093 _{0.005}
	TabDDPM (u)	0.060 _{0.005}	0.200 _{0.004}	0.199 _{0.002}	0.207 _{0.005}	0.087 _{0.004}
	CTGAN (u)	0.064 _{0.004}	0.202 _{0.005}	0.203 _{0.004}	0.210 _{0.002}	0.086 _{0.005}
	ADSGAN (u)	0.061 _{0.005}	0.207 _{0.004}	0.201 _{0.002}	0.208 _{0.005}	0.088 _{0.004}

Table 12: Event-time distribution quality metrics. Models conditioned on empirical t and e are highlighted (\dagger) and u refers to unconditional models. Subscripts are standard deviations for 5 repetitions.

Metric	Method	AIDS	METABRIC	SUPPORT	GBSG	FLCHAIN
Optimism ($\rightarrow 0$)	SurvivalGAN	0.021 _{0.005}	0.011 _{0.002}	0.016 _{0.004}	0.006 _{0.003}	0.134 _{0.005}
	TVAE \dagger	0.000 _{0.001}	0.000 _{0.002}	0.000 _{0.000}	0.003 _{0.001}	0.001 _{0.002}
	TabDDPM \dagger	0.000 _{0.001}	0.000 _{0.002}	0.000 _{0.000}	0.003 _{0.001}	0.001 _{0.002}
	CTGAN \dagger	0.000 _{0.001}	0.000 _{0.002}	0.000 _{0.000}	0.003 _{0.001}	0.001 _{0.002}
	ADSGAN \dagger	0.000 _{0.001}	0.000 _{0.002}	0.000 _{0.000}	0.003 _{0.001}	0.001 _{0.002}
	TVAE	0.001 _{0.001}	0.001 _{0.001}	-0.001 _{0.001}	-0.004 _{0.000}	-0.003 _{0.001}
	TabDDPM	0.001 _{0.001}	0.001 _{0.001}	-0.001 _{0.001}	-0.004 _{0.000}	-0.003 _{0.001}
	CTGAN	0.001 _{0.001}	0.001 _{0.001}	-0.001 _{0.001}	-0.004 _{0.000}	-0.003 _{0.001}
	ADS-GAN	0.001 _{0.001}	0.001 _{0.001}	-0.001 _{0.001}	-0.004 _{0.000}	-0.003 _{0.001}
	TVAE (u)	0.023 _{0.003}	-0.003 _{0.004}	-0.014 _{0.002}	0.004 _{0.005}	0.022 _{0.003}
	TabDDPM (u)	0.021 _{0.004}	0.001 _{0.005}	0.001 _{0.003}	0.026 _{0.002}	0.005 _{0.004}
	CTGAN (u)	-0.005 _{0.005}	0.017 _{0.002}	-0.038 _{0.004}	0.060 _{0.003}	-0.037 _{0.005}
	ADSGAN (u)	0.002 _{0.002}	-0.033 _{0.003}	-0.007 _{0.005}	0.010 _{0.004}	0.005 _{0.002}
Short Sightedness ($\rightarrow 0$)	SurvivalGAN	0.007 _{0.003}	0.124 _{0.004}	0.020 _{0.002}	0.019 _{0.005}	0.005 _{0.003}
	TVAE \dagger	0.001 _{0.000}	0.000 _{0.000}	0.000 _{0.001}	0.010 _{0.012}	0.002 _{0.001}
	TabDDPM \dagger	0.001 _{0.000}	0.000 _{0.000}	0.000 _{0.001}	0.010 _{0.012}	0.002 _{0.001}
	CTGAN \dagger	0.001 _{0.000}	0.000 _{0.000}	0.000 _{0.001}	0.010 _{0.012}	0.002 _{0.001}
	ADS-GAN \dagger	0.001 _{0.000}	0.000 _{0.000}	0.000 _{0.001}	0.010 _{0.012}	0.002 _{0.001}
	TVAE	0.001 _{0.001}	0.000 _{0.001}	0.000 _{0.002}	-0.012 _{0.001}	-0.002 _{0.002}
	TabDDPM	0.001 _{0.001}	0.000 _{0.001}	0.000 _{0.002}	-0.012 _{0.001}	-0.002 _{0.002}
	CTGAN	0.001 _{0.001}	0.000 _{0.001}	0.000 _{0.002}	-0.012 _{0.001}	-0.002 _{0.002}
	ADS-GAN	0.001 _{0.001}	0.000 _{0.001}	0.000 _{0.002}	-0.012 _{0.001}	-0.002 _{0.002}
	TVAE (u)	0.058 _{0.004}	0.148 _{0.005}	0.002 _{0.002}	0.017 _{0.003}	0.018 _{0.004}
	TabDDPM (u)	0.002 _{0.005}	0.000 _{0.002}	0.002 _{0.003}	0.015 _{0.004}	0.003 _{0.005}
	CTGAN (u)	0.071 _{0.002}	0.056 _{0.003}	0.010 _{0.004}	0.019 _{0.005}	0.017 _{0.002}
	ADSGAN (u)	0.040 _{0.003}	0.188 _{0.004}	0.002 _{0.005}	0.014 _{0.002}	0.006 _{0.003}
KM Divergence (\downarrow)	SurvivalGAN	0.021 _{0.004}	0.082 _{0.005}	0.064 _{0.002}	0.049 _{0.003}	0.134 _{0.004}
	TVAE \dagger	0.002 _{0.000}	0.008 _{0.001}	0.002 _{0.000}	0.005 _{0.001}	0.002 _{0.000}
	TabDDPM \dagger	0.002 _{0.000}	0.008 _{0.001}	0.002 _{0.000}	0.005 _{0.001}	0.002 _{0.000}
	CTGAN \dagger	0.002 _{0.000}	0.008 _{0.001}	0.002 _{0.000}	0.005 _{0.001}	0.002 _{0.000}
	ADS-GAN \dagger	0.002 _{0.000}	0.008 _{0.001}	0.002 _{0.000}	0.005 _{0.001}	0.002 _{0.000}
	TVAE	0.003 _{0.001}	0.013 _{0.002}	0.006 _{0.001}	0.007 _{0.001}	0.005 _{0.001}
	TabDDPM	0.003 _{0.001}	0.013 _{0.002}	0.006 _{0.001}	0.007 _{0.001}	0.005 _{0.001}
	CTGAN	0.003 _{0.001}	0.013 _{0.002}	0.006 _{0.001}	0.007 _{0.001}	0.005 _{0.001}
	ADS-GAN	0.003 _{0.001}	0.013 _{0.002}	0.006 _{0.001}	0.007 _{0.001}	0.005 _{0.001}
	TVAE (u)	0.031 _{0.005}	0.042 _{0.002}	0.025 _{0.003}	0.027 _{0.004}	0.031 _{0.005}
	TabDDPM (u)	0.021 _{0.002}	0.019 _{0.003}	0.011 _{0.004}	0.026 _{0.005}	0.007 _{0.002}
	CTGAN (u)	0.015 _{0.003}	0.028 _{0.004}	0.038 _{0.005}	0.061 _{0.002}	0.037 _{0.003}
	ADSGAN (u)	0.016 _{0.004}	0.039 _{0.005}	0.020 _{0.002}	0.030 _{0.003}	0.012 _{0.004}

Table 13: Quality, downstream and event-time metrics. Performance reported for our best model between ADS-GAN, TVAE, CTGAN and TabDDPM from Table 11 and Table 12 for both settings. Models conditioned on empirical t and e are highlighted (\dagger). Subscripts are standard deviations for 5 repetitions.

Metric	Method	AIDS	METABRIC	SUPPORT	GBSG	FLCHAIN
JS distance (\downarrow)	Ours	0.006 _{0.005}	0.006 _{0.002}	0.004 _{0.005}	0.004 _{0.002}	0.001 _{0.003}
	Ours \dagger	0.006 _{0.002}	0.007 _{0.001}	0.004 _{0.000}	0.003 _{0.002}	0.001 _{0.004}
WS distance (\downarrow)	Ours	0.063 _{0.004}	0.030 _{0.002}	0.032 _{0.002}	0.011 _{0.004}	0.016 _{0.002}
	Ours \dagger	0.061 _{0.004}	0.028 _{0.005}	0.032 _{0.002}	0.011 _{0.005}	0.016 _{0.005}
C-Index (\uparrow)	Ours	0.785 _{0.025}	0.652 _{0.005}	0.626 _{0.004}	0.682 _{0.002}	0.880 _{0.005}
	Ours \dagger	0.797 _{0.015}	0.655 _{0.005}	0.623 _{0.002}	0.684 _{0.004}	0.880 _{0.005}
Brier Score (\downarrow)	Ours	0.060 _{0.004}	0.197 _{0.005}	0.198 _{0.002}	0.210 _{0.004}	0.085 _{0.004}
	Ours \dagger	0.059 _{0.005}	0.197 _{0.004}	0.198 _{0.002}	0.213 _{0.005}	0.084 _{0.004}
Optimism ($\rightarrow 0$)	Ours	0.001 _{0.001}	0.001 _{0.001}	-0.001 _{0.001}	-0.004 _{0.000}	-0.003 _{0.001}
	Ours \dagger	0.000 _{0.001}	0.000 _{0.002}	0.000 _{0.000}	0.003 _{0.001}	0.001 _{0.002}
Short Sightedness ($\rightarrow 0$)	Ours	0.001 _{0.001}	0.000 _{0.001}	0.000 _{0.002}	-0.012 _{0.001}	-0.002 _{0.002}
	Ours \dagger	0.001 _{0.000}	0.000 _{0.000}	0.000 _{0.001}	0.010 _{0.012}	0.002 _{0.001}
KM Divergence (\downarrow)	Ours	0.003 _{0.001}	0.013 _{0.002}	0.006 _{0.001}	0.007 _{0.001}	0.005 _{0.001}
	Ours \dagger	0.002 _{0.000}	0.008 _{0.001}	0.002 _{0.000}	0.005 _{0.001}	0.002 _{0.000}

Table 14: Median value of Distance of closest record from the original. Models conditioned on empirical t and e are highlighted (\dagger) and u refers to unconditional models. Subscripts are standard deviations for 5 repetitions. The highest (best) values are in bold and the least (worst) values are underlined.

Metric	Method	AIDS	METABRIC	SUPPORT	GBSG	FLCHAIN
Median DCR	SurvivalGAN	1.035 _{0.001}	0.969 _{0.004}	1.589 _{0.002}	0.500 _{0.003}	0.796 _{0.004}
	TVAE \dagger	0.883 _{0.002}	0.877 _{0.003}	1.511 _{0.002}	0.476 _{0.002}	0.642 _{0.003}
	TabDDPM \dagger	1.172 _{0.031}	0.908 _{0.015}	1.612 _{0.035}	0.519 _{0.004}	0.572 _{0.013}
	CTGAN \dagger	0.918 _{0.013}	1.043 _{0.004}	1.594 _{0.001}	0.524 _{0.003}	0.695 _{0.004}
	ADS-GAN \dagger	1.133 _{0.170}	0.992 _{0.006}	1.691 _{0.004}	0.519 _{0.003}	0.667 _{0.014}
	SMOTE	<u>0.388</u> _{0.002}	<u>0.698</u> _{0.001}	<u>0.958</u> _{0.003}	<u>0.290</u> _{0.005}	<u>0.381</u> _{0.006}
	TVAE	0.992 _{0.002}	0.881 _{0.001}	1.523 _{0.004}	0.481 _{0.002}	0.648 _{0.003}
	TabDDPM	1.184 _{0.034}	0.915 _{0.015}	1.625 _{0.035}	0.522 _{0.004}	0.578 _{0.013}
	CTGAN	1.007 _{0.012}	1.052 _{0.001}	1.602 _{0.003}	0.521 _{0.002}	0.689 _{0.004}
	ADS-GAN	1.142 _{0.175}	0.998 _{0.003}	1.698 _{0.007}	0.523 _{0.003}	0.672 _{0.014}
	TVAE (u)	1.044 _{0.023}	0.813 _{0.004}	1.405 _{0.003}	0.432 _{0.003}	0.553 _{0.001}
	TabDDPM (u)	1.020 _{0.013}	1.087 _{0.004}	1.611 _{0.006}	0.477 _{0.004}	0.567 _{0.003}
	CTGAN (u)	1.112 _{0.014}	1.001 _{0.003}	1.586 _{0.001}	0.515 _{0.003}	0.641 _{0.004}
	ADS-GAN (u)	1.158 _{0.011}	0.945 _{0.002}	1.666 _{0.003}	0.475 _{0.004}	0.533 _{0.001}

Table 15: Minimum value of Distance of closest record from the original. Models conditioned on empirical t and e are highlighted (\dagger) and u refers to unconditional models. Subscripts are standard deviations for 5 repetitions. The highest (best) values are in bold and the least (worst) values are underlined.

Metric	Method	AIDS	METABRIC	SUPPORT	GBSG	FLCHAIN
Minimum DCR	SurvivalGAN	0.048 _{0.003}	0.172 _{0.004}	0.326 _{0.003}	0.062 _{0.002}	0.057 _{0.003}
	TVAE \dagger	0.077 _{0.034}	0.202 _{0.025}	0.370 _{0.023}	0.033 _{0.004}	0.026 _{0.003}
	TabDDPM \dagger	0.095 _{0.003}	0.193 _{0.055}	0.403 _{0.013}	0.065 _{0.004}	0.037 _{0.002}
	CTGAN \dagger	0.139 _{0.015}	0.215 _{0.014}	0.321 _{0.015}	0.045 _{0.013}	0.054 _{0.014}
	ADS-GAN \dagger	0.102 _{0.013}	0.185 _{0.043}	0.391 _{0.015}	0.053 _{0.015}	0.066 _{0.025}
	SMOTE	<u>0.000</u> _{0.000}	<u>0.000</u> _{0.001}	<u>0.000</u> _{0.002}	<u>0.000</u> _{0.001}	<u>0.000</u> _{0.002}
	TVAE	0.081 _{0.014}	0.208 _{0.011}	0.378 _{0.024}	0.035 _{0.003}	0.028 _{0.004}
	TabDDPM	0.098 _{0.004}	0.198 _{0.040}	0.412 _{0.015}	0.068 _{0.003}	0.039 _{0.003}
	CTGAN	0.135 _{0.011}	0.211 _{0.016}	0.325 _{0.017}	0.048 _{0.015}	0.058 _{0.016}
	ADS-GAN	0.106 _{0.015}	0.189 _{0.045}	0.395 _{0.017}	0.056 _{0.017}	0.069 _{0.021}
	TVAE (u)	0.083 _{0.013}	0.154 _{0.037}	0.171 _{0.014}	0.031 _{0.003}	0.028 _{0.004}
	TabDDPM (u)	0.090 _{0.035}	0.213 _{0.004}	0.337 _{0.013}	0.054 _{0.004}	0.033 _{0.003}
	CTGAN (u)	0.109 _{0.025}	0.194 _{0.023}	0.316 _{0.025}	0.046 _{0.015}	0.024 _{0.004}
	ADS-GAN (u)	0.062 _{0.037}	0.205 _{0.035}	0.429 _{0.004}	0.050 _{0.004}	0.055 _{0.003}