000
001UNIVERSALTIME-SERIESGENERATIONUSING002
003SCORE-BASEDGENERATIVEMODELS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

Paper under double-blind review

Abstract

Score-based generative models (SGMs) have demonstrated unparalleled sampling quality and diversity in numerous fields, such as image generation, voice synthesis, and tabular data synthesis, etc. Inspired by those outstanding results, we apply SGMs to synthesize time-series by learning its conditional score function. To this end, we present a conditional score network for time-series synthesis, deriving a denoising score matching loss tailored for our purposes. In particular, our presented denoising score matching loss is the conditional denoising score matching loss for time-series synthesis. In addition, our framework is such flexible that both regular and irregular time-series can be synthesized with minimal changes to our model design. Finally, we obtain exceptional synthesis performance on various time-series datasets, achieving state-of-the-art sampling diversity and quality.

1 INTRODUCTION

Time-series frequently occurs in our daily lives, e.g., stock data, climate data, and so on. Especially, time-series forecasting and classification are popular research topics in the field of machine learning (Ahmed et al., 2010; Fu, 2011; Ismail Fawaz et al., 2019). In many cases, however, time-series samples are incomplete and/or the number of samples is insufficient, in which case training machine learning models cannot be fulfilled in a robust way. To overcome the limitation, time-series synthesis has been studied actively recently (Chen et al., 2018; Dash et al., 2020). These synthesis models have been designed in various ways, including variational autoencoders (VAEs) and generative adversarial networks (GANs) (Desai et al., 2021; Yoon et al., 2019; Jeon et al., 2022).

Moreover, real-world time series often inevitably contain missing values because of privacy reasons or medical events. It has been noted that the missing values involve important information, called *informative missingness* (Rubin, 1976). To remedy the problem, several works have been developed to deal with irregular time-series, i.e., the inter-arrival time between observations is not fixed and/or some observations can be missing (Schafer and Graham, 2002; Che et al., 2016; Kidger et al., 2020), in which case synthesizing irregular time series is challenging (Jeon et al., 2022).

 Score-based generative models (SGMs) have shown good sampling quality and diversity in numerous fields, such as image generation, voice synthesis, and tabular data synthesis, etc (Yang et al., 2023). However, in time-series generation, SGMs should consider *autoregressiveness*, meaning each time-series observation is generated in consideration of its previously generated observations which makes time-series generation more difficult (Yoon et al., 2019; Jeon et al., 2022). To this end, we propose the method of <u>T</u>ime-series generation using conditional <u>S</u>core-based <u>G</u>enerative <u>M</u>odel (TSGM), which consists of three neural networks, i.e., an encoder, a score network, and a decoder (see Figure 2).

O47 Score-based time-series synthesis SGMs have its potential to deal with autoregressiveness (Meng et al., 2020). We design our own autoregressive denoising score matching loss for time-series generation and prove its correctness (see Section 2.1.3). Along with its performances from image generation domain (cf. Fig. 1), here is our main motivation of using SGMs:

- ⁰⁵¹ The mathematical structure of SGMs naturally aligns with the sequential nature of time-series data.
- 053 Besides, we design a conditional score network on time-series, which learns the gradient of the conditional log-likelihood.

Table 1: The table illustrates how many medals
each method gets across all datasets and evaluation metrics, based on the generation evaluation scores presented in Table 2 and Table 12. Our method with the two specific types, TSGM-VP and TSGM-subVP, achieves superior generation performance compared to baselines.

	Olympic Rankings						
Method	Gold			Silver		Bronze	
	Regular	Irregular	R	Ι	R	Ι	
TSGM-VP	4	11	4	11	0	1	
TSGM-subVP	6	16	1	7	1	0	
TimeGAN	1	0	0	0	1	0	
TimeVAE	0	0	0	0	1	4	
GT-GAN	0	1	1	1	2	16	



Figure 1: The KDE plots show the estimated distributions of original data and ones generated by several methods in the Air and AI4I datasets we ignore time stamps for drawing these distributions. Unlike baseline methods, the distribution of TSGM-VP is almost identical to the original one. These figures provide an evidence of the excellent generation quality and diversity of our method. For TSGM-subVP, similar results are observed. Refer to Appendix L for additional visualizations

Regular vs. irregular time-series synthesis Time-series are prevalent throughout our daily lives. 073 Although regular time-series are easy to be considered, it is common to have missing values because 074 of privacy issue (Che et al., 2016; Kidger et al., 2020), which makes irregularly sampled time-series¹. 075 This hinders time-series processing and we consider the hardest problem condition: providing com-076 plete time-series given not only regular time-series, but also irregular ones. To this end, our method 077 is considered 'universial' in that both both regular and irregular time-series samples can be treated with minimal changes to our model design. For synthesizing regular time series, we use a recurrent 079 neural network-based encoder and decoder. Continuous-time methods, such as neural controlled differential equations (Kidger et al., 2020) and GRU-ODE (Brouwer et al., 2019), can be used as 081 our encoder and decoder for synthesizing irregular time series (see Section 3.2 and Appendix I, J).

082 We conduct in-depth experiments with 4 real-world datasets under regular and irregular settings. 083 To be specific, for the irregular settings, we randomly drop 30%, 50%, and 70% of observations 084 from regular time-series on training, then generate complete time-series on evaluation. Therefore, 085 we test with 16 different settings, i.e., 4 datasets for one regular and three irregular settings. Our specific choices of 9 baselines include almost all existing types of time-series generative paradigms, 087 ranging from VAEs to GANs. In Table 1 and Figure 1, we compare our method to the baselines, ranking methods by their evaluation scores and estimating data distribution by kernel density esti-880 mation (KDE). We also visualize real and generated time-series samples onto a latent space using 089 t-SNE (van der Maaten and Hinton, 2008) in Figure 3. Our proposed method shows the best gen-090 eration quality in almost all cases. Furthermore, the t-SNE and KDE visualization results provide 091 intuitive evidence that our method's generation diversity is also superior to that of the baselines. Our 092 contributions are summarized as follows:

094 095

096

098

099

063

064

065

066

067

068

069

070 071

- 1. We, for the first time, propose an SGM-based universal time-series synthesis method.
- 2. We derive our own denoising score matching loss considering the autoregressive nature of sequential data, connecting SGMs to time-series generation domain.
- 3. We conduct comprehensive experiments with 4 real-world datasets and 9 baselines under one regular and three irregular settings since our method supports both regular and irregular time-series. Overall, our proposed method shows the best generation quality and diversity.

103 104 105

 ¹For example, Physionet (Goldberger et al., 2000 (June 13), a famous dataset for time series classification, deliberately removed 90% of observations to protect the privacy of patients, posing challenges for learning and analysis. The synthesized time series can be used instead.

108 2 RELATED WORK AND PRELIMINARIES

110 2.1 SCORE-BASED GENERATIVE MODELS

SGMs offer several advantages over other generative models, including their higher generation quality and diversity. SGMs follow a two-step process, wherein i) gaussian noises are continuously added to a sample and ii) then removed to recover a new sample. These processes are known as the forward and reverse processes, respectively. In this section, we provide a brief overview of the original SGMs in (Song et al., 2021), which will be adapted for the time-series generation tasks.

117 118 2.1.1 FORWARD AND REVERSE PROCESS

At first, SGMs add noises with the following stochastic differential equation (SDE):

139 140

148

153 154 $d\mathbf{x}^{s} = \mathbf{f}(s, \mathbf{x}^{s})ds + g(s)d\mathbf{w}, s \in [0, 1],$ (1)

where $\mathbf{w} \in \mathbb{R}^{\dim(\mathbf{x})}$ is a multi-dimensional Brownian motion, $\mathbf{f}(s, \cdot) : \mathbb{R}^{\dim(\mathbf{x})} \to \mathbb{R}^{\dim(\mathbf{x})}$ is a vectorvalued drift term, and $g : [0, 1] \to \mathbb{R}$ is a scalar-valued diffusion function. Hereafter, we define \mathbf{x}^s as a noisy sample diffused at time $s \in [0, 1]$ from an original sample $\mathbf{x} \in \mathbb{R}^{\dim(\mathbf{x})}$. Therefore, \mathbf{x}^s can be understood as a stochastic process following the SDE.

There are several options for **f** and *g*: variance exploding(VE), variance preserving(VP), and subVP. Song et al. (2021) proved that VE and VP are continuous generalizations of the two discrete diffusion methods: one in Song and Ermon (2019) and the other in Sohl-Dickstein et al. (2015); Ho et al. (2020). The subVP method shows, in general, better negative log-likelihood (NLL) according to Song et al. (2021). We describe the exact form of each SDE in Table 14 with detailed explanation in Appendix N. Note that we only use the subVP-based TSGM in our main experiments and exclude the VE and VP-based one for its inferiority for time series synthesis in our experiments, but checked the VP-based method from ablation study for its better performances than that of the VE.

SGMs run the forward SDE with a sufficiently large number of steps to make sure that the diffused sample converges to a Gaussian distribution at the final step. The score network $M_{\theta}(s, \mathbf{x}^s)$ learns the gradient of the log-likelihood $\nabla_{\mathbf{x}^s} \log p(\mathbf{x}^s)$, which will be used in the reverse process.

138 For the forward SDE, there exists the following corresponding reverse SDE (Anderson, 1982):

$$d\mathbf{x}^{s} = [\mathbf{f}(s, \mathbf{x}^{s}) - g^{2}(s)\nabla_{\mathbf{x}^{s}}\log p(\mathbf{x}^{s})]ds + g(s)d\bar{\mathbf{w}}.$$
(2)

The formula suggests that if knowing the score function, $\nabla_{\mathbf{x}^s} \log p(\mathbf{x}^s)$, we can recover real samples from the prior distribution $p_1(\mathbf{x}) \sim \mathcal{N}(\mu, \sigma^2)$, where μ, σ vary depending on the forward SDE type.

144 2.1.2 TRAINING AND SAMPLING

146 In order for the model M to learn the score function, the model has to optimize the following loss 147 function:

$$L(\theta) = \mathbb{E}_{s}\{\lambda(s)\mathbb{E}_{\mathbf{x}^{s}}[\|M_{\theta}(s, \mathbf{x}^{s}) - \nabla_{\mathbf{x}^{s}}\log p(\mathbf{x}^{s})\|_{2}^{2}]\},\tag{3}$$

where s is uniformly sampled over [0, 1] with an appropriate weight function $\lambda(s) : [0, 1] \rightarrow \mathbb{R}$. However, using the above formula is computationally prohibitive (Hyvärinen, 2005; Song et al., 2019). Thanks to Vincent (2011), the loss can be substituted with the following denoising score matching loss:

$$L^*(\theta) = \mathbb{E}_s\{\lambda(s)\mathbb{E}_{\mathbf{x}^0}\mathbb{E}_{\mathbf{x}^s|\mathbf{x}^0}[\left\|M_{\theta}(s,\mathbf{x}^s) - \nabla_{\mathbf{x}^s}\log p(\mathbf{x}^s|\mathbf{x}^0)\right\|_2^2]\}.$$
(4)

Since SGMs use an affine drift term, the transition kernel $p(\mathbf{x}^s | \mathbf{x}^0)$ follows a certain Gaussian distribution (Särkkä and Solin, 2019) and therefore, $\nabla_{\mathbf{x}^s} \log p(\mathbf{x}^s | \mathbf{x}^0)$ can be analytically calculated.

158 2.1.3 AUTOREGRESSIVENESS OF DIFFUSION MODEL159

Let $\mathbf{x}_{1:N}$ be a time-series sample which consists of N observations. In order to synthesize timeseries $\mathbf{x}_{1:N}$, unlike other generation tasks, we must consider autoregressiveness: generating each observation \mathbf{x}_n at sequential order $n \in \{2, ..., N\}$ considering its previous history $\mathbf{x}_{1:n-1}$. One



Figure 2: The overall workflow of TSGM (see Section 3.3). Our original learning objective is 175 to approximate $\nabla \log p(\mathbf{x}_{1:n}^s | \mathbf{x}_{1:n-1}^0)$, which is computationally prohibitive, with the conditional 176 score network $M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0})$ using an MSE loss. We then prove in Thm. 3.1 that learning $\nabla \log p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0})$ is equivalent to $\nabla \log p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n-1}^{0})$ for θ of M_{θ} in the MSE loss, i.e., their 177 178 optimal model parameter θ is identical. At the end, our score network $M_{\theta}(s, \mathbf{h}_n^s, \mathbf{h}_{n-1}^0)$ learns 179 $\nabla \log p(\mathbf{h}_n^s | \mathbf{h}_n)$ since RNNs can encode $\mathbf{x}_{1:n}^0$ and $\mathbf{x}_{1:n-1}^0$ into their hidden states \mathbf{h}_n^0 and \mathbf{h}_{n-1}^0 , respectively. 181

182 can train neural networks to learn the conditional likelihood $p(\mathbf{x}_n | \mathbf{x}_{1:n-1})$ and generate each \mathbf{x}_n 183 recursively using it, depicting so-called autoregressive property of time-series domain.

Up to our survey, there's no paper about diffusion models considering autoregressiveness in time-185 series generation. Instead of it, Meng et al. (2020) dealt with a more generalized case, autoregres-186 sive conditional score matching to make its score model, $M_{\theta}(s, \cdot, \cdot)$, to be $M_{\theta}(s, \mathbf{x}_n^s, \mathbf{x}_{1:n-1}^o) \sim$ 187 $\nabla_{\mathbf{x}_n^s} \log p(\mathbf{x}_n^s | \mathbf{x}_{1:n-1}^0)$. Inspired by the previous work, we enlighten a connection between SGMs 188 and time-series generation by deriving an autoregressive denoising score matching (see Theo-189 rem 3.1). We emphasize that TSGM is fundamentally different in that (i) TSGM uses its own 190 autoregressive denoising score matching loss, (ii) as a result, it optimizes not $\nabla_{\mathbf{x}_n^s} \log p(\mathbf{x}_n^s | \mathbf{x}_{1:n-1}^0)$ but $\nabla_{\mathbf{x}_{1:n}^s} \log p(\mathbf{x}_{1:n}^s | \mathbf{x}_{1:n-1}^0)$, which is fit for its RNN-based (markovian) encoder and decoder. 191

2.2 TIME-SERIES GENERATION 193

194 There are several time-series generation papers, and we introduce their ideas. TimeVAE (Desai 195 et al., 2021) is a variational autoencoder to synthesize time-series data. This model can provide 196 interpretable results by reflecting temporal structures such as trend and seasonality in the generation 197 process. CTFP (Deng et al., 2020) is a well-known normalizing flow model. It can treat both regular and irreugular time-series data by a deformation of the standard Wiener process. 199

TimeGAN (Yoon et al., 2019) uses a GAN architecture to generate time-series. First, it trains an 200 encoder and decoder, which transform a time-series sample $\mathbf{x}_{1:N}$ into latent vectors $\mathbf{h}_{1:N}$ and re-201 cover them by using a recurrent neural network (RNN). Next, it trains a generator and discriminator 202 pair on latent space, by minimizing the discrepancy between an estimated and true distribution, i.e. 203 $\hat{p}(\mathbf{x}_n|\mathbf{x}_{1:n-1})$ and $p(\mathbf{x}_n|\mathbf{x}_{1:n-1})$. Since it uses an RNN-based encoder, it can efficiently learn the 204 conditional likelihood $p(\mathbf{x}_n | \mathbf{x}_{1:n-1})$ by treating it as $p(\mathbf{h}_n | \mathbf{h}_{n-1})$, since $\mathbf{h}_n \sim \mathbf{x}_{1:n}$ under the regime 205 of RNNs. Therefore, it can generate each observation \mathbf{x}_n considering its previous history $\mathbf{x}_{1:n-1}$. 206 However, GAN-based generative models are vulnerable to the issue of mode collapse (Xiao et al., 207 2022) and unstable behavior problems during training (Chu et al., 2020). GT-GAN (Jeon et al., 208 2022) attempted to solve the problems by incorporating an invertible neural network-based generator into its framework. There also exist GAN-based methods to generate other types of sequential 209 data, e.g., video, sound, etc (Esteban et al., 2017; Mogren, 2016; Xu et al., 2020; Donahue et al., 210 2019). In our experiments, we also use them as our baselines for thorough evaluations. 211

212

- **PROPOSED METHOD** 3
- 213 214

- Our proposed TSGM consists of three networks: an encoder, a decoder, and a conditional score 215 network (cf. Fig. 2). Firstly, we train the encoder and the decoder to connect between time-series

samples and a latent space. Next, using the pre-trained encoder and decoder, we train the conditional score network on the latent space. The conditional score network will be used for sampling fake time-series on the latent space.

220 3.1 PROBLEM FORMULATION

Let \mathcal{X} and \mathcal{H} denote a data space and a latent space, respectively. We define $\mathbf{x}_{1:N}$ as a timeseries sample with a sequential length of N, and \mathbf{x}_n is a multi-dimensional observation of $\mathbf{x}_{1:N}$ at sequential order n. Similarly, $\mathbf{h}_{1:N}$ (resp. \mathbf{h}_n) denotes an embedded time series (resp. an embedded observation).

Each observation \mathbf{x}_n can be represented as a pair of time and features, i.e., $\mathbf{x}_n = (t_n, \mathbf{u}(t_n))$, where $t_n \in \mathbb{R}_{\geq 0}$ is a time stamp of feature $\mathbf{u}(t_n) \in \mathbb{R}^{\dim(\mathbf{u})}$, and $\dim(\mathbf{u})$ is a feature dimension. \mathcal{X} can be classified into two types: regular time-series and irregular time-series. For irregular setting, we randomly remove 30%, 50% and 70% of samples from regular time-series. Therefore, the only difference between these types is whether time intervals, $\{t_{n+1} - t_n\}_{n=1}^{N-1}$, are the same or not.

Our irregular setting considers generating complete time-series given training data with missing values. Consequently, we generate total time-series on both regular and irregular settings and compare the generated ones with the original regular data on evaluation process.

- 234 235
- 3.2 ENCODER AND DECODER

237 The encoder and decoder have the task of mapping time-series data to a latent space and vice versa. 238 We define e and d as an encoding function mapping \mathcal{X} to \mathcal{H} and a decoding function mapping \mathcal{H} to 239 \mathcal{X} , respectively. For simplicity but without loss of generality, we utilize an autoregressive autoencoder: RNN-based ones for regular time-series (Cho et al., 2014b) and Neural CDE, GRU-ODE for 240 irregular setting (Kidger et al., 2020; Brouwer et al., 2019). In this section, we describe the regular 241 time-series generation as representative one and leave the irregular case in Appendix I, depicting 242 how continuous-time methods can be used for the encoder and the decoder to better synthesize 243 irregular time-series. 244

The encoder e and the decoder d consist of recurrent neural networks, e.g., gated recurrent units (GRUs) (Cho et al., 2014b). Since we use RNNs, both e and d are defined recursively as follows:

247 248

$$\mathbf{h}_n = e(\mathbf{h}_{n-1}, \mathbf{x}_n), \qquad \hat{\mathbf{x}}_n = d(\mathbf{h}_n), \tag{5}$$

where $\hat{\mathbf{x}}_n$ denotes a reconstructed time-series sample at sequential order *n*. It is well-known that RNN was devised to efficiently handle variable sequences by summarizing past observations (Cho et al., 2014a). For example, Yoon et al. (2019) also used RNN-based encoder and decoder to provide a reversible mapping between features and latent representations, thereby reducing the highdimensionality of the adversarial learning space, which is well supported by $\mathbf{h}_n \sim \mathbf{x}_{1:n}$.

After embedding real time-series data onto a latent space, we can train the conditional score network with its conditional log-likelihood, whose architecture is described in Appendix J.2. The encoder and decoder are pre-trained before our main training.

256 257 258

254

255

3.3 TRAINING OBJECTIVE FUNCTION

Loss for autoencoder We use two training objective functions. First, we train the encoder and the decoder using L_{ed} . Let $\mathbf{x}_{1:N}^0 \sim p(\mathbf{x}_{1:N}^0)$ and $\hat{\mathbf{x}}_{1:N}^0$ denote an real time-series sample and its reconstructed copy by the encoder-decoder process, respectively. Then, L_{ed} denotes the following MSE loss between $\mathbf{x}_{1:N}^0$ and its reconstructed copy $\hat{\mathbf{x}}_{1:N}^0$:

$$L_{ed} = \mathbb{E}_{\mathbf{x}_{1:N}^{0}}[\|\hat{\mathbf{x}}_{1:N}^{0} - \mathbf{x}_{1:N}^{0}\|_{2}^{2}].$$
(6)

264 265

Loss for score network Next, we define another loss $L_{score}^{\mathcal{H}}$ in Eq. equation 11 to train the conditional score network M_{θ} , which is one of our main contributions. In order to derive the training loss $L_{score}^{\mathcal{H}}$ from the initial loss definition L_1 , we describe its step-by-step derivation procedure. At sequential order n in $\{1, ..., N\}$, we diffuse $\mathbf{x}_{1:n}^0$ through a sufficiently large number of steps of the forward SDE to a Gaussian distribution. Let $\mathbf{x}_{1:n}^s$ denotes a diffused sample at step $s \in [0, 1]$ from 270 $\mathbf{x}_{1:n}^0$. Then, the conditional score network $M_{\theta}(s, \mathbf{x}_{1:n}^s, \mathbf{x}_{1:n-1}^0)$ can be trained to learn the gradient of the conditional log-likelihood with the following L_1 loss:

$$L_1 = \mathbb{E}_s \mathbb{E}_{\mathbf{x}_{1:N}^0} \left[\sum_{n=1}^N \lambda(s) l_1(n,s) \right],\tag{7}$$

where

277

278 279

286

287

289 290 291

292 293 294

295

296

297

298

299

$$l_1(n,s) = \mathbb{E}_{\mathbf{x}_{1:n}^s} \left[\left\| M_{\theta}(s, \mathbf{x}_{1:n}^s, \mathbf{x}_{1:n-1}^0) - \nabla_{\mathbf{x}_{1:n}^s} \log p(\mathbf{x}_{1:n}^s | \mathbf{x}_{1:n-1}^0) \right\|_2^2 \right].$$
(8)

In the above definition, $\nabla_{\mathbf{x}_{1:n}^s} \log p(\mathbf{x}_{1:n}^s | \mathbf{x}_{1:n-1}^0)$, where \mathbf{x}_i^0 depends on $\mathbf{x}_{1:i-1}^0$ for each $i \in \{2, ..., n\}$, is designed specially for time-series generation. Note that for our training, $\mathbf{x}_{1:n}^s$ is sampled from $p(\mathbf{x}_{1:n}^s | \mathbf{x}_{1:n-1}^0)$, and s is uniformly sampled from [0, 1].

However, using the above formula, which is a naïve score matching on time-series, is computationally prohibitive (Hyvärinen, 2005; Song et al., 2019). Thanks to the following theorem, the more efficient denoising score loss L_{score} can be defined.

Theorem 3.1 (Autoregressive denoising score matching). $l_1(n, s)$ can be replaced with the following $l_2(n, s)$

$$L_{score} = \mathbb{E}_{s} \mathbb{E}_{\mathbf{X}_{1:N}^{0}} \left[\sum_{n=1}^{N} \lambda(s) l_{2}(n,s) \right],$$
(9)

(11)

where

$$l_{2}(n,s) = \mathbb{E}_{\mathbf{X}_{1:n}^{s}} \left[\left\| M_{\theta}(s, \mathbf{X}_{1:n}^{s}, \mathbf{X}_{1:n-1}^{0}) - \nabla_{\mathbf{X}_{1:n}^{s}} \log p(\mathbf{X}_{1:n}^{s} | \mathbf{X}_{1:n}^{0}) \right\|_{2}^{2} \right].$$
(10)

Then, $L_1 = L_{score}$ is satisfied.

Since we pre-train the encoder and decoder, the encoder can embed $\mathbf{x}_{1:n}^0$ into $\mathbf{h}_n^0 \in \mathcal{H}$. Ideally, \mathbf{h}_n^0 involves the entire information of $\mathbf{x}_{1:n}^0$. Therefore, L_{score} can be re-written as follows with the embeddings in the latent space:

 $L_{score}^{\mathcal{H}} = \mathbb{E}_{s} \mathbb{E}_{\mathbf{h}_{1:N}^{0}} \sum_{n=1}^{N} \left[\lambda(s) l_{3}(n,s) \right],$

300 301 302

303

304

306

307

308

310

with $l_3(n,s) = \mathbb{E}_{\mathbf{h}_n^s} \left[\left\| M_{\theta}(s, \mathbf{h}_n^s, \mathbf{h}_{n-1}^0) - \nabla_{\mathbf{h}_n^s} \log p(\mathbf{h}_n^s | \mathbf{h}_n^0) \right\|_2^2 \right]$. $L_{score}^{\mathcal{H}}$ is what we use for our experiments (instead of L_{score}). Until now, we introduced our target objective functions, L_{ed} and $L_{score}^{\mathcal{H}}$. We note that we use exactly the same weight $\lambda(s)$ as that in (Song et al., 2021). Related proofs are given in Appendix A.

3.4 TRAINING AND SAMPLING PROCEDURES

Training method We explain details of our training method. At first, we pre-train both the encoder and decoder using L_{ed} . After pre-training them, we train the conditional score network. When training the latter one, we use the embedded hidden vectors produced by the encoder. After encoding an input $\mathbf{x}_{1:N}^0$, we obtain its latent vectors $\mathbf{h}_{1:N}^0$ — we note that each hidden vector \mathbf{h}_n^0 has all the previous information from 1 to *n* for the RNN-based encoder's autoregressive property as shown in the Equation 3.2. We use the following forward process (Song et al., 2021), where *n* means the sequence order of the input time-series, and *s* denotes the time (or step) of the diffusion step :

$$d\mathbf{h}_n^s = \mathbf{f}(s, \mathbf{h}_n^s) ds + g(s) d\mathbf{w}, \qquad s \in [0, 1].$$

318 319 320

Note that we only use the VP and subVP-based TSGM in our experiments and exclude the VE-based one for its inferiority for time series synthesis in our experiments. During the forward process, the conditional score network reads the pair $(s, \mathbf{h}_n^s, \mathbf{h}_{n-1}^0)$ as input and thereby, it can learn the conditional score function $\nabla \log p(\mathbf{h}_n^s | \mathbf{h}_{n-1}^0)$ by using $L_{score}^{\mathcal{H}}$, where $\mathbf{h}_0^0 = \mathbf{0}$.



Figure 3: t-SNE plots for TSGM (1st and 2nd columns), TimeGAN (3rd columns), TimeVAE (4th columns), GT-GAN (5th columns) in Stocks and Air datasets. Red and blue dots mean original and synthesized samples, respectively. Refer to Appendix L for addition visualizations

Sampling method After the training procedure, we use the following conditional reverse process: $d\mathbf{h}_{n}^{s} = [\mathbf{f}(s, \mathbf{h}_{n}^{s}) - g^{2}(s)\nabla_{\mathbf{h}_{n}^{s}}\log p(\mathbf{h}_{n}^{s}|\mathbf{h}_{n-1}^{0})]ds + g(s)d\bar{\mathbf{w}}, \quad (12)$

where s is uniformly sampled over [0, 1]. The conditional score function in this process can be replaced with the trained score network $M_{\theta}(s, \mathbf{h}_{n}^{s}, \mathbf{h}_{n-1}^{0})$. The detailed sampling method is as follows:

- 1. At first, we sample \mathbf{z}_1 from a Gaussian prior distribution and set $\mathbf{h}_1^1 = \mathbf{z}_1$ and $\mathbf{h}_0^0 = \mathbf{0}$. We then generates an initial observation $\hat{\mathbf{h}}_1^0$ by denoising \mathbf{h}_1^1 following the conditional reverse process with $M_{\theta}(\mathbf{s}, \mathbf{h}_n^s, \mathbf{h}_0^0)$ via the *predictor-corrector* method (Song et al., 2021).
- 2. We repeat the following computation for every $2 \le n \le N$, i.e., recursive generation. We sample \mathbf{z}_n from a Gaussian prior distribution and set $\mathbf{h}_n^1 = \mathbf{z}_n$ for $n \in \{2, ..., N\}$. After reading the previously generated samples $\hat{\mathbf{h}}_{n-1}^0$, we then denoise \mathbf{h}_n^1 following the conditional reverse process with $M_{\theta}(s, \mathbf{h}_n^s, \mathbf{h}_{n-1}^0)$ to generate $\hat{\mathbf{h}}_n^0$ via the *predictor-corrector* method.

Once the sampling procedure is finished, we can reconstruct $\hat{\mathbf{x}}_{1:N}^0$ from $\hat{\mathbf{h}}_{1:N}^0$ using the trained decoder at once.

- 357 4 EXPERIMENTS
 - 4.1 EXPERIMENTAL ENVIRONMENTS
 - 4.1.1 BASELINES AND DATASETS

In the case of the regular time-series generation, we use 4 real-world datasets from various fields
with 9 baselines. For the irregular time-series generation, we randomly remove some observations
from each time-series sample with 30%, 50%, and 70% missing rates. Therefore, we totally treat 16
datasets, i.e., 4 datasets with one regular and three irregular settings, and 9 baselines.

Our collection of baselines covers almost all existing types of time-series synthesis methods, rang-ing from autoregressive generative models to VAEs and GANs. For the baselines, we reuse their released source codes in their official repositories and rely on their designed training and model selection procedures. If a baseline does not support irregular time-series synthesis, we replace its RNN encoder with GRU-D (Che et al., 2016) modified from GRUs to deal with irregular time-series. For those that do not use an RNN-based encoder, we add GRU-D in front of the encoder, such as TimeVAE and COT-GAN. Therefore, all baselines are tested for the regular and irregular environments. We refer to Appendix E for the detailed descriptions on our datasets, baselines, and Appendix G for other software/hardware environments.

- 375 4.1.2 EVALUATION METRICS376
- In the image generation domain, researchers have evaluated the *fidelity* and the *diversity* of models by using the Fréchet inception distance (FID) and inception score (IS). On the other hand, to measure

Table 2: The left and right ones denote experimental results on regular time-series and irregular
time-series with 30% missing rates, respectively. Results for higher missing rates are in Table 12.
Note that except for our method (TSGM), CTFP, and GT-GAN, the other methods cannot deal with
irregular time series, so we make it possible for them to operate on irregular settings by replacing
RNN encoder with GRU-D.

3	1	Mathad		Regular	Settings		Irregular Settings (Missing Rate: 30%)			
4		Method	Stocks	Energy	Ăir	AI4I	Stocks	Energy	Ăir	AI4I
-		TSGM-VP	$.022 \pm .005$.221±.025	$.122 {\pm} .014$	$.147 {\pm} .005$	$.062 \pm .018$.294±.007	$.190 {\pm} .042$	$.142 {\pm} .048$
C		TSGM-subVP	$.021 {\pm} .008$	$.198 {\pm} .025$	$.127 {\pm} .010$	$.150 {\pm} .010$.025±.009	$.326 {\pm} .008$	$.240 {\pm} .018$	$.121 {\pm} .082$
6		T-Forcing	.226±.035	$.483 {\pm} .004$	$.404 \pm .020$	$.435 \pm .025$.409±.051	$.347 \pm .046$.458±.122	.493±.018
7	•	P-Forcing	.257±.026	$.412 \pm .006$	$.484 {\pm} .007$	$.443 \pm .026$	$.480 \pm .060$	$.491 \pm .020$	$.494 \pm .012$	$.430 \pm .061$
ſ	ore	TimeGAN	.102±.031	$.236 \pm .012$	$.447 \pm .017$.070±.009	.411±.040	$.479 \pm .010$	$.500 \pm .001$	$.500 \pm .000$
8	sc	RCGAN	.196±.027	$.336 \pm .017$	$.459 \pm .104$	$.234 \pm .015$	$.500 \pm .000$	$.500 {\pm} .000$	$.500 \pm .000$	$.500 \pm .000$
2	sc.	C-RNN-GAN	.399±.028	$.499 {\pm} .001$	$.499 {\pm} .000$	$.499 {\pm} .001$	$.500 \pm .000$	$.500 {\pm} .000$	$.500 \pm .000$	$.450 \pm .150$
9	Di	TimeVAE	.175±.031	$.498 {\pm} .006$.381±.037	$.446 \pm .024$.423±.088	$.382 \pm .124$.373±.191	$.384 {\pm} .086$
0		COT-GAN	$.285 \pm .030$	$.498 {\pm} .000$	$.423 \pm .001$	$.411 \pm .018$.499±.001	$.500 \pm .000$	$.500 \pm .000$	$.500 \pm .000$
		CTFP	$.499 \pm .000$	$.500 \pm .000$	$.499 \pm .000$	$.499 \pm .001$	$.500 \pm .000$	$.500 \pm .000$	$.500 \pm .000$	$.499 {\pm} .001$
1		GT-GAN	.077±.031	$.221 \pm .068$	$.413 {\pm} .001$	$.394 \pm .090$.251±.097	$.333 {\pm} .063$	$.454 \pm .029$	$.435 {\pm} .018$
2		TSGM-VP	.037±.000	$.257 \pm .000$	$.005 {\pm} .000$	$.217 {\pm} .000$	$.012 {\pm} .002$	$.049 {\pm} .001$	$.042 \pm .002$	$.067 \pm .013$
-		TSGM-subVP	$.037 {\pm} .000$	$.252 {\pm} .000$	$.005 {\pm} .000$	$.217 {\pm} .000$	$.012 \pm .001$	$.049 {\pm} .001$	$.044 \pm .004$	$.061 {\pm} .001$
3		T-Forcing	$.038 \pm .001$	$.315 \pm .005$	$.008 \pm .000$	$.242 \pm .001$	$.027 \pm .002$	$.090 \pm .001$	$.112 \pm .004$	$.147 \pm .010$
4		P-Forcing	.043±.001	$.303 \pm .006$	$.021 \pm .000$	$.220 \pm .000$	$.079 \pm .008$	$.147 \pm .001$	$.101 \pm .003$	$.134 \pm .005$
-	re	TimeGAN	.038±.001	$.273 \pm .004$	$.017 \pm .004$	$.253 \pm .002$.105±.053	$.248 \pm .024$	$.325 \pm .005$	$.251 \pm .010$
C	sco	RCGAN	$.040 \pm .001$	$.292 \pm .005$	$.043 \pm .000$	$.224 \pm .001$	$.523 \pm .020$	$.409 \pm .020$	$.342 \pm .018$	$.329 \pm .037$
6	J. 5	C-RNN-GAN	$.038 \pm .000$	$.483 \pm .005$	$.111 \pm .000$	$.340 \pm .006$	$.345 \pm .002$	$.440 \pm .000$	$.354 \pm .060$	$.400 \pm .026$
_	ree	TimeVAE	$.042 \pm .002$	$.268 \pm .004$	$.013 \pm .002$	$.233 \pm .010$.207±.014	$.139 \pm .004$	$.105 \pm .002$	$.144 \pm .003$
7	H	COT-GAN	$.044 \pm .000$	$.260 \pm .000$	$.024 \pm .001$	$.220 \pm .000$.274±.000	$.427 \pm .000$	$.451 \pm .000$	$.570 \pm .000$
R		CTFP	$.084 \pm .005$	$.469 \pm .008$	$.476 \pm .235$	$.412 \pm .024$	$.070 \pm .009$	$.499 {\pm} .000$	$.060 \pm .027$	$.424 \pm .002$
		GT-GAN	$.040 \pm .000$	$.312 \pm .002$	$.007 \pm .000$	$.239 \pm .000$.077±.031	$.221 \pm .068$.064±.002	.087±.013
9		Original	$.036 \pm .001$	$.250 \pm .003$	$.004 \pm .000$	$.217 \pm .000$	$.011 \pm .002$	$.045 \pm .001$	$.044 \pm .006$	$.059 \pm .001$

the fidelity and the diversity of synthesized time-series samples, we use the following predictive score and the discriminative score as in (Yoon et al., 2019; Jeon et al., 2022). We strictly follow the evaluation protocol agreed by the time-series research community (Yoon et al., 2019; Jeon et al., 2022). Both metrics are designed in a way that lower values are preferred. We run each generative method 10 times with different seeds, and report its mean and standard deviation of the following discriminative and predictive scores:

i) *Predictive Score*: We use the predictive score to evaluate whether a generative model can successfully reproduce the temporal properties of the original data. To do this, we first train a popular LSTM-based sequence model for time-series forecasting with synthesized samples. However, the existing predictive score only predicts a last sample of total sequence, so Jeon et al. (2022) suggests a comprehensive approach that considers the entire time series. For fair comparison, we reuse evaluation metrics of Yoon et al. (2019), Jeon et al. (2022) for regular and irregular settings, respectively.

412
413 ii) *Discriminative Score*: In order to assess how similar the original and generated samples are, we train a 2-layer LSTM model that classifies the real/fake samples into two classes, real or fake. We use the performance of the trained classifier on the test data as the discriminative score. Therefore, lower discriminator scores mean real and fake samples are similar.

- 417 4.2 EXPERIMENTAL RESULTS
- 418

38

38

39 39 39

At first, on the regular time-series generation, Table 2 shows that our method achieves remarkable results, outperforming TimeGAN and GT-GAN except only for the discriminative score on AI4I. Especially, for Stock, Energy, and Air, TSGM exhibits overwhelming performance by large margins for the discriminative score. Moreover, for the predictive score, TSGM performs the best and obtains almost the same scores as that of the original data, which indicates that generated samples from TSGM preserve all the predictive characteristics of the original data.

Next, on the irregular time-series generation, we give the result with the 30% missing rate setting on
 Table 2 and other results in Appendix K. TSGM also defeats almost all baselines by large margins on
 both the discriminative and predictive scores. Interestingly, VP generates poorer data as the missing
 rate grows up, while subVP synthesizes better one.

We show t-SNE visualizations and KDE plots for the regular time-series generation in Figure 3 and
 Figure 1. TimeGAN, GT-GAN, and TimeVAE are representative GAN or VAE-based baselines. In
 the figures, unlike the baseline methods, the synthetic samples generated from TSGM consistently
 show successful recall from the original data.

Table 3: Sensitivity results on the depth of M_{θ} and the number of sampling steps. Our default TSGM has a depth of 4 and its number of sampling steps is 1,000. For other omitted datasets, we observe similar patterns.

Method	TSGM		Depth of 3		500 steps		250 steps		100 steps		
SDE	VP	subVP	VP	subVP	VP	subVP	VP	subVP	VP	subVP	
ý Stocks	.022±.005	.021±.008	.022±.004	$.020 \pm .007$.025±.005	.020±.004	.067±.009	.022±.009	.202±.013	.023±.005	
Energy	.221±.025	.198±.025	.175±.009	.182±.009	.259±.003	.248±.002	.250±.003	.247±.002	$.325 \pm .003$.237±.004	
ਦ੍ਹਾਂ Stocks	.037±.000	.037±.000	.037±.000	.037±.000	.037±.000	.037±.000	.037±.000	.037±.000	$.039 \pm .000$.037±.000	
Energy	$.257 \pm .000$.252±.000	.253±.000	$.253 \pm .000$.257±.000	.253±.000	.256±.000	.253±.000	$.256 \pm .000$.253±.000	

Furthermore, TSGM generates more diverse synthetic samples compared to the three representative baselines across all cases. Notably, TSGM achieves significantly higher diversity on the Energy and Air dataset, which exhibits the most complex correlations (cf. Fig. 3 and Fig. 6). Diffusion models are known to produce more diverse data than GANs (Bayat, 2023; Yang et al., 2023). As demonstrated by our results, TSGM synthesizes diverse data, highlighting another advantage of using diffusion models as anticipated.

4.3 SENSITIVITY AND ABLATION STUDIES

We conduct two sensitivity studies on regular time-series: i) reducing the depth of our score network,
ii) decreasing the sampling step numbers. The results are in Table 3. At first, we modify the depth of our score network from 4 to 3 to check the performance of the lighter conditional score network.
Surprisingly, we achieve a better discriminative score with a slight loss on the predictive score. Next, we decrease the number of sampling steps for faster sampling from 1,000 steps to 500, 250, and 100 steps, respectively. For VP, the case of 500 steps achieves almost the same results as that of original TSGM. Surprisingly, in the case of subVP, we achieve good results until 100 steps.

As an ablation study, we simultaneously train the conditional score network, encoder, and de-coder from scratch on regular time-series gen-eration (i.e., without the pre-training process). The results are in Table 4. These ablation mod-els are worse than the full model due to the in-creased training complexity, but they still out-perform many baselines. This ablation study shows the efficacy of pre-training our autoen-coder. Additionally, in Appendix M, we pro-

Table 4: Comparison between with and withoutpre-training the autoencoder

Disc.	Method	SDE	Stocks	Energy
	TSCM	VP	$.022 \pm .005$.221±.025
	13010	subVP	$.021 {\pm} .008$	$.198 {\pm} .025$
	w/o pro training	VP	$.022 \pm .004$	$.322 \pm .003$
	w/o pre-training	subVP	$.059 \pm .006$	$.284 \pm .004$
	TECM	VP	$.037 \pm .000$.257±.000
Pred.	ISGM	subVP	$.037 {\pm} .000$	$.252 \pm .000$
	who are training	VP	$.037 \pm .000$	$.252 \pm .000$
	w/o pre-training	subVP	$.037 \pm .000$	$.251 {\pm} .000$

vide an additional ablation study about the efficacy of our recursive structures.

5 CONCLUSIONS

We presented a score-based generative model framework for universal time-series generation. We combined an autoencoder and our score network into a single framework to accomplish the goal — our framework supports RNN-based or continuous-time method-based autoencoders. We also designed an appropriate denoising score matching loss for our generation task and achieved state-of-the-art results on various datasets in terms of the discriminative and predictive scores. In addition, we conducted rigorous ablation and sensitivity studies to prove the efficacy of our model design.

Limitations. Although our method achieves state-of-the-art sampling quality and diversity, there exists a fundamental problem that all SGMs have. That is, SGMs are slower than GANs for generating samples (see Appendix H). Since there are several accomplishments for faster sampling (Xiao et al., 2022; Jolicoeur-Martineau et al., 2021), however, one can apply them to our method and it would be much faster without any loss of sampling quality and diversity. Moreover, TSGM's generation is not too much slow since time-series data are much smaller than image data, which means TSGM still has its commercial value. We point out that TSGM generates a sample of the Energy dataset, the largest dataset among our baselines, using only 0.8 second (see Section H).

486 REFERENCES 487

501

507

517

527

528

- Nesreen K. Ahmed, Amir F. Atiya, Neamat El Gayar, and Hisham El-Shishiny. An empirical com-488 parison of machine learning models for time series forecasting. *Econometric Reviews*, 29(5-6): 489 594-621, 2010. 490
- 491 Brian D.O. Anderson. Reverse-time diffusion equation models. Stochastic Processes and their 492 Applications, 12(3):313-326, 1982. 493
- Reza Bayat. A study on sample diversity in generative models: GANs vs. diffusion models, 2023. 494 URL https://openreview.net/forum?id=BQpCuJoMykZ. 495
- 496 Edward De Brouwer, Jaak Simm, Adam Arany, and Yves Moreau. Gru-ode-bayes: Continuous 497 modeling of sporadically-observed time series. CoRR, abs/1905.12374, 2019. URL http: 498 //arxiv.org/abs/1905.12374. 499
- 500 Luis M. Candanedo, Véronique Feldheim, and Dominique Deramaix. Data driven prediction models of energy use of appliances in a low-energy house. Energy and Buildings, 140:81–97, 2017. ISSN 0378-7788. 502
- Zhengping Che, S. Purushotham, Kyunghyun Cho, David A. Sontag, and Yan Liu. Recurrent neural 504 networks for multivariate time series with missing values. Scientific Reports, 8, 2016. URL 505 https://api.semanticscholar.org/CorpusID:4900015. 506
- Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary dif-508 ferential equations, 2018. URL https://arxiv.org/abs/1806.07366.
- 509 KyungHyun Cho, Bart van Merrienboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties 510 of neural machine translation: Encoder-decoder approaches. CoRR, abs/1409.1259, 2014a. URL 511 http://arxiv.org/abs/1409.1259. 512
- 513 Kyunghyun Cho, Bart van Merrienboer, Caglar Gülcehre, Fethi Bougares, Holger Schwenk, and 514 Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical ma-515 chine translation. CoRR, abs/1406.1078, 2014b. URL http://arxiv.org/abs/1406. 516 1078.
- Casey Chu, Kentaro Minami, and Kenji Fukumizu. Smoothness and stability in GANs. In Interna-518 tional Conference on Learning Representations, 2020. 519
- Saloni Dash, Andrew Yale, Isabelle Guyon, and Kristin P. Bennett. Medical time-series data gen-521 eration using generative adversarial networks. In Artificial Intelligence in Medicine: 18th In-522 ternational Conference on Artificial Intelligence in Medicine, AIME 2020, Minneapolis, MN, USA, August 25-28, 2020, Proceedings, page 382-391, Berlin, Heidelberg, 2020. Springer-523 Verlag. ISBN 978-3-030-59136-6. doi: 10.1007/978-3-030-59137-3_34. URL https: 524 //doi.org/10.1007/978-3-030-59137-3_34. 525
 - S. De Vito, E. Massera, M. Piga, L. Martinotto, and G. Di Francia. On field calibration of an electronic nose for benzene estimation in an urban pollution monitoring scenario. Sensors and Actuators B: Chemical, 129(2):750-757, 2008. ISSN 0925-4005.
- Ruizhi Deng, Bo Chang, Marcus A Brubaker, Greg Mori, and Andreas Lehrmann. Mod-530 eling continuous stochastic processes with dynamic normalizing flows. In H. Larochelle, 531 M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, Advances in Neural In-532 formation Processing Systems, volume 33, pages 7805-7815. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/ 534 file/58c54802a9fb9526cd0923353a34a7ae-Paper.pdf. 535
- Abhyuday Desai, Cynthia Freeman, Zuhui Wang, and Ian Beaver. Timevae: A variational autoencoder for multivariate time series generation. arXiv:2111.08095, 2021. 538
- Chris Donahue, Julian McAuley, and Miller Puckette. Adversarial audio synthesis. In International Conference on Learning Representations, 2019.

551

554

566

567

568

569

580

581

582

583 584

585 586

- 540 Cristóbal Esteban, Stephanie L. Hyland, and Gunnar Rätsch. Real-valued (medical) time series 541 generation with recurrent conditional GANs. arXiv:1706.02633, 2017. 542
- Tak-chung Fu. A review on time series data mining. Engineering Applications of Artificial Intelli-543 gence, 24(1):164-181, 2011. 544
- A. L. Goldberger, L. A. N. Amaral, L. Glass, J. M. Hausdorff, P. Ch. Ivanov, R. G. 546 Mark, J. E. Mietus, G. B. Moody, C.-K. Peng, and H. E. Stanley. PhysioBank, 547 PhysioToolkit, and PhysioNet: Components of a new research resource for complex 548 Circulation, 101(23):e215-e220, 2000 (June 13). physiologic signals. Circulation 549 Electronic Pages: http://circ.ahajournals.org/content/101/23/e215.full PMID:1085218; doi: 550 10.1161/01.CIR.101.23.e215.
- Anirudh Goyal, Alex Lamb, Ying Zhang, Saizheng Zhang, Aaron Courville, and Yoshua Bengio. 552 Professor forcing: A new algorithm for training recurrent networks. In Advances in Neural Infor-553 mation Processing Systems, 2016.
- 555 Alex Graves. Generating sequences with recurrent neural networks. arXiv:1308.0850, 2013. 556
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Advances 558 in Neural Information Processing Systems, 2020.
- 559 Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. Journal of 560 Machine Learning Research, 6(24):695-709, 2005. URL http://jmlr.org/papers/v6/ 561 hyvarinen05a.html. 562
- 563 Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Deep learning for time series classification: a review. Data Mining and Knowledge 564 Discovery, 33:917–963, 2019. 565
 - Jinsung Jeon, Jeonghak Kim, Haryong Song, Seunghyeon Cho, and Noseong Park. Gt-gan: General purpose time series synthesis with generative adversarial networks, 2022. URL https: //arxiv.org/abs/2210.02040.
- 570 Alexia Jolicoeur-Martineau, Ke Li, Rémi Piché-Taillefer, Tal Kachman, and Ioannis Mitliagkas. 571 Gotta go fast when generating data with score-based models. CoRR, abs/2105.14080, 2021. URL https://arxiv.org/abs/2105.14080. 572
- 573 Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. An-574 alyzing and improving the image quality of stylegan. CoRR, abs/1912.04958, 2019. URL 575 http://arxiv.org/abs/1912.04958. 576
- 577 Patrick Kidger, James Morrill, James Foster, and Terry J. Lyons. Neural controlled differential 578 equations for irregular time series. CoRR, abs/2005.08926, 2020. URL https://arxiv. org/abs/2005.08926. 579
 - Jayoung Kim, Chaejeong Lee, Yehjin Shin, Sewon Park, Minjung Kim, Noseong Park, and Jihoon Cho. Sos: Score-based oversampling for tabular data. In Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, 2022.
 - Stephan Matzka. Ai4i 2020 predictive maintenance dataset. https://archive.ics.uci. edu/ml/datasets/AI4I+2020+Predictive+Maintenance+Dataset, 2020.
- Chenlin Meng, Lantao Yu, Yang Song, Jiaming Song, and Stefano Ermon. Autoregressive score 587 matching. CoRR, abs/2010.12810, 2020. URL https://arxiv.org/abs/2010.12810. 588
- 589 Olof Mogren. C-RNN-GAN: A continuous recurrent neural network with adversarial training. In 590 Constructive Machine Learning Workshop (CML) at NeurIPS 2016, 2016.
- Kashif Rasul, Calvin Seward, Ingmar Schuster, and Roland Vollgraf. Autoregressive denoising dif-592 fusion models for multivariate probabilistic time series forecasting. In International Conference on Machine Learning, 2021.

594 595 596	Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedi- cal image segmentation. In <i>Medical Image Computing and Computer-Assisted Intervention (MIC-CAI)</i> , pages 234–241, 2015.
598 599	Donald B. Rubin. Inference and missing data. <i>Biometrika</i> , 63(3):581–592, 1976. ISSN 00063444. URL http://www.jstor.org/stable/2335739.
600 601 602	Joseph L. Schafer and John W. Graham. Missing data: our view of the state of the art. <i>Psychological methods</i> , 7 2:147–77, 2002. URL https://api.semanticscholar.org/CorpusID: 7745507.
603 604 605 606	Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In <i>International Conference on Machine Learning</i> , 2015.
607 608	Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In Advances in Neural Information Processing Systems, 2019.
609 610 611 612	Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach to density and score estimation. <i>CoRR</i> , abs/1905.07088, 2019. URL http://arxiv.org/abs/1905.07088.
613 614 615	Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In <i>International Conference on Learning Representations</i> , 2021.
616 617 618 619	Ilya Sutskever, James Martens, and Geoffrey Hinton. Generating text with recurrent neural networks. In <i>Proceedings of the 28th International Conference on International Conference on Machine Learning</i> , page 1017–1024, 2011.
620 621	Simo Särkkä and Arno Solin, editors. <i>Applied stochastic differential equations</i> , volume 10. Cambridge University Press, 2019.
622 623 624	Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. Csdi: Conditional score-based diffusion models for probabilistic time series imputation. In <i>Advances in Neural Information Processing Systems</i> , 2021.
625 626 627	Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. Journal of machine learning research, 9(86):2579–2605, 2008.
628 629	Pascal Vincent. A connection between score matching and denoising autoencoders. <i>Neural Computation</i> , 23(7):1661–1674, 2011.
631 632 633	Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma with denoising diffusion GANs. In <i>International Conference on Learning Representations</i> , 2022. URL https://openreview.net/forum?id=JprM0p-q0Co.
634 635	Tianlin Xu, Li K. Wenliang, Michael Munn, and Beatrice Acciaio. Cot-gan: Generating sequential data via causal optimal transport. In <i>Advances in Neural Information Processing Systems</i> , 2020.
637 638 639	Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications, 2023.
640 641 642	Jinsung Yoon, Daniel Jarrett, and Mihaela van der Schaar. Time-series generative adversarial net- works. In Advances in Neural Information Processing Systems, 2019.
643 644 645	
646	

A PROOFS

 We introduce an additional lemma to prove Theorem 3.1. In the following lemma, we state about the denoising score matching on time-series.

Lemma A.1. In L_1 loss function, $l_1(n, s)$ can be replaced by the following $l_2^{\star}(n, s)$:

$$l_{2}^{\star}(n,s) = \mathbb{E}_{\mathbf{X}_{n}^{0}} \mathbb{E}_{\mathbf{X}_{1:n}^{s}} \left[\left\| M_{\theta}(s,\mathbf{X}_{1:n}^{s},\mathbf{X}_{1:n-1}^{0}) - \nabla_{\mathbf{X}_{1:n}^{s}} \log p(\mathbf{X}_{1:n}^{s}|\mathbf{X}_{1:n}^{0}) \right\|_{2}^{2} \right],$$
(13)

where \mathbf{X}_{n}^{0} and $\mathbf{X}_{1:n}^{s}$ are sampled from $p(\mathbf{X}_{n}^{0}|\mathbf{X}_{1:n-1}^{0})$ and $p(\mathbf{X}_{1:n}^{s}|\mathbf{X}_{1:n}^{0})$. Therefore, we can use an alternative objective, $L_{2} = \mathbb{E}_{s}\mathbb{E}_{\mathbf{X}_{1:N}}\left[\sum_{n=1}^{N}\lambda(s)l_{2}^{\star}(n,s)\right]$ instead of L_{1} .

Proof. At first, if n = 1, it can be substituted with the naive denoising score loss by Vincent (2011) since $\mathbf{x}_0^0 = \mathbf{0}$.

Next, let us consider n > 1. $l_1(n, s)$ can be decomposed as follows:

$$l_{1}(n,s) = -2 \cdot \mathbb{E}_{\mathbf{x}_{1:n}^{s}} \langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \nabla_{\mathbf{x}_{1:n}^{s}} \log p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n-1}^{0}) \rangle + \mathbb{E}_{\mathbf{x}_{1:n}^{s}} \left[\left\| M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}) \right\|_{2}^{2} \right] + C_{1}$$
(14)

Here, C_1 is a constant that does not depend on the parameter θ , and $\langle \cdot, \cdot \rangle$ means the inner product. Then, the first part's expectation of the right-hand side can be expressed as follows:

$$\begin{split} \mathbb{E}_{\mathbf{x}_{1:n}^{s}} [\langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \nabla_{\mathbf{x}_{1:n}^{s}} \log p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n-1}^{0}) \rangle] \\ = \int_{\mathbf{x}_{1:n}^{s}} \langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \nabla_{\mathbf{x}_{1:n}^{s}} \log p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n-1}^{0}) \rangle p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n-1}^{0}) d\mathbf{x}_{1:n}^{s} \\ = \int_{\mathbf{x}_{1:n}^{s}} \langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \frac{1}{p(\mathbf{x}_{1:n-1}^{0})} \frac{\partial p(\mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0})}{\partial \mathbf{x}_{1:n}^{s}} \rangle d\mathbf{x}_{1:n}^{s} \\ = \int_{\mathbf{x}_{n}^{0}} \int_{\mathbf{x}_{1:n}^{s}} \langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \frac{1}{p(\mathbf{x}_{1:n-1}^{0})} \frac{\partial p(\mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}, \mathbf{x}_{n}^{0})}{\partial \mathbf{x}_{1:n}^{s}} \rangle d\mathbf{x}_{1:n}^{s} \\ = \int_{\mathbf{x}_{n}^{0}} \int_{\mathbf{x}_{1:n}^{s}} \langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \frac{\partial p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0})}{\partial \mathbf{x}_{1:n}^{s}} \rangle \frac{\partial p(\mathbf{x}_{1:n-1}^{s}, \mathbf{x}_{n}^{0})}{\partial \mathbf{x}_{1:n}^{s}} \rangle d\mathbf{x}_{1:n}^{s} d\mathbf{x}_{n}^{0} \\ = \int_{\mathbf{x}_{n}^{0}} \int_{\mathbf{x}_{1:n}^{s}} \langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \frac{\partial p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0})}{\partial \mathbf{x}_{1:n}^{s}} \rangle p(\mathbf{x}_{1:n-1}^{0}) d\mathbf{x}_{1:n-1}^{s} d\mathbf{x}_{n}^{0} \\ = \int_{\mathbf{x}_{n}^{0}} \int_{\mathbf{x}_{1:n}^{s}} \langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \frac{\partial p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0})}{\partial \mathbf{x}_{1:n}^{s}} \rangle p(\mathbf{x}_{n}^{0} | \mathbf{x}_{1:n-1}^{0}) d\mathbf{x}_{1:n}^{s} d\mathbf{x}_{n}^{0} \\ = \mathbb{E}_{\mathbf{x}_{n}^{0}} \left[\int_{\mathbf{x}_{1:n}^{s}} \langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \frac{\partial p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0})}{\partial \mathbf{x}_{1:n}^{s}} \rangle d\mathbf{x}_{1:n}^{s}} \right] \\ = \mathbb{E}_{\mathbf{x}_{n}^{0}} \mathbb{E}_{\mathbf{x}_{1:n}^{s}} [\langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \nabla_{\mathbf{x}_{1:n}^{s}} \log p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0}) \rangle p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0}) \rangle d\mathbf{x}_{1:n}^{s}} \right] \\ = \mathbb{E}_{\mathbf{x}_{n}^{0}} \mathbb{E}_{\mathbf{x}_{1:n}^{s}} [\langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \nabla_{\mathbf{x}_{1:n}^{s}} \log p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0}) \rangle] \end{cases}$$

Similarly, the second part's expectation of the right-hand side can be rewritten as follows:

 $\mathbb{E}_{\mathbf{x}_{1:n}^{s}}[\|M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0})\|_{2}^{2}]$ $= \int_{\mathbf{x}_{1:n}^{s}} \|M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0})\|_{2}^{2} \cdot \mathbf{p}(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n-1}^{0}) d\mathbf{x}_{1:n}^{s}$ $= \int_{\mathbf{x}_{n}^{0}} \int_{\mathbf{x}_{1:n}^{s}} \|M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0})\|_{2}^{2} \cdot \frac{\mathbf{p}(\mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}, \mathbf{x}_{n}^{0})}{\mathbf{p}(\mathbf{x}_{1:n-1}^{0})} d\mathbf{x}_{1:n}^{s} d\mathbf{x}_{n}^{0}$ $= \int_{\mathbf{x}_{n}^{0}} \int_{\mathbf{x}_{1:n}^{s}} \|M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0})\|_{2}^{2} \cdot \mathbf{p}(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0}) \frac{\mathbf{p}(\mathbf{x}_{1:n-1}^{0}, \mathbf{x}_{n}^{0})}{\mathbf{p}(\mathbf{x}_{1:n-1}^{0})} d\mathbf{x}_{1:n}^{s} d\mathbf{x}_{n}^{0}$ $= \int_{\mathbf{x}_{n}^{0}} \int_{\mathbf{x}_{1:n}^{s}} \|M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0})\|_{2}^{2} \cdot \mathbf{p}(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0}) \mathbf{p}(\mathbf{x}_{n}^{0} | \mathbf{x}_{1:n-1}^{0}) d\mathbf{x}_{1:n}^{s} d\mathbf{x}_{n}^{0}$ $= \int_{\mathbf{x}_{n}^{0}} \int_{\mathbf{x}_{1:n}^{s}} \|M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0})\|_{2}^{2} \cdot \mathbf{p}(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0}) \mathbf{p}(\mathbf{x}_{n}^{0} | \mathbf{x}_{1:n-1}^{0}) d\mathbf{x}_{1:n}^{s} d\mathbf{x}_{n}^{0}$ $= \mathbb{E}_{\mathbf{x}_{n}^{0}} \mathbb{E}_{\mathbf{x}_{1:n}^{s}} [\|M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0})\|_{2}^{2}]$ (16)

Finally, by using above results, we can derive following result:

$$l_{1} = \mathbb{E}_{\mathbf{x}_{n}^{0}} \mathbb{E}_{\mathbf{x}_{1:n}^{s}} \left[\left\| M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}) \right\|_{2}^{2} \right] + C_{1} - 2 \cdot \mathbb{E}_{\mathbf{x}_{n}^{0}} \mathbb{E}_{\mathbf{x}_{1:n}^{s}} \langle M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}), \nabla_{\mathbf{x}_{1:n}^{s}} \log p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0}) \rangle = \mathbb{E}_{\mathbf{x}_{n}^{0}} \mathbb{E}_{\mathbf{x}_{1:n}^{s}} \left[\left\| M_{\theta}(s, \mathbf{x}_{1:n}^{s}, \mathbf{x}_{1:n-1}^{0}) - \nabla_{\mathbf{x}_{1:n}^{s}} \log p(\mathbf{x}_{1:n}^{s} | \mathbf{x}_{1:n}^{0}) \right\|_{2}^{2} \right] + C$$
(17)

C is a constant that does not depend on the parameter θ .

Theorem A.1 (Autoregressive denoising score matching). $l_1(n,s)$ can be replaced with the following $l_2(n,s)$

$$L_{score} = \mathbb{E}_{s} \mathbb{E}_{\mathbf{x}_{1:N}^{0}} \left[\sum_{n=1}^{N} \lambda(s) l_{2}(n,s) \right],$$
(18)

where

$$l_2(n,s) = \mathbb{E}_{\mathbf{X}_{1:n}^s} \left[\left\| M_{\theta}(s, \mathbf{X}_{1:n}^s, \mathbf{X}_{1:n-1}^0) - \nabla_{\mathbf{X}_{1:n}^s} \log p(\mathbf{X}_{1:n}^s | \mathbf{X}_{1:n}^0) \right\|_2^2 \right].$$
(19)

746 Then, $L_1 = L_{score}$ is satisfied.

747 proof. By Lemma A.1, it suffices to show that $L_2 = L_{score}$. Whereas one can use the law 748 of total expectation, which means E[X] = E[E[X|Y]] if X,Y are on an identical probability 749 space to show the above formula, we calculate directly. At first, let us simplify the expecta-750 tion of the inner part with a symbol $f(\mathbf{x}_{1:n}^0)$ for our computational convenience, i.e., $f(\mathbf{x}_{1:n}^0) =$ 751 $\mathbb{E}_s \mathbb{E}_{\mathbf{x}_{1:n}^s} \left[\lambda(s) \| M_{\theta}(s, \mathbf{x}_{1:n}^s, \mathbf{x}_{1:n-1}^0) - \nabla_{\mathbf{x}_{1:n}^s} \log p(\mathbf{x}_{1:n}^s |\mathbf{x}_{1:n}^0) \|_2^2 \right]$. Then we have the following defi-753 nition:

$$L_{2} = \mathbb{E}_{s} \mathbb{E}_{\mathbf{x}_{1:N}^{0}} \left[l_{2}^{\star} \right] = \mathbb{E}_{\mathbf{x}_{1:N}^{0}} \left[\sum_{n=1}^{N} \mathbb{E}_{\mathbf{x}_{n}^{0}} \left[f(\mathbf{x}_{1:n}^{0}) \right] \right] = \sum_{n=1}^{N} \mathbb{E}_{\mathbf{x}_{1:N}^{0}} \mathbb{E}_{\mathbf{x}_{n}^{0}} \left[f(\mathbf{x}_{1:n}^{0}) \right]$$
(20)

At last, the expectation part can be further simplified as follows:

$$\mathbb{E}_{\mathbf{x}_{1:N}^{0}} \mathbb{E}_{\mathbf{x}_{n}^{0}} [f(\mathbf{x}_{1:n}^{0})]$$

$$= \int_{\mathbf{x}_{1:N}^{0}} \int_{\mathbf{x}_{n}^{0}} f(\mathbf{x}_{1:n}^{0}) p(\mathbf{x}_{n}^{0} | \mathbf{x}_{1:n-1}^{0}) d\mathbf{x}_{n}^{0} \cdot p(\mathbf{x}_{1:n-1}^{0}) p(\mathbf{x}_{n:N}^{0} | \mathbf{x}_{1:n-1}^{0}) d\mathbf{x}_{1:N}^{0}$$

$$= \int_{\mathbf{x}_{1:N}^{0}} \int_{\mathbf{x}_{n}^{0}} f(\mathbf{x}_{1:n}^{0}) p(\mathbf{x}_{1:n}^{0}) d\mathbf{x}_{n}^{0} \cdot p(\mathbf{x}_{n:N}^{0} | \mathbf{x}_{1:n-1}^{0}) d\mathbf{x}_{1:N}^{0}$$

$$= \int_{\mathbf{x}_{n:N}^{0}} \left(\int_{\mathbf{x}_{1:n}^{0}} f(\mathbf{x}_{1:n}^{0}) p(\mathbf{x}_{1:n}^{0}) d\mathbf{x}_{1:n}^{0} \right) p(\mathbf{x}_{n:N}^{0} | \mathbf{x}_{1:n-1}^{0}) d\mathbf{x}_{n:N}^{0}$$

$$= \int_{\mathbf{x}_{1:n}^{0}} f(\mathbf{x}_{1:n}^{0}) p(\mathbf{x}_{1:n}^{0}) d\mathbf{x}_{1:n}^{0}$$

 Since $\sum_{n=1}^{N} \mathbb{E}_{\mathbf{x}_{1:n}^{0}}[f(\mathbf{x}_{1:n}^{0})] = \mathbb{E}_{\mathbf{x}_{1:n}^{0}}[\sum_{n=1}^{N} f(\mathbf{x}_{1:n}^{0})] = L_{score}$, we prove the theorem.

B EXISTING TIME-SERIES DIFFUSION MODELS

There exist time-series diffusion models for forecasting and imputation (Rasul et al., 2021; Tashiro et al., 2021). However, our approach to time-series synthesis is technically distinct. While these models aim to generate fulfilled samples given partially known time-series, TSGM focuses on producing diverse samples. Furthermore, TSGM is designed to provide diverse regular time-series given regular or irregular data unlike time-series forecasting, which predicts future observations based on past data, and time-series imputation, which infers missing elements within a given time-series sample. This intuition is reflected in the experimental results in Section C.

We refer to a detailed analysis of these experiments in Section C.3.

B.1 DIFFUSION MODELS FOR TIME-SERIES FORECASTING AND IMPUTATION

TimeGrad (Rasul et al., 2021) is a diffusion-based method for time-series forecasting, and CSDI (Tashiro et al., 2021) is for time-series imputation.

In TimeGrad (Rasul et al., 2021), they used a diffusion model for forecasting future observations given past observations. On each sequential order $n \in \{2, ..., N\}$ and diffusion step $s \in \{1, ..., T\}$, they train a neural network $\epsilon_{\theta}(\cdot, \mathbf{x}_{1:n-1}, s)$ with a time-dependent diffusion coefficient $\bar{\alpha}_s$ by minimizing the following objective function:

$$\mathbb{E}_{\mathbf{x}_{n}^{0},\epsilon,s}\left[\left\|\epsilon-\epsilon_{\theta}\left(\sqrt{\bar{\alpha}_{s}}\mathbf{x}_{n}^{0}+\sqrt{1-\bar{\alpha}_{s}}\epsilon,\mathbf{x}_{1:n-1},s\right)\right\|_{2}^{2}\right],$$
(22)

where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The above formula assumes that we already know $\mathbf{x}_{1:n-1}$, and by using an RNN encoder, $\mathbf{x}_{1:n-1}$ can be encoded into \mathbf{h}_{n-1} . After training, the model forecasts future observations recursively. More precisely speaking, $\mathbf{x}_{1:n-1}$ is encoded into \mathbf{h}_{n-1} and the next observation \mathbf{x}_n is forecast from the previous condition \mathbf{h}_{n-1} .

807 CSDI (Tashiro et al., 2021) proposed a general diffusion framework which can be applied mainly 808 to time-series imputation. CSDI reconstructs an entire sequence at once, not recursively. Let $\mathbf{x}^0 \in \mathbb{R}^{\dim(\mathbf{X}) \times N}$ be an entire time-series sequence with N observations in a matrix form. They define 809 \mathbf{x}_{co}^0 and \mathbf{x}_{ta}^0 as conditions and imputation targets which are derived from \mathbf{x}^0 , respectively. They then



Figure 4: Graphical representation of TimeGrad (left) and TSGM (right). We adapt TimeGrad to our generation task but its results are not comparable even to other baselines' results (see Appendix C.1).

train a neural network $\epsilon_{\theta}(\cdot, \mathbf{x}_{co}^0, s)$ with a corresponding diffusion coefficient $\bar{\alpha}_s$ and a diffusion step $s \in \{1, ..., T\}$ by minimizing the following objective function:

$$\mathbb{E}_{\mathbf{x}^{0},\epsilon,s}[\left\|\epsilon - \epsilon_{\theta}(s, \mathbf{x}_{ta}^{s}, \mathbf{x}_{co}^{0})\right\|_{2}^{2}],\tag{23}$$

where $\mathbf{x}_{ta}^s = \sqrt{\bar{\alpha}_s} \mathbf{x}_{ta}^0 + (1 - \bar{\alpha}_s) \epsilon$. By training the network using the above loss, it generates missing elements from the partially filled matrix \mathbf{x}_{co}^0 .

B.2 DIFFERENCE BETWEEN EXISTING AND OUR WORKS

Although they have earned state-of-the-art results for forecasting and imputation, we found that they are not suitable for our generative task due to the fundamental mismatch between their model designs and our task (cf. Table 5 and Fig. 4).

Table 5: Comparison among various recent GAN, diffusion, and SGM-based methods for time-series. \mathbf{x}_t (resp. $\hat{\mathbf{x}}_t$) means a raw (resp. synthesized) observation at time t. For CSDI, \mathbf{x}_{co} means a set of known values and \mathbf{x}_{ta} means a set of target missing values — it is not necessary that \mathbf{x}_{co} precedes \mathbf{x}_{ta} in time in CSDI.

Method	Туре	Task Description
TimeGrad	Diffusion	From $\mathbf{x}_{1:N-K}$, infer $\hat{\mathbf{x}}_{N-K+1:N}$.
CSDI	Diffusion	Given known values \mathbf{x}_{co} , infer missing values $\hat{\mathbf{x}}_{ta}$.
TimeGAN	GAN	Synthesize $\hat{\mathbf{x}}_{1:N}$ from scratch.
GT-GAN	GAN	Synthesize $\hat{\mathbf{x}}_{1:N}$ from scratch.
TSGM	SGM	Synthesize $\hat{\mathbf{x}}_{1:N}$ from scratch.

TimeGrad generates future observations given the hidden representation of past observations \mathbf{h}_{n-1} , i.e., a typical forecasting problem. Since our task is to synthesize from scratch, past known observations are not available. Thus, TimeGrad cannot be directly applied to our task.

In CSDI, there are no fixed temporal dependencies between \mathbf{x}_{co}^0 and \mathbf{x}_{ta}^0 since its task is to impute missing values, i.e., \mathbf{x}_{ta}^0 , from known values, i.e., \mathbf{x}_{co}^0 , in the matrix \mathbf{x}^0 . It is not necessary that \mathbf{x}_{co}^0 precedes \mathbf{x}_{ta}^0 in time, according to the CSDI's method design. Our synthesis task can be considered on $\mathbf{x}_{co}^0 = \mathbf{0}$, which is the metric action of the CSDI's task is the matrix \mathbf{x}_{co}^0 . as $\mathbf{x}_{co}^0 = \emptyset$, which is the most extreme case of the CSDI's task. Therefore, it is not suitable to be used for our task.

To our knowledge, we are the first proposing an SGM-based time-series synthesis method. We propose to train a conditional score network by using the denoising score matching loss proposed by us, which is denoted as $L_{score}^{\mathcal{H}}$. Unlike other methods (Rasul et al., 2021; Tashiro et al., 2021) that resort to existing known proofs, we design our denoising score matching loss in Eq. equation 11 and prove its correctness. Meanwhile, TimeGrad and CSDI can be somehow modified for time-series synthesis but their generation quality is mediocre (see Appendix C).

С EXPERIMENTAL RESULTS FOR INAPPLICABILITY OF EXISTING TIME-SERIES DIFFUSION MODELS TO OUR WORK

In this section, we provide experimental results to show inapplicability of the existing time-series diffusion models, TimeGrad and CSDI, to the time-series generation task.

Table 6: Comparison between TSGM and modified TimeGrad in Energy for its regular time-series setting

Method	Disc.	Pred.
TSGM-VP	.221±.025	.257±.000
TSGM-subVP	.198±.025	.252±.000
Modified TimeGrad	.500±.000	.287±.003

Table 7: Comparison between TSGM and modified CSDI in Stock, Air, Energy and AI4I for its regular time-series setting

Method	Stock		Air		Energy		AI4I	
	Disc.	Pred.	Disc.	Pred.	Disc.	Pred.	Disc.	Pred.
TSGM-VP	$.022 \pm .005$	$.037 {\pm} .000$.122±.014	$.005 {\pm} .000$.147±.005	$.217 {\pm} .000$.221±.025	$.257 \pm .000$
TSGM-subVP	.021±.008	$.037 {\pm} .000$.127±.010	$.005 {\pm} .000$.150±.010	$.217 {\pm} .000$.198±.025	$.252 \pm .000$
Modified CSDI	.379±.008	$.045 \pm .001$.437±.144	$.040 {\pm} .001$.427±.081	$.217 \pm .000$	$.500 \pm .000$.251±.000

C.1 ADAPTING TIMEGRAD TOWARD GENERATION TASK

In this section, TimeGrad (Rasul et al., 2021) is modified for the generation task. We simply add an artificial zero vector **0** in front of the all time-series samples of Energy. Therefore, TimeGrad's task becomes given a zero vector, forecasting (or generating) all other remaining observations. For the stochastic nature of its forecasting process, it can somehow generate various next observations given the sample input **0**. Table 6 shows the experimental comparison between modified TimeGrad and TSGM in Energy for its regular time-series setting. TSGM gives outstanding performance, compared to modified TimeGrad. When checked in Table 2, modified TimeGrad is even worse than some baselines. Therefore, unlike TSGM, TimeGrad is not appropriate for the generation task.

C.2 ADAPTING CSDI TOWARD GENERATION TASK

In this section, we apply CSDI to the unconditional time-series generation task by regarding all observations as missing values (i.e., $\mathbf{x}_{co}^0 = \mathbf{0}$) varying i) its kernel size from 1 to 7 and ii) the number of diffusion steps from 50 to 250. However, as demonstrated in Table 7, CSDI fails to generate reliable time series samples in the datasets for its regular time series setting. In particular, TSGM with 250 steps in the ablation study section significantly outperforms it. Hence, we conclude that CSDI's unconditional generation is unsuitable for the time-series generation task.

C.3 ANALYSIS ON EXPERIMENTAL RESULTS

Until now, we have demonstrated the inferior results of forecasting and imputation on time-series generation. As shown by the results, these two models achieve relatively good predictive score but poor discriminative score, indicating a lack of diversity. This is because the primary goal of imputation and forecasting is to generate precise values that closely match the ground truth, as we discussed in Section B. For example, CSDI takes its imputed time-series by averaging synthesized samples 100 times. Therefore, as the quality of forecasting and imputation improves, the diversity of the generated samples decreases, which means they are not suitable for time-series generation.

D DETAILED TRAINING PROCEDURE

We train the conditional score network and the encoder-decoder pair alternately after the pre-training step. For some datasets, we found that training only the conditional score network achieves better results after pre-training the autoencoder. Therefore, $use_{alt} = \{True, False\}$ is a hyperparameter to set whether we use the alternating training method. We give the detailed training procedure in Algorithm 1.

918 Algorithm 1: Training algorithm 919 **Input:** $\mathbf{x}_{1:N}^0$; use_{alt} is a Boolean parameter to set whether to use the alternating training method; $iter_{pre}$ 920 is the number of iterations for pre-training; itermain is the number of iterations for training. 921 1 for $iter \in \{1, ..., iter_{pre}\}$ do 922 Train Encoder and Decoder by using L_{ed} 2 923 3 end 924 4 for $iter \in \{1, ..., iter_{main}\}$ do Train M_{θ} by using $L_{score}^{\mathcal{H}}$ 925 5 if use_{alt} then 6 926 Train the Encoder and Decoder by using L_{ed} 7 927 8 end 928 9 end 929 10 return Encoder, Decoder, M_{θ} 930 931 932 E

933

934

935 936

937

938

939

940

941

942 943

944

945

946

947

948

DATASETS AND BASELINES

We use 4 datasets from various fields as follows. We summarize their data dimensions, the number of training samples, and their time-series lengths (window sizes) in Table 8.

- Stock (Yoon et al., 2019): The Google stock dataset was collected irregularly from 2004 to 2019. Each observation has (volume, high, low, opening, closing, adjusted closing prices), and these features are correlated.
- Energy (Candanedo et al., 2017): This dataset is from the UCI machine learning repository for predicting the energy use of appliances from highly correlated variables such as house temperature and humidity conditions.
 - Air (De Vito et al., 2008): The UCI Air Quality dataset was collected from 2004 to 2005. Hourly averaged air quality records are gathered using gas sensor devices in an Italian city.
- A141 (Matzka, 2020): A14I means the UCI A14I 2020 Predictive Maintenance dataset. This data reflects the industrial predictive maintenance scenario with correlated features including several physical quantities.

949 We use several types of generative methods for time-series as baselines. At first, we consider autoregressive generative methods: T-Forcing (teacher forcing) (Graves, 2013; Sutskever et al., 950 2011) and P-Forcing (professor forcing) (Goyal et al., 2016). Next, we use GAN-based methods: 951 TimeGAN (Yoon et al., 2019), RCGAN (Esteban et al., 2017), C-RNN-GAN (Mogren, 2016), COT-952 GAN (Xu et al., 2020), GT-GAN (Jeon et al., 2022). We also test VAE-based methods into our 953 baselines: TimeVAE (Desai et al., 2021). Finally, we treat flow-based methods. Among the array 954 of flow-based models designed for time series generation, we have chosen to compare our TSGM 955 against CTFP (Deng et al., 2020). This choice is informed by the fact that CTFP possesses the capa-956 bility to handle both regular and irregular time series samples, aligning well with the nature of our 957 task which involves generating both regular and irregular time series data.

958 959 960

Table 8: Characteristics of the datasets we use for our experiments

Dataset	Dimension	#Samples	Length
Ctul	6	2605	Zengu
Stocks	0	3085	
Energy	28	19735	24
Air	13	9357	24
AI4I	5	10000	

965 966 967

HYPERPARAMETERS AND ITS SEARCH SPACE F

968 969

Table 9 shows the best hyperparameters for our conditional score network M_{θ} on regular and irreg-970 ular time-series, and we explain its neural network architecture in Appendix J.2. M_{θ} has various 971 hyperparameters and for key hyperparameters, we set them as listed in Table 9. For other common Table 9: The best hyperparameter setting for TSGM on regular and irregular time-series . D_{hidden}
denotes the hidden dimension of GRU-ODE-decoder.

975		1	Pagulo	r Sottings		1	In	ogular Satti	200	
976	Dataset	dim(h)	use_{alt}	iter _{pre}	$iter_{main}$	D _{hidden}	dim(h)	use _{alt}	iter _{pre}	$iter_{main}$
977	Stocks	24	True	50000		48	24	True		
978	Air	40	True	50000	40000	40	30 40	True	50000	40000
979	AI4I	24	True	50000		48	24	True		
980										
981	hyperparat	notore wit	h hasalin	ac warai	ica tha dafe	ault configu	irations of	fTimeG	AN (Voor	a = 1 - 2010
982	and VPSD	E (Song e	t al 202	1) to con	duct the re	oular time-	series gei	r ration		l ct al., 2019)
983		E (boing c	t ul., 202	1) to com		Sului tille	Series ger	forution.		
984	We give ou	ir search s	pace for	the hyper	parameter	s of TSGM	$I.iter_{pre}$	is in {50	000,1000	00}. The di-
985	mension of	t hidden te	eatures, d	hidden, ra	anges from	2 times to	5 times the	he dimen	sion of in	put features.
986	On regular	E (Song o	es generation de la g	1) Eor in	regular tir	e default v	alues in	nmeGA	N (1001) hiddon (dimension of
987	decoder fr	E (Solig c om 2 time	s to 4 tim	1). FUI II	nension of	f input dim	ension at	real follow	r GTGAN	Inficiision of al
988	2022) for c	other settir	ngs of NC	DE-enco	oder and G	RU-ODE-	decoder	iu ionow	UIUAN	(Jeon et al.,
989	2022) 101 (Julier Settin	155 01 110				accoact.			
990	For baselir	ies, we ch	eck their	hyperpar	ameters as	s follow:				
991	• Т	foreing (Groups 1	$(12), W_{2}$	control b	tah siza ar	nong (25	6 512 1	024]	
992	• 1.	-toteing (C		015). we			nong {23	0, 512, 1	024}.	
993	• P.	-forcing (C	Goyal et a	al., 2016)	: We conti	ol batch si	ze among	{256, 5]	12, 1024	•
994	• T	imeGAN	(Yoon et	al., 2019): The din	nension of	hidden fe	atures ra	nge from	2 times to 4
995	ti	mes the di	mension	of input	features.					
996	• R	CGAN (E	Esteban et	al., 2017	7): We con	trol learnin	ng rate of	generato	or's optim	izer and dis-
997	cı	riminator's	s optimiz	er from {	1e-4, 2e-4	} and {1e-	3, 5e-3},	respectiv	ely.	
998	• C-RNN-GAN (Mogren 2016): We control learning rate of generator's optimizer and dis									
999	ci	riminator's	s optimiz	er from {	1e-4, 2e-4	$and {3e}$	4.4e-4	respectiv	elv. We a	lso use label
1000	sr	noothing	which is a	stated in	the paper.)	, .) ,		J	
1001	• T	imeVAE (Desai et :	al 2021)	• We conti	rol its laten	t dimensi	on amon	σ{5 10	20}
1002			(Vrs. et. al.	2020).	We select	-4		on anton	5(0, 10)	20j.
1003 1004	• C	et the best	experime	., 2020): ental resu	ilts.	ate score e	very 250	epoch du	ning 1000) epochs and
1005	• G	T-GAN (.	Jeon et a	1., 2022)	: For enco	oder-decod	ler pair, v	ve test fr	om exact	tly the same
1006	se	earch spac	e as TSC	GM. We	calculate s	core every	5000 ite	ration du	ring 400	00 iterations
1007	aı	nd get the	best score	e.						
1008	Especially	for COT-	-GAN, si	nce it is	on video	generation	, modifyi	ng the a	rchitectur	e to one di-
1009	mensional	form was	difficult	. So, we	augment	our time-se	eries data	into two	dimensio	onal ones by
1010	stacking th	nem. Afte	r generat	ing two-o	dimension	al data, we	extract th	he first ro	ow of the	synthesized
1011	one and ca	lculate the	e score. W	Ve search	every hyp	erparamete	r from {0	.5, 1, 2} 1	times of d	lefault value.
1012	Through th	ie experim	ient, we a	acquire co	ompatible	result but lo	ower than	TimeGA	N in seve	eral datasets.
1013	We follow	default va	lues for r	niscellan	eous settin	gs which a	re not exp	lained or	1 the abov	ve. Addition-
1014	ally, to dea	al with irre	gular tim	ne-series,	we search	the hyperp	parameter	s of GRU	J-D, whic	h substitutes
1015	for RNN c	or are adde	ed to the	head of	baselines.	We test th	ne hidden	dimensi	on of GR	U-D from 2
1016	times to 4	times the o	dimensio	n of inpu	t features.					
1017										
1018	G MIS	CELLAN	IEOUS F	EXPERT	MENTAL		NMENT	S		
1019	- 1110		I					-		
1020	We give de	etailed ext	periment	al enviror	ments. T	he followir	ig softwa	re and ha	ardware e	nvironments
1021	were used	for all exr	periments	: UBUNI	U 18.04 L	TS, Pytho	ON 3.9.12	, CUDA	9.1, NV	IDIA Driver
1022	470.141, i	J CPU, an	d GEFOF	RCE RTX	2080 TI.	,				
1023										

1024 In the experiments, we report only the VP and subVP-based TSGM and exclude the VE-based 1025 one for its lower performance. For baselines, we reuse their released source codes in their official repositories and rely on their designed training and model selection procedures. For our method, we Table 10: The memory usage of TimeGAN, GTGAN, and TSGM for training and the sampling time of them for generating 100 samples on each dataset. Note that The original score-based model (Song et al., 2021) requires 3,214 seconds for sampling 1000 CIFAR-10 images while StyleGAN (Karras et al., 2019) needs 0.4 seconds, which is similar to the case between TSGM and TimeGAN.

1031		Memory V	Usage (GiB)	Inference Time (s)		
1032	Method	Stock	Energy	Stock	Energy	
1033	TimeGAN	1.1	1.6	0.43	0.47	
1034	GTGAN	2.3	2.3	0.43	0.47	
1035	TSGM	1.9	1.9	86.32	85.89	

1036 1037

1030

select the best model for every 5000 iterations. For this, we synthesize samples and calculate the mean and standard deviation scores of the discriminative and predictive scores. Furthermore, TSGM requires 100 sampling steps, which is relatively lower compared to other diffusion-based models—
CSDI (Tashiro et al., 2021) requires 50 steps but averages samples over 100 iterations, effectively requiring 5000 steps.

1042 1043

1044 H EMPIRICAL SPACE AND TIME COMPLEXITY ANALYSES

1045

1054 1055

1056 1057

1062 1063 1064

1046 We report the memory usage during training and the wall-clock time for generating 100 time-series 1047 samples in Table 10. We compare TSGM to TimeGAN (Yoon et al., 2019) and GTGAN (Jeon 1048 et al., 2022). Our method is relatively slower than TimeGAN and GTGAN, which is a fundamental 1049 drawback of all SGMs. For example, the original score-based model (Song et al., 2021) requires 1050 3,214 seconds for sampling 1,000 CIFAR-10 images while StyleGAN (Karras et al., 2019) needs 0.4 seconds. However, we also emphasize that this problem can be relieved by using the techniques 1051 suggested in (Xiao et al., 2022; Jolicoeur-Martineau et al., 2021) as we mentioned in the conclusion 1052 section. 1053

I ENCODER AND DECODER FOR IRREGULAR TIME-SERIES

To process irregular time-series, one can use continuous-time methods for constructing the encoder and the decoder. In our case, we use neural controlled differential equations (NCDEs) for designing the encoder and GRU-ODEs for designing the decoder, respectively (Kidger et al., 2020; Brouwer et al., 2019). Our encoder based on NCDEs can be defined as follows:

$$\mathbf{h}(t_n) = \mathbf{h}(t_{n-1}) + \int_{t_{n-1}}^{t_n} f(t, \mathbf{h}(t); \theta_f) \frac{dX(t)}{dt} dt,$$
(24)

where X(t) is an interpolated continuous path from $\mathbf{x}_{1:N}$ — NCDEs typically use the natural cubic spline algorithm to define X(t), which is twice differentiable and therefore, there is not any problem to be used for forward inference and backward training. In other words, NCDEs evolve the hidden state $\mathbf{h}(t)$ by solving the above Riemann-Stieltjes integral.

For the decoder, one can use the following GRU-ODE-based definition:

1071
1072
$$\overline{\mathbf{d}}(t_n) = \mathbf{d}(t_{n-1}) + \int_{t_{n-1}}^{t_n} g(t, \mathbf{d}(t); \theta_g) dt, \qquad \mathbf{d}(t_n) = \mathrm{GRU}(\mathbf{h}(t_n), \overline{\mathbf{d}}(t_n)), \qquad \hat{\mathbf{x}}_n = FC(\mathbf{d}(t_n)),$$
(25)

1074

where FC denotes a fully-connected layer-based output layer. The intermediate hidden representation $\overline{\mathbf{d}}(t_n)$ is jumped into the hidden representation $\mathbf{d}(t_n)$ by the GRU-based jump layer. At the end, there is an output layer.

1079 For our irregular time-series experiments, i.e, dropping 30%, 50%, and 70% of observations from regular time-series, we use the above encoder and decoder definitions and have good results.

¹⁰⁸⁰ J NEURAL NETWORK ARCHITECTURE

1082 J.1 ARCHITECTURAL DETAILS OF NCDES AND GRU-ODES

Laver

1

2

3

4

Gate

 r_t

 z_t

 u_t

Layer

1

1084 As mentioned in Appendix I, we take the following architecture for functions f, g of (24) and (25) 1085 in Table 11.

Table 11: Architecture of functions f(upper) and g(lower). Each layer of encoder and gate of decoder takes ($\sigma \circ Linear$) form where σ denotes activation function. We describe which activation and Linear function are used.

Activation function

ReLU

ReLU

Tanh

Linear

 $\dim(\mathbf{x}) \rightarrow 4 \dim(\mathbf{x})$

 $4 \dim(\mathbf{x}) \rightarrow 4 \dim(\mathbf{x})$

 $4 \dim(\mathbf{x}) \rightarrow 4 \dim(\mathbf{x})$

 $4 \dim(\mathbf{x}) \rightarrow \dim(\mathbf{x})$

Linear

 $\dim(\mathbf{h}) \rightarrow \dim(\mathbf{h})$

Activation function

ReLU

ReLU

ReLU

Tanh

1089 1090

> 091 092

1093 1094

1	095	
1	096	

097

- 1098 1099
- 1100
- 1101

1102 J.2 CONDITIONAL SCORE NETWORK

Unlike other generation tasks, e.g., image generation (Song et al., 2021) and tabular data synthesis (Kim et al., 2022), where each sample is independent, time-series observations are dependent to their past observations. Therefore, the score network for time-series generation must be designed to learn the conditional log-likelihood given past generated observations, which is more complicated than that in image generation.

In order to learn the conditional log-likelihood, we modify the popular U-net (Ronneberger et al., 2015) architecture for our purposes. Since U-net has achieved various excellent results for other generative tasks (Song and Ermon, 2019; Song et al., 2021), we modify its 2-dimensional convolution layers to 1-dimensional ones for handling time-series observations. The modified U-net, denoted M_{θ} , is trained to learn our conditional score function (cf. Eq. equation 11). More details on training and sampling with M_{θ} are in Sec. 3.4.

1115 1116

K ADDITIONAL EXPERIMENTAL RESULTS

We give additional experimental results for irregular time-series generation with 50% and 70% missing rates in Table 12.

1120

1121 L ADDITIONAL VISUALIZATIONS

In this section, we provide additional visualization results in each dataset. Figure 5 illustrates the density function of each feature estimated by KDE in original and generated data. Figure 6 shows original and generated data points projected onto a latent space using t-SNE (van der Maaten and Hinton, 2008)

1127

1128 M EFFICACY OF OUR RECURSIVE GENERATION

1129

In this section, we investigate the efficacy of our proposed recursive design. We compare TSGM to an method using *one-shot generation*. we call *one-shot generation* when a generation method generates all time-series observations at once, not recursively. In other words, $D \times N$ matrices, where D means the number of features and N means the sequence length, are synthesized at once. CSDI (Tashiro et al., 2021) is one of the most famous one-shot imputation model for time-series. 1156

1166 1167 1168

1174

Table 12: The left and right ones denote experimental results on irregular time-series with 30% and 50% missing rates, respectively. Note that except for our method (TSGM), CTFP, and GT-GAN, the other methods cannot deal with irregular time series, so we make it possible for them to operate on irregular settings by replacing RNN encoder with GRU-D.

1138										
1139			Irregular Settings (Missing Rate: 50%)				Irregular Settings (Missing Rate: 70%)			
1140		Method	Stocks	Energy	Air	AI4I	Stocks	Energy	Air	AI4I
1141		TSGM-VP TSGM-subVP	.051±.014 .031+.012	$.398 \pm .003$ $421 \pm .008$.272±.012	.156±.106	.065±.010 .035±.009	.482±.003	.337±.025	$.327 \pm .104$ $.235 \pm .123$
1142		T-Forcing	.407±.034	.376±.046	.499±.001	.473±.045	.404±.068	.336±.032	.499±.001	.493±.010
1174		P-Forcing	$.500 \pm .000$	$.500 \pm .000$	$.494 {\pm} .012$	$.437 {\pm} .079$.449±.150	$.494 {\pm} .011$	$.498 {\pm} .002$	$.440 \pm .125$
1143	ore	TimeGAN	.477±.021	.473±.015	$.500 {\pm} .001$	$.500 {\pm} .000$.485±.022	$.500 {\pm} .000$	$.500 {\pm} .000$	$.500 {\pm} .000$
1144	SC	RCGAN	$.500 \pm .000$	$.500 \pm .000$	$.500 {\pm} .000$	$.500 {\pm} .000$.500±.000	$.500 \pm .000$	$.500 {\pm} .000$	$.500 \pm .000$
	sc.	C-RNN-GAN	$.500 \pm .000$	$.500 {\pm} .000$	$.500 {\pm} .000$	$.450 \pm .150$.500±.000	$.500 {\pm} .000$	$.500 {\pm} .000$	$.500 {\pm} .000$
1145	D	TimeVAE	.411±.110	$.436 {\pm} .088$	$.423 \pm .153$.389±.113	.444±.148	$.498 {\pm} .003$	$.426 \pm .148$	$.371 \pm .092$
1146		COT-GAN	.499±.001	$.500 \pm .000$	$.500 {\pm} .000$	$.500 {\pm} .000$.498±.001	$.500 {\pm} .000$	$.500 {\pm} .000$	$.500 \pm .000$
		CTFP	.499±.000	$.500 \pm .000$	$.500 {\pm} .000$	$.499 {\pm} .001$.500±.000	$.500 {\pm} .000$	$.500 {\pm} .000$	$.499 {\pm} .000$
1147		GT-GAN	.265±.073	$.317 {\pm} .010$	$.434 {\pm} .035$	$.276 \pm .033$.230±.053	$.325 \pm .047$	$.444 {\pm} .019$	$.362 \pm .043$
1148		TSGM-VP	.011±.000	$.051 {\pm} .001$.041±.001	$.060 {\pm} .001$.011±.000	$.053 \pm .001$	$.043 \pm .000$	$.092 \pm .024$
1140		TSGM-subVP	.011±.000	$.051 {\pm} .001$	$.042 \pm .002$	$.065 \pm .013$.012±.000	$.042 {\pm} .002$	$.042 \pm .001$	$.097 {\pm} .020$
1149		T-Forcing	$.038 \pm .003$	$.090 \pm .000$.121±.003	$.143 \pm .005$.031±.002	$.091 \pm .000$	$.116 \pm .003$	$.144 \pm .004$
1150		P-Forcing	.089±.010	$.198 \pm .005$	$.101 \pm .003$	$.116 \pm .007$.107±.009	$.193 \pm .006$	$.107 \pm .002$	$.125 \pm .007$
	e	TimeGAN	.254±.047	$.339 {\pm} .029$	$.325 \pm .005$	$.251 \pm .010$.228±.000	$.443 {\pm} .000$	$.425 \pm .008$	$.323 \pm .011$
1151	00	RCGAN	.333±.044	$.250 \pm .010$	$.335 \pm .023$	$.276 \pm .066$.441±.045	$.349 \pm .027$	$.359 \pm .008$	$.346 \pm .029$
1152		C-RNN-GAN	$.273 \pm .000$	$.438 {\pm} .000$	$.289 \pm .033$	$.373 \pm .037$.281±.019	$.436 \pm .000$	$.306 \pm .040$	$.262 \pm .053$
	rec	TimeVAE	.195±.012	$.143 \pm .007$	$.103 \pm .002$	$.144 \pm .004$.199±.009	$.134 \pm .004$	$.108 \pm .004$	$.142 \pm .008$
1153	щ	COT-GAN	$.246 \pm .000$	$.475 \pm .000$	$.557 \pm .000$	$.449 \pm .000$.278±.000	$.456 \pm .000$	$.556 \pm .000$	$.435 \pm .000$
1154		CTFP	$.084 \pm .005$	$.469 {\pm} .008$	$.476 \pm .235$	$.412 \pm .024$.084±.005	$.469 {\pm} .008$.476±.235	$.412 \pm .024$
		GT-GAN	.018±.002	$.064 \pm .001$	$.061 \pm .003$	$.113 \pm .024$.020±.005	$.076 \pm .001$	$.059 \pm .004$	$.124 \pm .003$
1155		Original	.011±.002	$.045 \pm .001$	$.044 \pm .006$	$.059 \pm .001$.011±.002	$.045 \pm .001$	$.044 \pm .006$	$.059 \pm .001$



Figure 5: Additional KDE plots for each feature in Air and AI4I datasets.

We convert our TSGM for the one-shot generation by removing the RNN-based encoder. In Table 13, TSGM-oneshot shows poor generation quality in Stock and Energy. TSGM-oneshot achieves comparable predictive scores but its discriminative score gets worse a lot. From these results, we can support the efficacy of our recursive structures, compared to one-shot generation. One can also check the one-shot generation result by CSDI in Appendix C.2.



Figure 6: Additional t-SNE plots in Energy and AI4I datasets.

Table 13: Comparison between TSGM and one-shot generations. We give representative results.For other datasets, the results are similar or worse than the table.

Mathod	S	tock	Energy		
Methou	Disc.	Pred.	Disc.	Pred.	
TSGM-VP	.022±.005	5 .037±.000	.221±.025	$.257 \pm .000$	
TSGM-subV	P .021±.008	8 .037±.000	.198±.025	$.252 {\pm} .000$	
TSGM-onesh	ot .029±.018	8 .037±.000	.494±.001	$.258 {\pm} .000$	

Table 14: Comparison of drift and diffusion terms. $\sigma(s)$ means positive noise values which are increasing, and $\beta(s)$ denotes noise values in [0,1], which are used in Song and Ermon (2019); Ho et al. (2020).

1200 1201

1202

1203

1205 1206 1207

1208

N DRIFT AND DIFFUSION TERMS IN VE, VP AND SUBVP SDE

In this section, we describe the detailed form of each SDE. In (Song et al., 2021), the authors investigated that SMLD (Song and Ermon, 2019) and DDPM (Ho et al., 2020) can be extended to continuous forms and as a result, suggested VE and VP SDEs. Furthermore, the author proposed an additional SDE form, called subVP SDE, which has a smaller variance than VP SDE but the same expectation. The exact calculation is not a main subject of this paper, so we only explain the form of these terms in Table 14. Please refer to (Song et al., 2021) for the detailed computation.

1215 Along with already mentioned notations in Section 2.1, we define noise scales. $\sigma(s)$ means positive 1216 noise values which are increasing, and $\beta(s)$ denotes noise values in [0,1] which are used in SMLD 1217 and DDPM. Although we give the exact form of the three SDEs, we report only the VP and subVP-1218 based TSGM in our experiments and exclude the VE-based one for its lower performance.

1219 1220

1221

O DETAILED DESCRIPTION OF GRU-D

GRU-D (Che et al., 2016) is a modified GRU model which is for learning time-series data with missing values. This concept is similar with our problem statement, so we apply it to our baseline for irregular case. GRU-D needs to learn decaying rates along with the values of GRU. First, GRU-D learns decay rates which depict vagueness of data as time passed. After calculating the decay rates, each value is composed of decay rate, mask, latest observed data, and predicted empirical mean that of GRU. The code can be utilized in the following link: https://github.com/zhiyongc/GRU-D

- 1228
- 1229
- 1230
- 1231
- 1232 1233
- 1234
- 1235
- 1236
- 1237
- 1238
- 1239

1240