

L-GTA: Latent Generative Modeling for Time Series Augmentation

Luis Roque

luis_roque@live.com

LIACC/Faculty of Engineering, University of Porto
Porto, Portugal

Vitor Cerqueira

LIACC/Faculty of Engineering, University of Porto
Porto, Portugal

Carlos Soares

LIACC/Faculty of Engineering, University of Porto
Porto, Portugal
Fraunhofer AICOS Portugal
Porto, Portugal

Luís Torgo

Dalhousie University
Halifax, Canada

Abstract

Data augmentation is gaining importance across various aspects of time series analysis, from forecasting to classification and anomaly detection tasks. We introduce the Latent Generative Transformer Augmentation (L-GTA) model, a generative approach using a transformer-based variational recurrent autoencoder. This model uses controlled transformations within the latent space of the model to generate new time series that preserve the intrinsic properties of the original dataset. L-GTA enables the application of diverse transformations, ranging from simple jittering to magnitude warping, and combining these basic transformations to generate more complex synthetic time series datasets. Our evaluation of several real-world datasets demonstrates the ability of L-GTA to produce more reliable, consistent, and controllable augmented data. This translates into significant improvements in predictive accuracy and similarity measures compared to direct transformation methods.

CCS Concepts

• **Computing methodologies** → **Machine learning**; • **Mathematics of computing** → *Time series analysis*.

Keywords

time series augmentation, generative models, variational autoencoder, transformer, synthetic data generation, anomaly detection, classification, forecasting

ACM Reference Format:

Luis Roque, Carlos Soares, Vitor Cerqueira, and Luís Torgo. 2025. L-GTA: Latent Generative Modeling for Time Series Augmentation. In *Proceedings of Machine Learning in Time Series (MILETS) Workshop at the ACM SIGKDD Conference on Knowledge Discovery and Data Mining (MILETS @ KDD '25)*. ACM, New York, NY, USA, 9 pages. <https://doi.org/XXXXXXX.XXXXXXX>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
MILETS @ KDD '25, Toronto, ON, Canada

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/2018/06
<https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

In the era of big data, time series analysis has emerged as a critical tool in various domains ranging from financial markets [3] to health monitoring [6], from climate forecasting [5] to retail sales predictions [13]. Accurate and insightful time series analysis can lead to predictive models that provide crucial information for decision-making processes. Nevertheless, the performance of these models significantly depends on the quality and quantity of the available data. Also, collecting labeled time series data for a specific task or domain can be challenging.

Similarly, real-world time series data often exhibit complex dependencies, non-linear dynamics, and high dimensionality. They may also be affected by noise, irregular sampling, and missing values, further complicating their analysis and modeling. As a result, models trained on such data often fail to generalize well, leading to poor predictive performance when applied to unseen data.

In such scenarios, data augmentation techniques, which generate synthetic data samples from the original data, offer a promising solution. They increase the quantity of the data and improve model robustness by providing diversified data instances [21]. The problem with traditional data augmentation techniques for time series data, such as jittering, scaling, and warping, is that they are relatively simple and may not adequately capture the complexities often found in relevant datasets. Furthermore, these techniques are difficult to apply in a controlled manner, which may introduce artificial distortions that deviate significantly from real-world scenarios. Consequently, this reduces the practical utility of the models trained on such data.

We introduce the Latent Generative Transformer Augmentation (L-GTA) model to address these limitations (see Figure 1). L-GTA generates semi-synthetic time series data that merges the capabilities of Transformers, Bidirectional Long Short-Term Memory Networks (Bi-LSTMs), conditional Variational Autoencoders (CVAEs), and traditional time series augmentation techniques. We introduce a Variational Multi-Head Attention (VMHA) mechanism to capture more nuanced and long-term temporal dependencies in the data. The variational component introduces randomness in a controlled manner, making the attention mechanisms more robust and flexible. Alongside VMHA, we use Bi-LSTMs to capture the more immediate temporal dependencies within the data, while CVAEs construct a probabilistic, low-dimensional latent space that accurately represents the dataset. In this learned latent space, we apply traditional

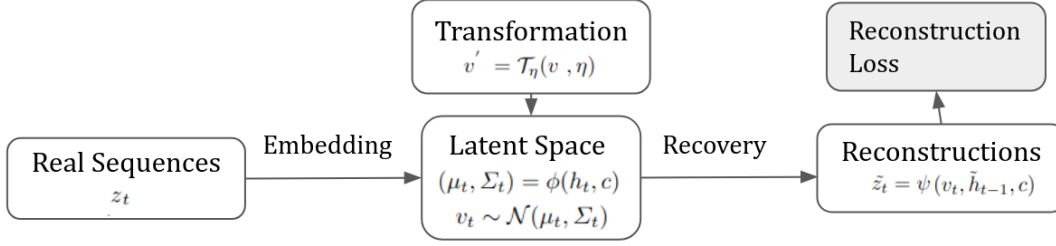


Figure 1: Diagram depicting the L-GTA framework for augmenting time series data. This method combines Bi-LSTMs and CVAEs with a Variational Multi-Head Attention (VMHA) mechanism to generate semi-synthetic time series data. The learned latent space is transformed by applying time series augmentation techniques such as jittering and magnitude warping. The augmented latent representations are then decoded, producing data that maintains the statistical fidelity of the original sequences. L-GTA allows for robust and diverse data transformed data.

time series augmentation techniques, like jittering and magnitude warping. The transformed latent space is then passed through the learned decoder. This process consistently yields augmented data that mirrors the effects of direct transformations while more strictly respecting the statistical properties of the original data. As a result, the produced data is less prone to exhibit artificial patterns and extreme values. For example, we demonstrate that the residuals of the datasets generated after applying jittering directly or using L-GTA are different. Direct transformations produce residuals that are more skewed than those generated by L-GTA. L-GTA also offers additional variety while controlling for the amount of change in augmented data that is generated. In the current set of generative models (e.g., [22]), we have to rely on randomly sampling the latent space to generate new data. With our proposed approach, the latent space is manipulated in a controlled way, which creates the ability to generate data with certain characteristics. This is relevant to increase robustness and diversity in time series analysis.

We evaluate the L-GTA method against traditional direct augmentation techniques using three real-world datasets, focusing on its ability to preserve the statistical properties and predictive characteristics of the original data. Through a series of experiments, we assess the fidelity of the transformed data by comparing the generated patterns and the control and consistency of L-GTA. We measured the Wasserstein distance and reconstruction error alongside the Train-on-Synthetic, Test-on-Real (TSTR) framework. Our findings highlight the superior performance of L-GTA in maintaining the original data integrity, as demonstrated by lower Wasserstein distances, minimal deviation in reconstruction error, and prediction accuracies in line with those achieved using the original datasets.

The primary contributions of this paper are as follows:

- We propose a novel generative model, L-GTA, which leverages a transformer-based CVAE for controlled time series augmentation;
- We introduce a flexible method to apply transformations, such as jittering and magnitude warping, to lower dimensional embeddings of a dataset. These transformations are then propagated coherently to generate new transformed time series that incorporate both the transformation and the intrinsic characteristics of the original dataset. These transformations can be combined in various ways to create a

large set of diverse synthetic time series, thereby improving the versatility of the augmentation process;

- We empirically tested L-GTA, comparing it with direct transformation methods. We analyzed the patterns of the transformed data and the control and consistency of L-GTA when generating it. We also discussed three perspectives to assess the fidelity of the generated data compared to the original: the distribution of distances, the reconstruction error, and the prediction performance.

All experiments are fully reproducible, and the methods and time series data are available as a publicly available code repository.¹

2 Notation and Background

2.1 Time Series

Let us consider a set of S related univariate time series, represented as $\mathcal{Z} = \{z_t^i : t \in \mathbb{N}, i = 1, \dots, S\}$, where the sequence $z_{1:T}^i = [z_1^i, z_2^i, \dots, z_T^i]$ denotes the observed values of the i -th time series up to the final observation time T . Each observation $z_t^i \in \mathbb{R}$ indicates the value of series i at time t . To simplify discussions in certain contexts, we denote $\mathbf{z}^i = z_{1:T}^i$ as the complete observed time series for the i -th series. The process of time series augmentation involves generating new time series data, represented as $\tilde{\mathbf{z}}^i$, which is derived from, yet distinct from, the original series \mathbf{z}^i . The goal is to enhance or modify the dataset for further analysis or model training while preserving the intrinsic properties of the time series.

To quantitatively assess the closeness of the transformed time series $\tilde{\mathbf{z}}^i$ to the original \mathbf{z}^i , we compute the first Wasserstein distance, denoted by W_1 . Our decision to use this measure follows the approach used by [17] for training and evaluating a conditional generative model. The Wasserstein distance reflects the true geometric structure of probability distributions. It considers not just simple statistical metrics but the actual distribution of probability mass across the domain. Hence, it offers a more nuanced perspective of the comparison between time series. We apply W_1 to assess how closely the distributional characteristics of the transformed time series $\tilde{\mathbf{z}}^i$ align with those of the original series \mathbf{z}^i . For a pair of time series \mathbf{z}^i and $\tilde{\mathbf{z}}^i$, the Wasserstein distance is computed as:

¹<https://github.com/luisroque/latent-generative-modeling-time-series-augmentation>

$$W_1(\mathbf{z}^i, \tilde{\mathbf{z}}^i) = \inf_{\gamma \in \Gamma(\mathbf{z}^i, \tilde{\mathbf{z}}^i)} \int_{\mathbb{R} \times \mathbb{R}} |z_t^i - \tilde{z}_t^i| d\gamma(z_t^i, \tilde{z}_t^i) \quad (1)$$

where $\Gamma(\mathbf{z}^i, \tilde{\mathbf{z}}^i)$ is the set of all joint distributions γ whose marginals are the empirical distributions of \mathbf{z}^i and $\tilde{\mathbf{z}}^i$.

In our approach, we employ a transformation function \mathcal{T} to modify the latent space represented by v_t within the context of a CVAE. By applying \mathcal{T} , we manipulate the latent space to generate diverse augmented versions of the data. Specifically, we define the transformation of the latent space as $v'_t = \mathcal{T}_\eta(v_t)$, where v'_t denotes the transformed latent space and η represents the parameters guiding the transformation. The transformed latent space, v'_t , when passed through the decoder of the CVAE, produces the augmented time series data. This process enables us to explore a variety of augmentations by changing the transformation \mathcal{T} and adjusting the parameters η .

2.2 Bi-directional Long Short-Term Memory Networks

Long Short-Term Memory (LSTM) networks [9] are a type of Recurrent Neural Network (RNN) that effectively handle the issue of long-term dependencies in sequences. They accomplish this with a gating mechanism that selectively forgets and updates the cell state at each time step.

Bidirectional LSTM (Bi-LSTM) networks [8] are an extension of LSTMs that process the data in both forward and backward directions. Given an input sequence $\mathbf{z} = (z_1, z_2, \dots, z_T)$, a Bi-LSTM consists of a forward LSTM and a backward LSTM. The forward LSTM reads the sequence in the forward direction to produce a sequence of hidden states \vec{h}_t , and the backward LSTM reads the sequence in the backward direction to produce a sequence of hidden states \overleftarrow{h}_t . For a time step t , the forward and backward hidden states are given by $\vec{h}_t = \text{LSTM}(z_t, \vec{h}_{t-1})$ and $\overleftarrow{h}_t = \text{LSTM}(z_t, \overleftarrow{h}_{t+1})$.

The hidden states of the Bi-LSTM at each time step t , h_t , are then obtained by concatenating the forward and backward hidden states $h_t = [\vec{h}_t; \overleftarrow{h}_t]$. Bi-LSTMs capture both past and future information for a given time step.

2.3 Conditional Variational Autoencoders

Autoencoders (AEs) are neural network architectures that aim to reconstruct their input by first encoding it into a latent space and then decoding it back to the original space. Formally, an autoencoder consists of two components: an encoder function ϕ and a decoder function ψ . Given an input $z \in \mathbb{R}^d$, the encoder maps z to a latent representation, $v = \phi(z)$, $v \in \mathbb{R}^p$. The decoder then maps v back to the original space to produce a reconstruction $\tilde{z} = \psi(v)$, $\tilde{z} \in \mathbb{R}^d$. Autoencoders are typically trained by minimizing the reconstruction error between the original input z and the reconstruction \tilde{z} .

Variational Autoencoders (VAEs) [14] are a generative variant of autoencoders that introduce a probabilistic approach to the encoding process. Instead of directly mapping an input z to a deterministic latent representation v , a VAE maps z to a distribution over the latent space. The encoder of a VAE, also known as the recognition model, thus outputs the parameters of a Gaussian distribution, typically the mean μ and the diagonal covariance Σ . Given an input z ,

the encoder outputs $(\mu, \Sigma) = \phi(z)$. A latent variable v is then sampled from this Gaussian distribution, $v \sim \mathcal{N}(\mu, \Sigma)$. The decoder of a VAE maps v back to the original space to produce \tilde{z} , a reconstruction of z (a process similar to AEs). The VAE is trained by maximizing the Evidence Lower Bound (ELBO) on the log-likelihood of the data:

$$\log p(z) \geq \mathbb{E}q(v|z)[\log p(z|v)] - D_{KL}(q(v|z)||p(v)) \quad (2)$$

where $q(v|z)$ is the approximate posterior, $p(v)$ is the prior, and D_{KL} is the Kullback-Leibler divergence.

Conditional Variational Autoencoders (CVAEs) [20] extend VAEs by introducing a conditioning variable c into the model. This allows the generation process to be conditioned on specific features or characteristics. The encoder of a CVAE thus maps an input z and a condition c to a distribution over the latent space. Given z and c , the encoder outputs $(\mu, \Sigma) = \phi(z, c)$. A latent variable v is then sampled from this Gaussian distribution and the decoder of a CVAE takes the sampled v and the condition c to produce \tilde{z} , a reconstruction of z .

2.4 Variational Self-Attention Mechanism

The Variational Self-Attention Mechanism (VSAM) [18] was introduced to integrate a variational approach within the self-attention framework. More concretely, it allows leveraging attention in a VAE context.

VSAM can be integrated with a Bi-LSTM by processing the encoded hidden states $\mathbf{h}_e^i = [h_{e1}^i, h_{e2}^i, \dots, h_{eT}^i]$ to produce a sequence of context vectors \mathbf{c}_t^i . Each context vector is computed as a weighted sum of all encoded hidden states with the same sequence length T . Firstly, the relevance score s_{jk} between every pair of encoded hidden states h_{ej}^i and h_{ek}^i is computed using the scaled dot-product similarity:

$$s_{jk} = \frac{(h_{ej}^i)^T h_{ek}^i}{\sqrt{d_{he}}} \quad (3)$$

where d_{he} represents the dimensionality of the encoder Bi-LSTM state. Subsequently, the attention weights a_{jk} are calculated by normalizing these relevance scores across all pairs, ensuring the sum of the weights is 1 at each timestep:

$$\mathbf{a}_t = \text{softmax}(\mathbf{s}_t) \quad (4)$$

Having computed the attention weights, we can calculate the deterministic context vectors $\mathbf{c}_{\text{det}}^i$ as if we were not using a variational approach:

$$\mathbf{v}_{\text{det}}^i = \sum_{j=1}^T a_{tj} h_{ej}^i \quad (5)$$

Then, the context vectors \mathbf{v}_t^i are modeled as random variables with a prior normal distribution. Finally, the parameters of the approximate posterior distribution $\mu_{c_t^i}, \Sigma_{c_t^i}$ are derived using the reparameterization trick [14] to allow the gradient backpropagation through the stochastic sampling process. The prior and new context vectors can be defined as:

$$p(\mathbf{v}_t^i) = \mathcal{N}(0, \mathbf{I}), \quad \mathbf{v}_t^i \sim \mathcal{N}(\mu_{\mathbf{v}_t^i}, \Sigma_{\mathbf{v}_t^i}) \quad (6)$$

The final context vectors \mathbf{v}_t^i are sampled from this approximate posterior, thereby increasing the capacity of the model to encode expressive representations into its latent space.

3 Data Augmentation for Time Series

Data augmentation involves generating synthetic data that covers unexplored input space while preserving accurate target values or labels. This process enriches the dataset, reducing the need for additional real-world data collection [21]. This technique has proven effective in domains such as computer vision, where methods like AlexNet [15] have leveraged augmented data for image classification. However, the specific characteristics of time series data, such as its temporal dependencies and intricate dynamics, pose different challenges for data augmentation. Traditional methods for computer vision or speech processing do not directly translate to the time series domain. Moreover, the augmentation methods often need to be adapted for the specific task at hand, such as classification, forecasting, or anomaly detection [21].

Transformations based on jittering, scaling, and magnitude warping have been shown to help in specific tasks, for example, time series classification [19]. They are simple to implement and increase data variety but may not capture complex patterns. Additionally, they may introduce unrealistic distortions.

Another common approach relies on pattern mixing. This method combines two or more existing time series to create new patterns that theoretically incorporate features from all the inputs. However, this technique can miss global dependencies and may lead to overfitting due to repetitive patterns. A third possibility is to use decomposition methods. They extract features from the dataset, such as trend components [2], and generate new patterns from those extracted features. The main downside of decomposition methods is that they cannot provide much variety.

Autoregressive models represent another frequently used approach. They leverage past observations to predict future values and have been used to generate synthetic time series data [12]. These methods are effective for linear and stationary time series. Still, they often struggle with complex, non-linear, and non-stationary time series data. Also, they can often be computationally expensive. Additionally, they can be used to generate synthetic data but can hardly generate semi-synthetic data (synthetic data that follows a distribution similar to a real dataset).

Finally, generative models have rapidly increased in popularity and performance over the last few years, particularly in the image and natural language processing fields. They have also recently started showing significant potential for time series data augmentation [21]. They aim to capture the underlying data distribution and generate new data samples from this learned distribution. These models provide an appealing approach to augment time series data as they can create diverse, realistic synthetic samples that capture the complex temporal dependencies characteristic of time series data. They also have limitations, including the high computation demand; the challenge of ensuring the quality and consistency of

the generated data; and the need to align the data distribution assumptions of the models with the actual underlying distribution of the original dataset.

One commonly used generative approach applied to time series data is Generative Adversarial Networks (GANs). They consist of two neural networks, a generator, and a discriminator, that compete against each other. The generator tries to produce realistic synthetic samples while the discriminator attempts to distinguish between real and synthetic data. GANs have been adapted for time series data with TimeGAN [22] and other more recent variations such as SigWGAN [17]. The model consists of four neural networks: an encoder, a decoder, a generator, and a discriminator. The encoder and decoder are used to learn the underlying temporal dynamics of the data, allowing the model to capture the sequential nature of time series. The generator and discriminator then follow the traditional GAN setup, where the generator produces synthetic time series data, and the discriminator attempts to differentiate between real and synthetic sequences. The main downside of using GANs is that they can be challenging to train due to unstable dynamics between the generator and the discriminator. At the same time, we have to rely on sampling points randomly from the latent space of the generator. This approach is inherently stochastic, meaning the specific characteristics of the generated data can be unpredictable. Additionally, the randomness in sampling does not allow for precise control over specific features of the generated data. For instance, in a time series context, generating data with particular attributes, such as following certain trends, might be desirable, which random sampling does not guarantee. Finally, because the sampling is random, the range of variations in the generated samples is bounded by the regions of the latent space that are more densely sampled. This often leads to a concentration of similar samples and, thus, a lack of diversity in the generated data.

VAEs and CVAEs have also been growing in popularity [7, 21] as generative alternatives. They map input data to a lower-dimensional latent space and then sample from this space to generate new instances. VAEs have been used for time series data, with modifications to handle temporal dependencies (inspired by variational Recurrent Neural Networks (RNNs) [4]). Even so, standard VAEs and CVAEs lack control and variety over the generated data for the same reasons we have identified for GANs. The fact that we can only random sample from the latent space does not give a mechanism to ensure diverse and controllable samples.

4 L-GTA

In this section, we introduce L-GTA, a novel method for generating semi-synthetic time series data. It integrates the capabilities of Transformers, Bidirectional Long Short-Term Memory Networks (Bi-LSTMs), Conditional Variational Autoencoders (CVAEs), and traditional time series augmentation techniques. This architecture facilitates the efficient generation of diverse and controllable time series data that retain the statistical properties of the original dataset. Our approach combines transformations and generative models [11] in a single method. Thus, we benefit from the simplicity and variety of transformation methods and the ability to generate realistic samples from generative models.

4.1 Generative Model

L-GTA is a generative model that combines Bi-LSTMs and CVAEs with a Variational Multi-Head Attention (VMHA) mechanism to generate semi-synthetic time series data. We extended VSAM to VMHA to capture more nuanced temporal dependencies and introduce a richer variational component into the architecture of the model.

The L-GTA model processes an input time series z alongside a contextual condition c , utilizing a Bi-LSTM encoder to encapsulate the short-term temporal dependencies within the data. Including VMHA within this encoder allows a more detailed analysis of the sequence, enabling the model to attend to various aspects of the temporal data simultaneously. It also makes the model capable of capturing the long-term dynamics in the data. This process can be formalized as follows:

$$(\mu_t, \Sigma_t) = \phi(\text{VMHA}(h_t, c)), \quad (7)$$

where h_t represents the hidden state of the Bi-LSTM at time t , and $\text{VMHA}(h_t, c)$ denotes the application of a Variational Multi-Head Attention mechanism, enriching the representation obtained from the Bi-LSTM according to the value of the condition c (recall that c represents context features used to condition our VAE). The function ϕ then maps this enriched representation to the parameters μ_t and Σ_t of a Gaussian distribution, from which the latent space representation v_t is sampled:

$$v_t \sim \mathcal{N}(\mu_t, \Sigma_t). \quad (8)$$

This enables the production of various plausible time series variations. The Bi-LSTM's bidirectional structure, augmented with VMHA, is particularly adept at learning complex temporal structures in the data. These are the most relevant and challenging patterns to learn in order to generate high-quality semi-synthetic instances of time series data.

The decoding phase employs another Bi-LSTM, which reconstructs the input time series \tilde{z} from the sampled latent sequence v and the contextual condition c :

$$\tilde{z}_t = \psi(\tilde{y}_t); \quad \tilde{y}_t = \text{Bi-LSTM}(v_t, \tilde{h}_{t-1}, c), \quad (9)$$

where \tilde{y}_t is the output of the decoder Bi-LSTM at time t , and ψ is a transformation function that maps \tilde{y}_t to the reconstructed time series point \tilde{z}_t , completing the generative process.

A distinctive feature of L-GTA is its ability to generate datasets with controlled variability, enabling the systematic variation of the similarity between the original and generated data. The process involves applying parametric transformations to the latent space of the model, described as $v'^i = \mathcal{T}_\eta(v^i, \eta)$. Here, \mathcal{T} represents the transformation function, and η its parameters. The transformed latent representation v'^i is then decoded to yield the transformed time series $\tilde{z}'^i = \mathcal{D}(v'^i)$, maintaining consistency with the statistical properties of the original series.

Furthermore, L-GTA supports the sequential chaining of multiple transformations in the latent space, thus enabling the generation of an extensive and diverse set of semi-synthetic time series. This is formulated as:

$$v'^i = \mathcal{T}_n(\dots \mathcal{T}_2(\mathcal{T}_1(v^i, \eta_1), \eta_2) \dots, \eta_n), \quad (10)$$

where $\{\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_n\}$ are the transformation functions, applied in sequence with their respective parameters $\{\eta_1, \eta_2, \dots, \eta_n\}$.

4.2 Time Series Transformations

The transformations applied to the latent representation are meant to affect both individual time series components and relationships between series in a dataset. However, to be realistic, they should preserve a meaningful relationship between the original and the transformed series. For that to happen, the transformation functions need to be smooth and continuous. This ensures that similar parameters produce closely related transformed series.

We employ three magnitude domain transformations — jittering, scaling, and magnitude warping [11] — to demonstrate the applicability of L-GTA. These magnitude transformations affect the values of elements while preserving the integrity of the time steps. L-GTA can be easily extended to incorporate other transformations.

Jittering is defined as the addition of a random noise component to the values of a time series. It can be defined as $v'_{r,t} = v_t^i + \epsilon_{r,t}^i$ where $\epsilon_{r,t}^i \sim \mathcal{N}(0, \sigma_v^2)$, and the standard deviation $\sigma_{r,i}$ of the added noise is a parameter of the transformation. Applying jittering to the time series via the addition of i.i.d. Gaussian noise is widely used in the literature to simulate realistic sources of measurement error or irregularity in time series data [11]. For example, consider a retail store where sales data suddenly becomes more erratic due to unexpected external factors such as local construction affecting customer traffic.

Another transformation we used was scaling, which involves modifying the amplitude of the series by a random scalar value. It is defined as $v'_{r,t} = \alpha_r^i v_t^i$ where $\alpha_r^i \sim \mathcal{N}(0, \sigma_{r,i}^2)$. Once again, $\sigma_{r,i}$ is a parameter defining the standard deviation of the multiplicative effect for each version of the dataset and series. Scaling the time series data via a multiplicative factor simulates realistic changes in the magnitude of the time series. For example, consider a retail store experiencing increased variability in sales due to various promotional campaigns executed by the store itself and its competitors.

Finally, we applied magnitude warping [11]. This method causes a smooth, continuous, nonlinear transformation of time series data. It is written as $v'_{r,t} = S_{u_k}^i(v_{r,t}^i)$ where $u_k^i \sim \mathcal{N}(1, \sigma_{r,i}^2)$. Note that $S_{u_k}^i$ interpolates a cubic spline with knots $\mathbf{u} = u_1, \dots, u_k$. Each knot u_k comes from a distribution $\mathcal{N}(1, \sigma_{r,i}^2)$, with the number of knots k and the standard deviation $\sigma_{v,i}$ as parameters. The cubic spline is fitted to the original data points, and the transformed data is obtained by scaling the original magnitude using the evaluated cubic spline function values. The proposed transformation method modulates the magnitude of the time series, maintaining its overall structure and smoothness. In retail, for example, magnitude warping could accentuate summer season sales for a store close to a beach area during unusually high temperatures.

5 Experiments

In this section, we present the experiments that were carried out to validate L-GTA. These address the following research questions:

Q1: Does L-GTA produce individual time series showing the expected pattern of the transformation applied to its latent space?

Time Series Transformations Comparison

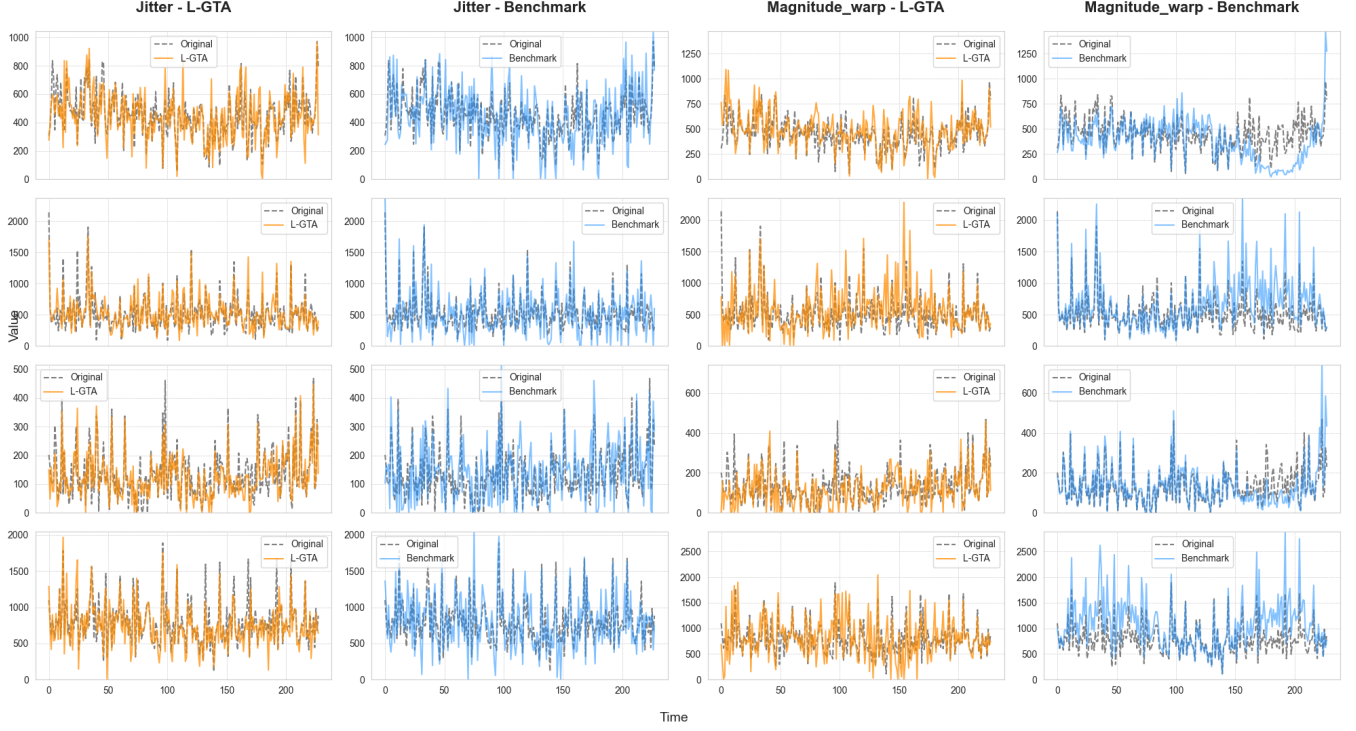


Figure 2: This figure compares the original and transformed time series of the tourism dataset over four examples (rows) for two transformation (columns) types. The first two columns display the original data alongside the data transformed using the L-GTA and direct methods using jittering. We are essentially adding random noise to the series. The last two columns show the original data compared with transformations by L-GTA and direct methods through magnitude warping. The transformation produces non-linear amplitude adjustments in the original series.

Q2: Does the L-GTA approach result in more controlled and consistent transformations of time series data compared to traditional direct methods?

Q3: How do the L-GTA and direct methods compare in preserving the original data distribution?

Q4: How well the data generated by L-GTA and direct methods preserve the predictive characteristics of the original?

Our empirical evaluation uses three public datasets: the Tourism [1] dataset from the Australian Bureau of Statistics, the M5 [16] dataset based on Walmart sales, and the Police [10] dataset from Houston police criminal reports. These datasets have various time granularities, frequencies and number of observations.

To efficiently conduct the experiments across all datasets and algorithms, we downsampled M5 [16] by reducing its frequency from daily to weekly. Additionally, for the M5 [16] and police [10] datasets, we selected a subset of 500 time series with high count levels. Preliminary experiments indicate that this selection has no significant impact on evaluating the performance of L-GTA. In the interest of space, we do not discuss them here.

5.1 Experimental Setup

Each dataset was transformed using jittering, scaling, and magnitude warping by applying both the L-GTA and traditional direct methods. The transformations were selected to mimic common variations and fluctuations that could naturally occur within time series datasets.

To address research questions **Q1** and **Q2**, our approach included a detailed visual analysis of the transformed time series data and an analysis of the residuals from the jittering transformation.

The effectiveness of the L-GTA method compared to direct methods was quantitatively assessed from three different perspectives:

- **Wasserstein distance** to measure the distributional similarity between the original and transformed datasets, providing insight into how well each method preserves the underlying data distribution and helping answer **Q3**.
- **Reconstruction error** as a percentage of the original reconstruction error, evaluating the fidelity of the transformed data in retaining the essential characteristics of the original dataset. This also contributes to answer **Q3**.
- **Prediction error** using a simple RNN model where we apply the Train-on-Synthetic, Test-on-Real (TSTR) framework, also

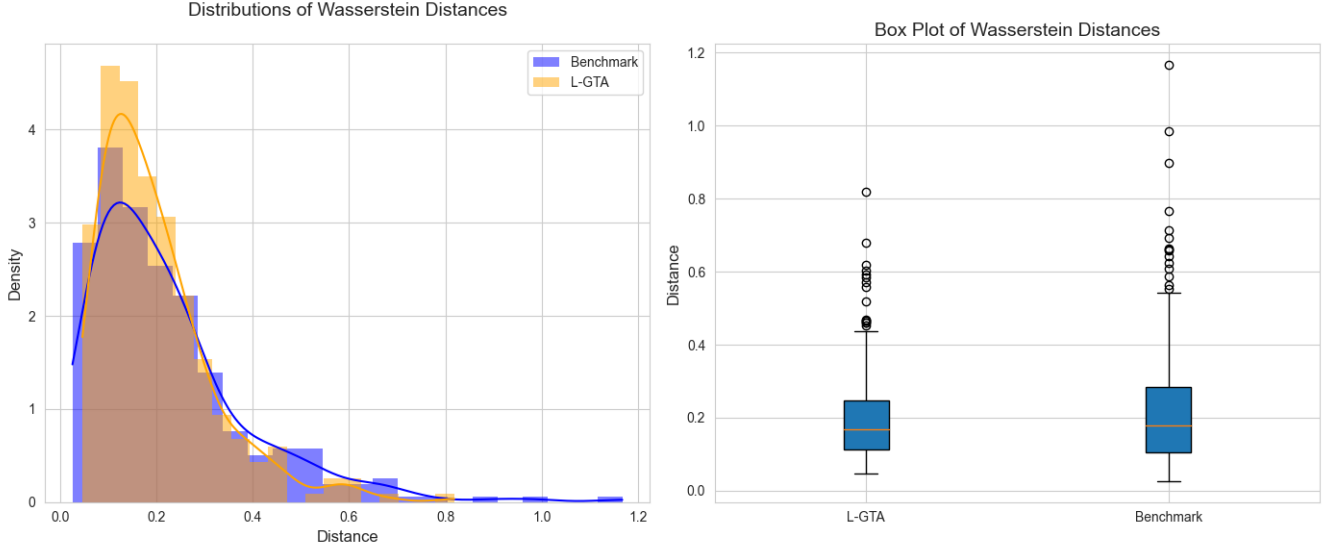


Figure 3: Comparison of Wasserstein distances for time series data transformed by L-GTA and direct methods for the tourism dataset and the magnitude warping transformation. The left panel displays density plots of the distances, highlighting the distribution of differences between transformed and original data. The right panel presents the corresponding box plots.

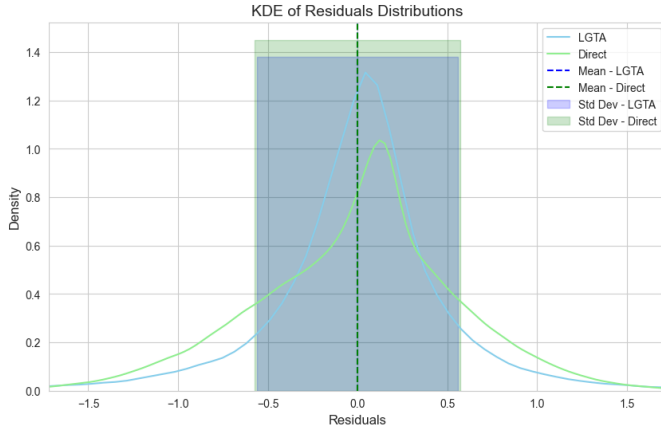


Figure 4: Kernel Density Estimation (KDE) of residuals distributions for time series data transformed by L-GTA compared to the direct method following the jittering transformation. The shaded areas denote the standard deviation, and the dashed lines mark the mean of the residuals for both methods.

used in [22, 23]. With TSTR, the basic idea is to use synthetic data generated by either L-GTA or direct methods to train the model and then test the model on a held-out set from the real data. This approach evaluates how well the generated data preserves the predictive characteristics of the original, helping us answer Q4.

The primary challenge in comparing augmentation methods lies in guaranteeing that the transformations applied by different

techniques are comparable in their impact on the dataset. To ensure a proper comparison, we tuned the parameters of the latent space transformation of L-GTA and the direct method to achieve a similar average Wasserstein distance. This calibration ensures that each method alters the data to a comparable degree. Thus, it facilitates a fair assessment of their respective strengths in preserving the essential characteristics and predictive performance of the original dataset.

For L-GTA, we used a combination of VMHA, Bi-LSTMs, and CVAEs. The encoder configuration included a transformer block integrating positional encoding and multi-head attention to capture complex temporal dependencies. The encoder also employed Bidirectional LSTMs with 256 units and dropout layers to manage overfitting. The decoder mirrored the encoder structure, using Repeat Vector and Dense layers to reconstruct the time series from the latent space.

In terms of optimization, we utilized the ADAM optimizer with a learning rate of 0.001. The loss function combined reconstruction loss, calculated as mean square error, and KL divergence for the variational component.

5.2 Results and Discussion

To answer Q1, Figure 2 compares the effects of L-GTA and the direct method on the same time series data. The plot consists of four rows, each representing a distinct dataset. For each row, the first two columns represent the addition of jittering using L-GTA and the direct method. The third and fourth columns do the same but for magnitude warping. The jittering transformation introduces random noise into the time series, emulating variations that could occur naturally within the data. The effects of jittering are evident as slight fluctuations superimposed on the original series, which

can be seen when comparing the original series and the transformation obtained with each of the L-GTA and direct methods. The magnitude warping transformation stretches or compresses the time series magnitude non-linearly. The transformation aims to mimic natural signal strength or intensity fluctuations. Once again, we observe similar patterns between L-GTA and the direct transformation.

Figure 4 offers additional insights into **Q1** and a more quantitative perspective into the jittering transformation. It showcases the distribution of residuals — the differences between the predicted and actual values — after applying L-GTA and direct transformations. Notably, the direct application of jittering results in skewed residuals, indicating a bias absent from the L-GTA method. The skewness may arise from a non-uniform application of jittering, possibly due to unaccounted variations in the amplitude, frequency, or phase of the data. Moreover, if the original data has inherent patterns, such as trends that the jittering does not adjust for, this could also introduce bias into the augmented data. Thus, it shows that in order to apply direct transformations, we need to account for several aspects to produce the expected results. This is not the case with L-GTA, where the model already learned those patterns from the original data.

Figure 2 also contributes to answer **Q2**. The first time series in the top right shows that magnitude warping can sometimes lead to extreme distortions, resulting in an erratic series that deviates significantly from the original pattern. On the other hand, L-GTA produces more consistent transformations without showing extreme values or artificial distortions.

In **Q3**, we assess the behavior of the proposed transformations by examining the distribution of distances between the transformed and the original datasets. Specifically, Figure 3 illustrates the pairwise Wasserstein distances between the original time series and those generated with L-GTA and the direct methods. This distance offers a quantitative measure to assess the discrepancy between the data distributions of transformed datasets and the original dataset. In this example for the tourism dataset, L-GTA exhibits a more concentrated distribution with a less pronounced right tail, which translates to a significantly reduced skewness. The box plot on the right reinforces this conclusion, showing a lower median, a more compact interquartile range (IQR), and fewer outliers for L-GTA. The outliers in the transformations of the direct method also corroborate the findings from Figure 2, where certain applications of magnitude warping led to erratic outcomes.

Analyzing the results presented in Table 1 allows us to address research question **Q3** for all datasets, demonstrating the consistency of L-GTA. They compare the Wasserstein distance between the original and the generated data between the L-GTA and direct methods. L-GTA achieves a lower median Wasserstein distance than the direct method across all datasets and transformations. It also exhibits consistent variability across all transformations and datasets, as indicated by similar IQR values. Also, these values are consistently smaller than the ones obtained with the direct method.

The results presented in Table 2 offer a new perspective on **Q3**. Our goal is to create a transformed time series that keeps the intrinsic information in the original data. The reconstruction error serves as a metric to evaluate how closely the transformed data retains the characteristics of the original dataset. For all transformations —

Table 1: Wasserstein distance results (median and interquartile range) for the Tourism, M5, and Police datasets across different transformations and methods.

Metric	Jitter		Magnitude Warp		Scaling	
	L-GTA	Direct	L-GTA	Direct	L-GTA	Direct
Tourism Dataset						
Median	0.123	0.181	0.116	0.182	0.111	0.135
IQR	0.110	0.059	0.098	0.194	0.106	0.113
M5 Dataset						
Median	0.108	0.194	0.108	0.411	0.109	0.310
IQR	0.108	0.063	0.107	0.405	0.109	0.321
Police Dataset						
Median	0.123	0.302	0.165	0.157	0.158	0.156
IQR	0.062	0.037	0.063	0.082	0.063	0.047

jitter, scaling, and magnitude warping — the L-GTA method consistently yields reconstruction errors close to or slightly above the original. In contrast, the direct method shows large oscillations in the reconstruction error for all transformations and across datasets. This indicates that L-GTA produces transformations without losing the key underlying information of the original dataset.

Table 2: Reconstruction error as a percentage of the original reconstruction error for the Tourism, M5, and Police datasets across different transformations.

Model	Transformation		
	Jitter	Magnitude Warp	Scaling
Tourism Dataset			
L-GTA	100.1%	102.6%	107.6%
Direct	168.7%	185.4%	244.1%
M5 Dataset			
L-GTA	101.0%	101.0%	100.9%
Direct	156.9%	180.6%	254.1%
Police Dataset			
L-GTA	99.6%	98.5%	97.0.4%
Direct	107.4%	126.9%	129.9%

Finally, another way to assess the transformed data is by performing a TSTR analysis, as shown in Table 3. It illustrates that the L-GTA method achieves an error very close to the original dataset across the three datasets. Thus, it shows that our method is capable of generating data that preserves the predictive characteristics of the original. This consistent behavior answers **Q4**. Conversely, the direct method deviates more significantly from the original performance except for the jittering transformation.

6 Conclusions and Future Work

We propose L-GTA, a novel method for generating semi-synthetic time series datasets in a controlled way. L-GTA is a transformer-based CVAE that we use to create a latent space where traditional time series augmentation methods, such as jittering and magnitude

Table 3: Prediction error using a simple RNN model on the original and transformed datasets for the Tourism, M5, and Police datasets.

Model	Transformation		
	Jitter	Magnitude Warp	Scaling
Tourism Dataset			
Original	0.022	0.022	0.022
L-GTA	0.022	0.022	0.021
Direct	0.023	0.025	0.025
M5 Dataset			
Original	0.049	0.049	0.049
L-GTA	0.048	0.049	0.046
Direct	0.051	0.064	0.069
Police Dataset			
Original	0.037	0.037	0.037
L-GTA	0.037	0.037	0.037
Direct	0.036	0.038	0.038

warping, are applied. The augmented latent representations are then decoded, producing data preserving the statistical properties and temporal dynamics of the original sequences. L-GTA enables the generation of robust and diverse transformed time series datasets.

We demonstrated that L-GTA can generate individual time series that exhibit the expected patterns resulting from transformations applied to its latent space while minimizing the production of extreme values and artificial distortions. This was evidenced from several perspectives. First, we plotted visual examples of the generated time series, highlighting the presence of extreme values in the direct method, which were absent in L-GTA. By analyzing the residuals of both transformations, we observed that the direct methods could introduce unexpected effects on the data distribution. The generally lower Wasserstein distances provided a more comprehensive perspective on the extreme values and artificial patterns across several real-world datasets. Lastly, the closer reconstruction and prediction errors relative to the originals confirmed the robustness and reliability of L-GTA.

Future research can focus on developing finer transformation controls for the proposed transformations while offering increased variety in the generated data. Another interesting pathway is exploring how to adapt L-GTA to study and evaluate the performance and robustness of algorithms to controlled data changes.

References

- [1] George Athanasopoulos and Rob Hyndman. 2006. Modeling and forecasting Australian domestic tourism. *Monash University, Department of Econometrics and Business Statistics, Monash Econometrics and Business Statistics Working Papers* 29 (01 2006). doi:10.1016/j.tourman.2007.04.009
- [2] Christoph Bergmeir, Rob J. Hyndman, and José M. Benítez. 2016. Bagging exponential smoothing methods using STL decomposition and Box–Cox transformation. *International Journal of Forecasting* 32, 2 (2016), 303–312. doi:10.1016/j.ijforecast.2015.07.002
- [3] Dawei Cheng, Fangzhou Yang, Sheng Xiang, and Jin Liu. 2022. Financial time series forecasting with multi-modality graph neural network. *Pattern Recognition* 121 (2022), 108218. doi:10.1016/j.patcog.2021.108218
- [4] Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron Courville, and Y. Bengio. 2015. A Recurrent Latent Variable Model for Sequential Data. 8.
- [5] Marzieh Fathi, Mostafa Haghi Kashani, Seyed Mahdi Jamei, and Ebrahim Mahdipour. 2022. Big Data Analytics in Weather Forecasting: A Systematic Review. *Archives of Computational Methods in Engineering* 29, 2 (Mar. 2022), 1247–1275. doi:10.1007/s11831-021-09616-4
- [6] Simone Gitto, Carmela Di Mauro, Alessandro Ancarani, and Paolo Mancuso. 2021. Forecasting national and regional level intensive care unit bed demand during COVID-19: The case of Italy. *Plos one* 16, 2 (2021), e0247726.
- [7] Maxime Goubeaud, Philipp Joußen, Nicolla Gmyrek, Farzin Ghorban, Lucas Schelkes, and Anton Kummert. 2021. Using Variational Autoencoder to augment Sparse Time series Datasets. In *2021 7th International Conference on Optimization and Applications (ICOA)*. 1–6. doi:10.1109/ICOA51614.2021.9442619
- [8] Alex Graves and Jürgen Schmidhuber. 2005. Frameworks for phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Networks* 18, 5 (2005), 602–610. doi:10.1016/j.neunet.2005.06.042 IJCNN 2005.
- [9] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Computation* 9, 8 (11 1997), 1735–1780. doi:10.1162/neco.1997.9.8.1735 arXiv:https://direct.mit.edu/neco/article-pdf/9/8/1735/813796/neco.1997.9.8.1735.pdf
- [10] Houston Police Department. 2022. Open data from the Houston Police Department for criminal reports. https://www.houstontx.gov/police/public_information.htm. Accessed: 2022-03-01.
- [11] Brian Kenji Iwana and Seiichi Uchida. 2021. An empirical survey of data augmentation for time series classification with neural networks. *PLOS ONE* 16, 7 (jul 2021), e0254841. doi:10.1371/journal.pone.0254841
- [12] Yanfei Kang, Rob Hyndman, and Feng Li. 2020. GRATIS: GeneRAting Time Series with diverse and controllable characteristics. *Statistical Analysis and Data Mining: The ASA Data Science Journal* 13 (05 2020). doi:10.1002/sam.11461
- [13] Juan Pablo Karmy and Sebastián Maldonado. 2019. Hierarchical time series forecasting via Support Vector Regression in the European Travel Retail Industry. *Expert Systems with Applications* 137 (2019), 59–73. doi:10.1016/j.eswa.2019.06.060
- [14] Diederik P Kingma and Max Welling. 2022. Auto-Encoding Variational Bayes. arXiv:1312.6114 [stat.ML]
- [15] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. [n. d.]. ImageNet Classification with Deep Convolutional Neural Networks. In *Advances in Neural Information Processing Systems*, F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger (Eds.). Curran Associates, Inc.
- [16] Spyros Makridakis, Evangelos Spiliotis, and Vasilios Assimakopoulos. 2021. The M5 competition: Background, organization, and implementation. *International Journal of Forecasting* (2021). doi:10.1016/j.ijforecast.2021.07.007
- [17] Hao Ni, Lukasz Szpruch, Marc Sabate-Vidales, Baoren Xiao, Magnus Wiese, and Shujian Liao. 2021. Sig-Wasserstein GANs for Time Series Generation. arXiv:2111.01207 [cs.LG]
- [18] João Pereira and Margarida Silveira. 2018. Unsupervised Anomaly Detection in Energy Time Series Data Using Variational Recurrent Autoencoders with Attention. In *2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)*. 1275–1282. doi:10.1109/ICMLA.2018.00207
- [19] Khandakar M. Rashid and Joseph Louis. 2019. Time-Warping: A Time Series Data Augmentation of IMU Data for Construction Equipment Activity Identification. In *Proceedings of the 36th International Symposium on Automation and Robotics in Construction (ISARC)*, Mohamed Al-Hussein (Ed.). International Association for Automation and Robotics in Construction (IAARC), Banff, Canada, 651–657. doi:10.22260/ISARC2019/0087
- [20] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. 2015. Learning Structured Output Representation using Deep Conditional Generative Models. In *Advances in Neural Information Processing Systems*, C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (Eds.), Vol. 28. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2015/file/8d55a249e6baa5c06772297520da2051-Paper.pdf
- [21] Qingsong Wen, Liang Sun, Fan Yang, Xiaomin Song, Jingkun Gao, Xue Wang, and Huan Xu. 2021. Time Series Data Augmentation for Deep Learning: A Survey. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence (IJCAI-2021)*. International Joint Conferences on Artificial Intelligence Organization. doi:10.24963/ijcai.2021/631
- [22] Jinsung Yoon, Daniel Jarrett, and Mihaela van der Schaar. 2019. Time-series Generative Adversarial Networks. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2019/file/c9efe5f26cd17ba6216bbe2a7d26d490-Paper.pdf
- [23] Jinsung Yoon, James Jordon, and Mihaela van der Schaar. 2019. PATE-GAN: Generating Synthetic Data with Differential Privacy Guarantees. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=S1zk9IRqF7>