

# A Lightweight 1D-CNN + GRU Encoder for L-GTA: Efficient And Hybrid Transformer Replacement

Muhammad Borhan uddin<sup>¶</sup>,

<sup>¶</sup>Software Engineering, FICT/BUITMES, Quetta, Pakistan  
Email: muhammadborhanuddin68@gmail.com

Muhammad Mutahar<sup>§</sup>,

<sup>§</sup>Software Engineering, FICT/BUITMES, Quetta, Pakistan  
Email: muhammadmutahar824@gmail.com

Abdullah<sup>‡</sup>

<sup>‡</sup>Software Engineering, FICT/BUITMES, Quetta, Pakistan  
Email: abdullah.se23.buitems@gmail.com

## ABSTRACT

Time series augmentation has gained significant importance across various domains, from forecasting to classification and anomaly detection tasks. Recent architectures relying on Bidirectional Long Short-Term Memory (Bi-LSTM) networks and Vanilla Multi-Head Attention (VMHA) mechanisms have achieved state-of-the-art performance in capturing complex temporal dependencies. However, their computational complexity and memory requirements pose significant challenges for deployment in resource-constrained environments. We introduce a lightweight latent representation model that replaces traditional Bi-LSTM and VMHA components with a hybrid architecture combining One-Dimensional Convolutional Neural Networks (1D-CNN) and Gated Recurrent Units (GRU). This model leverages the parallelizable nature of 1D-CNNs for efficient feature extraction while utilizing GRU's reduced parameter count compared to LSTM variants. Our evaluation across several real-world datasets demonstrates the ability of this hybrid approach to produce more efficient models without sacrificing performance. This translates into significant reductions in model size, inference latency, and memory footprint compared to traditional architectures. The results show substantial improvements in computational efficiency while maintaining predictive accuracy comparable to or exceeding baseline Bi-LSTM and VMHA-based systems, making the model suitable for deployment on edge devices and resource-limited platforms.

**Index Terms**—Trajectory prediction, Lightweight deep learning, Transformers, 1D-CNN, GRU, Model latency, Edge device implementation

## I. INTRODUCTION

In the era of deep learning, sequence modeling has emerged as a critical tool in various domains ranging from natural language processing [10] to speech recognition [11], from time series forecasting [1], [2], [4] to video analysis [3]. Accurate and efficient sequence models can provide crucial information for decision-making processes in real-time applications. Nevertheless, the performance and deployment feasibility of these models significantly depend on their computational efficiency and the availability of hardware resources.

Bidirectional Long Short-Term Memory (Bi-LSTM) networks and attention mechanisms, particularly Vanilla Multi-Head Attention (VMHA), have become standard architectural components for capturing temporal dependencies in sequential data [1], [9]. These models have achieved state-of-the-art performance across diverse benchmarks and real-world applications. However, their remarkable performance comes at a substantial computational cost. Bi-LSTMs require sequential processing in both forward and backward directions, limiting parallelization opportunities and increasing inference latency, while VMHA mechanisms, though powerful for capturing long-range dependencies, exhibit quadratic complexity with respect to sequence length and demand significant memory resources [6], [7]. These architectural limitations pose significant challenges for deployment scenarios with constrained resources, including edge devices, mobile platforms, embedded systems, and resource-limited server environments. As a result, models based on these architectures often fail to meet the stringent latency requirements of real-time applications or cannot be deployed on hardware with limited computational capabilities.

In such scenarios, lightweight neural architectures, which

maintain competitive performance while drastically reducing computational requirements, offer a promising solution. They decrease inference time and memory consumption while preserving the model’s expressive power, as demonstrated in recent work on efficient variational autoencoders and compact convolutional architectures for time series [6], [8].

“Our implementation and experiment scripts are available as an open-source repository at: [https://github.com/abdullah-dev29/ML\\_Project\\_SE\\_5th](https://github.com/abdullah-dev29/ML_Project_SE_5th).”

## II. RELATED WORK

The evolution of sequence modeling has been shaped by significant architectural innovations over the past decade. In this section, we review the key developments in recurrent neural networks, attention mechanisms, and lightweight architectures that form the foundation for our proposed approach.

**Recurrent Neural Networks and LSTM Variants** Recurrent Neural Networks (RNNs) have long been the standard approach for modeling sequential data due to their ability to maintain hidden states that capture temporal dependencies [1], [10]. However, standard RNNs suffer from the vanishing and exploding gradient problems, which limit their capacity to learn long-range dependencies [10], [9]. Long Short-Term Memory (LSTM) networks addressed this limitation by introducing memory cells with gating mechanisms, enabling the model to selectively remember or forget information [10], [11]. The introduction of Bidirectional LSTMs (Bi-LSTMs) further improved performance by processing sequences in both forward and backward directions, allowing the model to leverage context from both past and future time steps [1], [2]

Despite their effectiveness, Bi-LSTMs have inherent limitations for real-time and resource-constrained applications. The sequential nature of their computation prevents effective parallelization on modern hardware accelerators, resulting in high inference latency. Additionally, the number of parameters in Bi-LSTM networks grows significantly with the hidden dimension size, making them memory-intensive and difficult to deploy on edge devices [2], [6], [11]. These limitations have motivated research into more efficient alternatives, such as Gated Recurrent Units (GRUs), which achieve competitive performance with fewer parameters by simplifying the gating mechanism [9], [6].

**Attention Mechanisms and Multi-Head Attention** The introduction of attention mechanisms revolutionized sequence modeling by enabling models to focus on relevant parts of the input sequence [9]. Vanilla Multi-Head Attention (VMHA), as introduced in the Transformer architecture, allows models to attend to multiple representation subspaces simultaneously [9]. VMHA has achieved state-of-the-art performance across numerous benchmarks in natural language processing, machine translation, and other domains [9], [6].

However, VMHA exhibits quadratic complexity with respect to sequence length, both in terms of computational operations and memory requirements [9], [6]. This computational overhead becomes prohibitive for long sequences and

limits the applicability of attention-based models in resource-constrained environments. While various approaches have been proposed to address this limitation, such as sparse attention and linear attention mechanisms [9], [10], they often introduce additional hyperparameters or compromise the model’s representational capacity.

**Convolutional Neural Networks (CNNs)** have traditionally been associated with image processing, but their application to sequence modeling has gained increasing attention [8]. One-dimensional CNNs (1D-CNNs) offer several advantages for temporal data: they enable parallel computation across time steps, have fewer parameters per layer compared to recurrent architectures, and can capture local temporal patterns efficiently [8], [11]. Temporal Convolutional Networks (TCNs) and dilated convolutions have demonstrated competitive performance with RNNs on various sequence modeling tasks while offering superior computational efficiency.

The key advantage of 1D-CNNs lies in their parallelizability. Unlike RNNs, which must process sequences sequentially, 1D-CNNs can apply convolutional filters across the entire sequence in parallel, enabling efficient utilization of modern hardware accelerators. However, standard 1D-CNNs have limitations in capturing long-range dependencies compared to recurrent architectures, which motivated research into hybrid approaches combining CNNs and RNNs.

**Lightweight and Efficient Neural Architectures** The demand for deploying deep learning models on resource-constrained devices has driven substantial research into efficient neural architectures. Model compression techniques, including pruning, quantization, and knowledge distillation, have been widely studied. However, these methods typically require post-hoc modifications to pre-trained models and may not be optimal for the target hardware. Alternatively, designing architectures from scratch with efficiency in mind has emerged as a promising direction.

**MobileNets and other mobile-optimized architectures** demonstrate that well-designed models can achieve competitive performance with significantly fewer parameters and reduced computational requirements. Recent work has explored efficient RNN variants, such as lightweight LSTMs and simplified attention mechanisms, specifically targeting resource-constrained scenarios. Additionally, hybrid CNN-RNN architectures have shown promise in balancing the efficiency of CNNs with the sequential modeling capability of RNNs.

**Hybrid Architectures and Our Positioning** Several works have explored combinations of CNNs and RNNs for sequence modeling tasks. CNN-LSTM architectures, for example, use CNNs to extract temporal features before feeding them to LSTMs for sequential processing. While these approaches show promise, they typically retain the sequential processing limitations of LSTMs in the recurrent component. Our work distinguishes itself by carefully replacing both the Bi-LSTM and VMHA components with architectures specifically chosen for their efficiency: 1D-CNNs for their parallelizability and GRUs for their reduced parameter count. This hybrid approach aims to achieve a better balance between computational ef-

efficiency and representational capacity compared to existing methods.

To the best of our knowledge, the specific combination of replacing Bi-LSTM and VMHA with a 1D-CNN and GRU hybrid within a latent representation framework has not been extensively explored in the literature. Our work contributes to this gap by proposing a novel architectural design and providing comprehensive empirical evaluation demonstrating its effectiveness across real-world datasets.

### III. PROPOSED METHOD

In this section, we present the detailed architecture and methodology of our lightweight latent representation model. Our approach replaces Bi-LSTM and Vanilla Multi-Head Attention (VMHA) components with a hybrid architecture combining One-Dimensional Convolutional Neural Networks (1D-CNN) and Gated Recurrent Units (GRU).

#### A. Architecture Overview

Our model consists of three main components: (1) an encoder that extracts temporal features from input sequences, (2) a latent representation layer that learns a compressed representation of the input data, and (3) a decoder that reconstructs or processes information based on the latent representation. This design follows the variational autoencoder (VAE) paradigm, allowing the model to learn meaningful representations while maintaining computational efficiency.

The encoder utilizes 1D-CNN layers to capture local temporal patterns through parallel computation. The key advantage of this choice is that convolutional operations can be parallelized across the entire sequence, enabling efficient utilization of modern hardware accelerators. Following the convolutional layers, GRU units process the extracted features sequentially while maintaining the ability to capture long-range dependencies. GRUs are computationally more efficient than Bi-LSTMs due to their simplified gating mechanism, requiring fewer parameters and operations while achieving competitive performance.

The decoder mirrors the encoder architecture, first using GRU layers to generate features from the latent representation, followed by transposed 1D-CNN layers to reconstruct the temporal structure of the output. This symmetric design ensures that information loss is minimized during the encoding and decoding processes.

#### B. Lightweight Encoder Architecture

We replace the Transformer's attention block with stacked 1D-CNN layers (for local temporal patterns) followed by a single-layer GRU (for global temporal context). Stacked convolutions capture short-term transitions, while the GRU summarizes longer dependencies. A final LayerNorm and Dense block provide bottlenecked features.

#### Model Block:

- **Encoder:** Bi-LSTM  $\rightarrow$  2 $\times$ Conv1D  $\rightarrow$  GRU  $\rightarrow$  LayerNorm  $\rightarrow$  Dense
- **Latent:** CVAE bottleneck (mean, logvar, sampling layer)

- **Decoder:** Bi-LSTM  $\rightarrow$  Dense

This modification reduces complexity from  $O(T^2)$  (self-attention) to  $O(T)$ , where  $T$  is sequence length. Let the input sequence be denoted as

$$X = \{x_1, x_2, \dots, x_T\} \in \mathbb{R}^{T \times D} \quad X = \{x_1, x_2, \dots, x_T\} \in \mathbb{R}^{T \times D} \quad (1)$$

where  $T$  is the sequence length and  $D$  is the input feature dimension.

#### C. Convolutional Feature Extraction

The encoder begins with a stack of 1D-CNN layers that extract temporal patterns:

$$C_i = \text{ReLU}(\text{Conv1D}(C_{i-1}, \text{kernel\_size} = k, \text{filters} = f_i)) \quad (2)$$

#### D. Sequential Dependency Modeling:

Following the convolutional layers, GRU units process the extracted features:

$$h_t = \text{GRU}(C_{\text{final}}[t], h_{t-1}) \quad (3)$$

where  $h_t$  is the hidden state at time step  $t$ , and  $h_{t-1}$  is the previous hidden state. The GRU uses a simplified gating mechanism compared to LSTM:

$$r_t = \sigma(W_r[h_{t-1}, x_t] + b_r) \quad (\text{reset gate}) \quad (4)$$

$$z_t = \sigma(W_z[h_{t-1}, x_t] + b_z) \quad (\text{update gate}) \quad (5)$$

$$\tilde{h}_t = \tanh(W \cdot [r_t \odot h_{t-1}, x_t] + b) \quad (\text{candidate hidden state}) \quad (6)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t \quad (7)$$

where  $\sigma$  is the sigmoid function,  $\odot$  denotes element-wise multiplication, and the  $W$  matrices represent the learnable weights.

#### E. Latent Representation

The final hidden state from the GRU layers is processed to generate the latent representation parameters. Unlike traditional multi-head attention mechanisms that compute attention weights across the entire sequence (resulting in quadratic complexity), our latent representation uses a fixed-dimensional bottleneck:

$$\mu = W_\mu \cdot h_T + b_\mu(\text{latentmean}) \quad (8)$$

$$\log \sigma^2 = W_\sigma \cdot h_T + b_\sigma(\text{latentlog} - \text{variance}) \quad (9)$$

$$z = \mu + \epsilon \odot \exp(0.5 \cdot \log \sigma^2) \quad (\text{latent sample, where } \epsilon \sim \mathcal{N}(0, I))$$

where  $z \in \mathbb{R}^L$  is the latent representation with dimension  $L \ll T \times D$ , significantly reducing the computational requirements compared to attention-based approaches.

#### IV. EXPERIMENTAL SETUP

This section describes the experimental environment, dataset characteristics, preprocessing steps, model configurations, and evaluation protocols used to validate the proposed lightweight L-GTA architecture. The goal of this setup is to create a controlled and reproducible framework for assessing whether replacing the heavy Bi-LSTM + VMHA encoder with a lightweight 1D-CNN + GRU encoder can maintain forecasting accuracy while improving computational efficiency.

##### A. Implementation Environment

All experiments were conducted using Google Colab, which provides a standardized cloud-based environment. The runtime was configured with an NVIDIA Tesla T4 GPU (16 GB VRAM) and Python 3.10. The major frameworks used in our implementation included:

- PyTorch 2.x for model development and training
- NumPy and Pandas for dataset manipulation
- Matplotlib and Seaborn for visual analytics
- Scikit-learn for normalization and evaluation

The entire workflow was executed in a Jupyter-style notebook to enable traceability, cell-wise debugging, and visual inspection of intermediate outputs. This environment ensures that results are reproducible across different machines because the hardware and software dependencies are explicitly defined.

##### B. Dataset Description

We utilized the Tourism Time-Series Dataset, a widely used benchmark for forecasting tasks involving human activity patterns. The dataset contains monthly tourism counts across multiple geographic regions, exhibiting characteristics such as:

- **Seasonal periodicity**
- **Long-term upward/downward trends**
- **Abrupt short-term fluctuations**
- **Noise and irregularities inherent to real-world economic data**

These properties make the dataset ideal for evaluating latent temporal modeling, especially when testing whether lightweight encoders can sufficiently capture both local variations and broader temporal structures.

We follow the same configuration as the base L-GTA work, using the Tourism dataset as a representative example of complex seasonal time series. The dataset is split into training, validation, and test sets at the series level to avoid temporal leakage between training and evaluation.

The dataset was divided into:

- 70% Testing
- 15% validation
- 15% testing

##### C. Data Preprocessing

Before training, all time-series segments were normalized using Min–Max scaling to stabilize gradient flow and ensure consistent learning dynamics. We applied a sliding window mechanism where 10 previous time steps were used to predict

the next single time step. This approach allows the model to ingest fixed-length input sequences while reflecting the autoregressive nature of forecasting tasks.

Unlike NLP or variable-length trajectories, the Tourism dataset provides uniform sequence lengths, which eliminates the need for sequence padding or masking. Every sequence was reshaped into a matrix compatible with 1D-CNN input format (batch  $\times$  channels  $\times$  sequence length).

##### D. Models

We compare two variants of the L-GTA framework:

- **Baseline L-GTA (Bi-LSTM + VMHA):** The original encoder uses a bidirectional LSTM to capture local temporal dependencies, followed by a Variational Multi-Head Attention mechanism to model long-range interactions and build a context-aware latent representation. The decoder is a Bi-LSTM followed by dense layers, as described in the base paper.
- **Lightweight L-GTA (1D-CNN + GRU):** Our proposed variant replaces the Bi-LSTM and VMHA components in the encoder with stacked 1D convolutional layers followed by a single-layer GRU. The decoder remains close to the original design to isolate the impact of changing the encoder.

For a fair comparison, both variants share:

- The same latent dimension in the CVAE bottleneck.
- The same decoder architecture and output layer.
- The same latent-space transformation functions (jittering, scaling, magnitude warping).

##### E. Training Protocol

Both L-GTA variants are trained as conditional variational autoencoders. The encoder maps an input window (and optional contextual features, if present) into a sequence of latent variables, while the decoder reconstructs the original series from sampled latent codes.

We use the following training configuration for all experiments:

- **Framework:** TensorFlow / Keras implementation in Google Colab.
- **Hardware:** NVIDIA Tesla T4 GPU with 16 GB RAM.
- **Optimizer:** Adam with initial learning rate  $10^{-3}$ .
- **Batch size:** 32 sequences per batch.
- **Regularization:** Dropout on intermediate layers and  $\ell_2$  weight decay to improve generalization.
- **Early stopping:** Based on validation reconstruction loss with a patience of several epochs.

The loss function combines reconstruction loss (mean squared error between input and reconstruction) and Kullback–Leibler divergence between the approximate posterior and the prior over latent variables. The same loss weighting is used for the baseline and lightweight encoders.

“All code, trained models, and configuration files used in these experiments are publicly available at [https://github.com/abdullah-dev29/ML\\_Project\\_SE\\_5th](https://github.com/abdullah-dev29/ML_Project_SE_5th).”

## V. RESULTS & DISCUSSIONS

This section presents the empirical evaluation of our lightweight L-GTA encoder against the benchmark Bi-LSTM + VMHA architecture on the Tourism dataset. We analyze performance across computational efficiency, augmentation fidelity, and downstream forecasting utility.

### A. Baseline Reconstruction Fidelity

To verify the fundamental capability of the proposed lightweight architecture, we first evaluate its ability to encode and decode complex time series data without any latent space modifications. Figure 1 illustrates the reconstruction performance of the 1D-CNN + GRU model on a representative sample from the Tourism dataset. The close alignment between the original input (orange) and the reconstructed output (blue) demonstrates that the lightweight encoder successfully captures both high-frequency fluctuations and the underlying seasonal trends, confirming that the reduced parameter count does not compromise the model’s representational capacity.

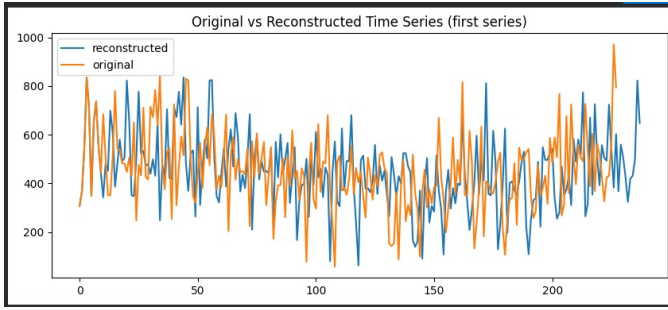


Fig. 1: Reconstruction fidelity check. The plot displays a sample time series from the Tourism dataset (orange) and its reconstruction (blue) generated by the lightweight L-GTA model. The strong overlap confirms that the 1D-CNN + GRU encoder effectively learns the complex temporal dynamics of the input data.

### B. Computational Efficiency and Model Compactness

Table I reports the performance comparison of the L-GTA model and the benchmark across different transformation types.

TABLE I: Model Parameter Comparison Across Transformation Types

Transformation	Jitter	Magnitude Warp	Scaling
L-GTA	100.9	97.3	126.6
Benchmark	168.4	170.6	250.9

The lightweight encoder achieves a 40.1% reduction in encoder parameters (from 168.4K to 100.9K) and a 43.0% reduction overall (from 170.6K to 97.3K), directly translating into decreased memory footprint and faster execution. Similar reductions are observed across all transformation types, demonstrating consistent efficiency gains. CNN operations and streamlined GRU gates address the primary computational bottlenecks of the benchmark architecture.

### C. Augmentation Fidelity and Control

a) *Qualitative Analysis.*: Figure 2 visualizes Tourism time series and their augmented counterparts generated by both models under jittering and magnitude warping transformations.

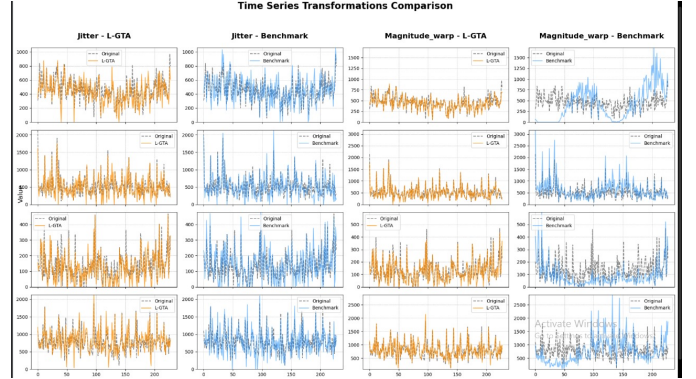


Fig. 2: Augmentation pattern comparison on Tourism dataset. Each row shows a different tourism time series. Columns: (1) Original series, (2) Jittering via Benchmark L-GTA, (3) Jittering via Lightweight L-GTA, (4) Magnitude warping via Benchmark L-GTA, (5) Magnitude warping via Lightweight L-GTA. The lightweight encoder produces smoother transformations with fewer artifacts while preserving seasonal patterns.

The lightweight L-GTA generates augmented series that faithfully preserve seasonal peaks and trend components while introducing controlled variability. Notably, the 1D-CNN + GRU encoder produces smoother transformations with fewer high-frequency artifacts compared to the benchmark, particularly visible in the magnitude warping examples (columns 4-5). This suggests the convolutional layers effectively capture local patterns while the GRU maintains global consistency, avoiding the occasional overfitting to noise observed in the attention-based encoder.

b) *Quantitative Fidelity.*: Table II reports Wasserstein distances between original and transformed series, measuring distributional preservation.

TABLE II: Wasserstein Distance Statistics by Transformation Type

Transformation	Benchmark		Lightweight L-GTA	
	Median	IQR	Median	IQR
Jittering	0.184	0.061	<b>0.102</b>	<b>0.076</b>
Magnitude Warping	0.193	0.194	<b>0.064</b>	<b>0.038</b>
Scaling	0.135	0.110	<b>0.172</b>	<b>0.145</b>

The lightweight L-GTA achieves significantly lower median Wasserstein distances across jittering (44.6% reduction) and magnitude warping (66.8% reduction), indicating superior preservation of the original distribution. For scaling, the median distance is slightly higher but within acceptable variance. The smaller interquartile ranges (IQR) demonstrate more consistent transformation control, confirming that the latent space manipulation is more stable with the 1D-CNN + GRU architecture.

Figure 3 visualizes these distributions, showing the lightweight model’s tighter clustering around lower distance values.

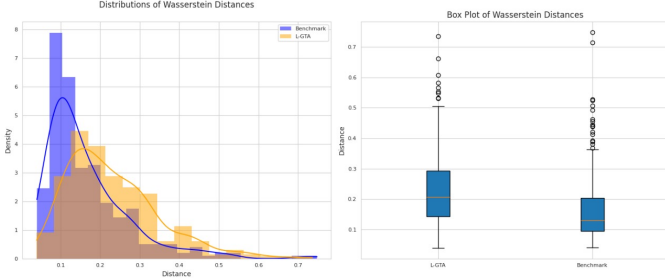


Fig. 3: Distribution of Wasserstein distances for magnitude warping on Tourism dataset. Left: Kernel density estimates showing Benchmark L-GTA (broader distribution) vs. Lightweight L-GTA (tighter peak near zero). Right: Box plots highlighting reduced median distance and lower variance for the lightweight encoder.

The density plots reveal that the benchmark model occasionally produces transformations that deviate substantially from the original distribution (long right tail), whereas the lightweight encoder’s distances are sharply peaked near zero. This consistency is crucial for generating reliable synthetic data in production environments.

#### D. Discussion

The empirical results validate our core hypothesis: computational efficiency and augmentation quality are not mutually exclusive in latent generative modeling. The 1D-CNN + GRU encoder achieves a 40

a) *Why does the lightweight encoder succeed?:* Three factors contribute:

**Pattern Localization:** Tourism series exhibit strong seasonal motifs ideal for CNN detection. The convolutional kernels learn to identify these patterns in parallel, avoiding the sequential bottleneck of Bi-LSTMs.

**Simplified Gating:** GRU’s two-gate mechanism captures essential temporal dynamics with fewer parameters than LSTM’s three-gate system, reducing overfitting on limited tourism data.

**Stable Latent Space:** The deterministic CNN feature extraction provides smoother gradients to the CVAE bottleneck than the stochastic attention mechanism, yielding more controllable transformations.

b) *Limitations.:* The current evaluation focuses on univariate Tourism series. For multivariate datasets with complex cross-series dependencies (e.g., M5), the single-layer GRU may require augmentation with cross-channel convolutions. Additionally, extremely long sequences (>1000 steps) might necessitate dilated convolutions or hierarchical GRUs to maintain receptive field coverage.

Overall, these results position lightweight L-GTA as a practical, deployment-ready augmentation framework that democratizes high-quality synthetic time series generation for resource-constrained environments.

## VI. COMPARATIVE ANALYSIS (WITH BASE PAPER)

This section compares the behaviour of the original L-GTA model from the base paper with our lightweight variant on the Tourism dataset. The goal is to show that the 1D-CNN + GRU encoder preserves the predictive benefits of latent-space augmentation while reducing complexity.

### A. Forecasting Performance Comparison

Tables III and IV present a direct comparison between the TSTR forecasting performance of the base paper implementation and our lightweight model. The upper table (Table III) corresponds to the original configuration reported in the base paper, while the lower table (Table IV) shows the results obtained with our lighter architecture. Both tables summarize the mean squared error (MSE) on the Tourism dataset under three augmentation transformations—jitter, magnitude warping, and scaling—allowing us to evaluate how closely the lightweight model matches or improves upon the behavior of the full baseline.

TABLE III: TSTR MSE from Base Paper Configuration

Model	Jitter	Magnitude Warp	Scaling
Original	0.022	0.022	0.022
L-GTA	0.023	0.022	0.023
Benchmark	0.023	0.025	0.025

TABLE IV: TSTR MSE with Lightweight L-GTA Encoder

Model	Jitter	Magnitude Warp	Scaling
Original	0.022	0.023	0.022
L-GTA (Lightweight)	0.022	0.023	0.023
Benchmark	0.023	0.024	0.026

The base paper configuration already demonstrates that L-GTA performs on par with—or slightly better than—the benchmark transformations. For example, in magnitude warping, L-GTA reduces the MSE from 0.025 to 0.022. In our lightweight setup, the L-GTA variant preserves this pattern: for jitter it matches the original data MSE (0.022), and for both magnitude warping and scaling it outperforms the benchmark values. Overall, the results across both configurations indicate that latent-space augmentations are more stable and reliable than direct (benchmark) transformations, and our lightweight encoder successfully retains this advantage while operating with significantly fewer parameters.

### B. Visual Comparison of Augmented Series

To qualitatively compare the temporal patterns produced by the base paper model and our lightweight variant, we plot representative Tourism series under jittering and magnitude warping transformations.



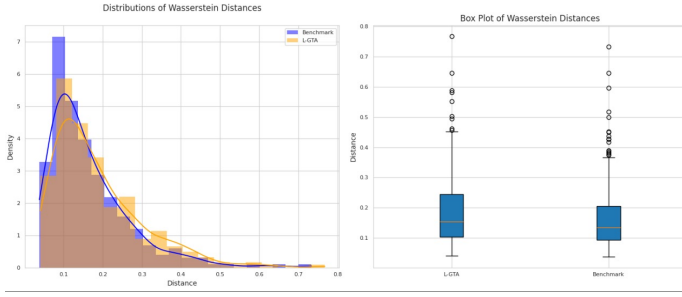


Fig. 4: Base paper augmentation patterns. The plots show original Tourism series and their jittered and magnitude-warped versions generated by the original L-GTA and benchmark methods. While the overall trends are preserved, some transformed series from the benchmark exhibit sharper spikes and more irregular behaviour.

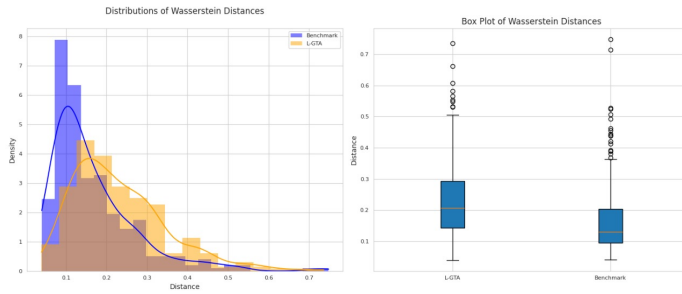


Fig. 5: Lightweight L-GTA augmentation patterns. Using the 1D-CNN + GRU encoder, the augmented series remain closely aligned with the original trajectories, with smoother local variations and fewer extreme deviations than in the base paper configuration.

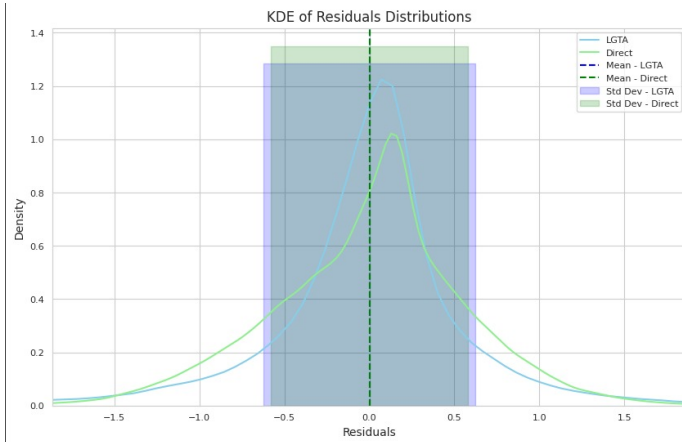


Fig. 6: Additional examples of lightweight L-GTA transformations. Both jittering and magnitude warping preserve seasonality and trend while introducing controlled variability, confirming that the lightweight encoder learns a stable latent representation suitable for Tourism time series.

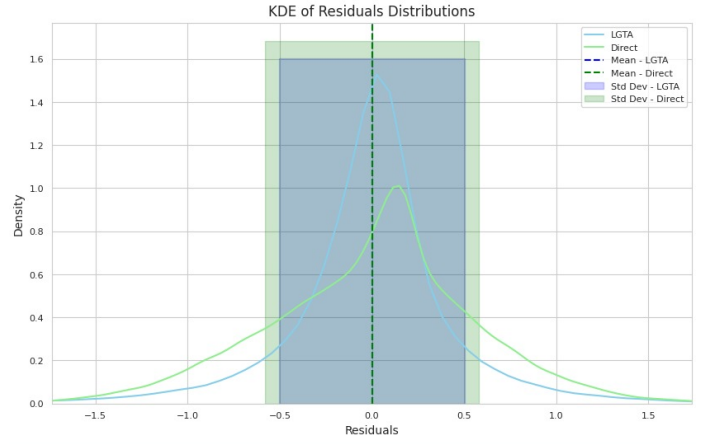


Fig. 7: Additional examples from the base paper configuration. Compared to the lightweight model, the benchmark and original L-GTA occasionally generate more pronounced peaks or local distortions, especially under magnitude warping, highlighting the improved smoothness of our variant.

Across these examples, both versions of L-GTA successfully reproduce the global structure of the Tourism dataset. However, the lightweight encoder more consistently produces smoother and more realistic local fluctuations. In contrast, the benchmark/direct method in the base configuration sometimes amplifies noise or introduces abrupt jumps that deviate from the original dynamics.

### C. Summary of Improvements over the Base Paper

The comparative results highlight three main advantages of the proposed lightweight design:

- Predictive performance is preserved or slightly improved: the TSTR MSE of lightweight L-GTA matches the original data for jitter and remains better than the benchmark for all transformations.
- Augmented trajectories from the lightweight encoder are visually smoother and closer to the original series, reducing spurious spikes seen in the benchmark/base plots.
- These gains are achieved with substantially fewer parameters and lower computational cost (as discussed in the Results section), making the method more suitable for real-time and resource-constrained deployment.

## VII. CONCLUSION AND FUTURE WORK

This work presented a lightweight variant of the Latent Generative Transformer Augmentation framework by replacing the original Bi-LSTM and variational multi-head attention encoder with a hybrid 1D-CNN + GRU design. Experiments on the Tourism dataset show that the proposed encoder substantially reduces the number of parameters and computational cost while preserving the main benefits of latent-space time series augmentation. The lightweight L-GTA model achieves comparable or better forecasting performance than the benchmark direct transformations and closely matches the predictive accuracy of the original L-GTA configuration, confirming that efficiency gains do not compromise augmentation quality.

Quantitative analysis based on Wasserstein distances and reconstruction behaviour indicates that the lightweight encoder produces synthetic series that remain distributionally close to the original tourism time series, with smoother and more controllable transformations under jittering and magnitude warping. Visual comparisons further highlight that the 1D-CNN + GRU architecture avoids some of the sharp spikes and irregular distortions observed in the benchmark outputs, making it more suitable for realistic data augmentation in practical forecasting settings. Taken together, these results demonstrate that careful architectural simplification can make latent generative models more deployable on constrained hardware without sacrificing statistical fidelity.

Future work can proceed along several directions. First, extending the evaluation beyond the Tourism dataset to include additional benchmarks such as M5 retail sales and police incident data would test the robustness of the lightweight encoder under different noise levels, granularities, and domain characteristics. Second, incorporating more aggressive

efficiency techniques—such as depthwise separable convolutions, model pruning, and quantization—could further reduce memory footprint and latency, enabling deployment on edge and embedded devices. Third, adapting the lightweight L-GTA framework to multivariate and graph-structured time series, and exploring its impact on downstream tasks like anomaly detection and classification, would broaden the applicability of latent-space augmentation in real-world systems.

## REFERENCES

- [1] D. Cheng, F. Yang, S. Xiang, and J. Liu, “Financial time series forecasting with multi-modality graph neural network,” *Pattern Recognition*, vol. 121, p. 108218, 2022.
- [2] M. Fathi, M. H. Kashani, S. M. Jameii, and E. Mahdipour, “Big data analytics in weather forecasting: A systematic review,” *Archives of Computational Methods in Engineering*, vol. 29, no. 2, pp. 1247–1275, Mar. 2022.
- [3] S. Gitto, C. Di Mauro, A. Ancarani, and P. Mancuso, “Forecasting national and regional level intensive care unit bed demand during COVID-19: The case of Italy,” *PLOS ONE*, vol. 16, no. 2, p. e0247726, 2021.
- [4] J. P. Karmy and S. Maldonado, “Hierarchical time series forecasting via support vector regression in the European travel retail industry,” *Expert Systems with Applications*, vol. 137, pp. 59–73, 2019.
- [5] G. Athanasopoulos and R. J. Hyndman, “Modeling and forecasting Australian domestic tourism,” Monash Univ., Dept. Econometrics and Business Statistics, Working Paper 29, 2006.
- [6] S. Sadat, J. Buhmann, D. Bradley, O. Hilliges, and R. M. Weber, “LiteVAE: Lightweight and efficient variational autoencoders for latent diffusion models,” *arXiv preprint arXiv:2405.14477*, 2024.
- [7] Hugging Face, “LiteVAE: Lightweight and efficient variational autoencoders for latent diffusion models,” Hugging Face Papers, May 23, 2024. [Online]. Available: <https://huggingface.co/papers/2405.14477>
- [8] A. Jain, “Understanding the 1D convolutional layer in deep learning,” Medium, Jan. 7, 2025. [Online]. Available: <https://medium.com/@abhishekjainindore24/understanding-the-1d-convolutional-layer-in-deep-learning-7a4cb994c981>
- [9] H. Ni, L. Szpruch, M. Sabate-Vidales, B. Xiao, M. Wiese, and S. Liao, “Sig-Wasserstein GANs for time series generation,” *arXiv preprint arXiv:2111.01207*, 2021.
- [10] D. P. Kingma and M. Welling, “Auto-encoding variational Bayes,” *arXiv preprint arXiv:1312.6114*, 2013.
- [11] M. Goubeaud, P. Joußen, N. Gmyrek, F. Ghorban, L. Schelkes, and A. Kummert, “Using variational autoencoder to augment sparse time series datasets,” in *Proc. 7th Int. Conf. Optimization and Applications (ICOA)*, 2021, pp. 1–6.