# INSTRUCTION-FOLLOWING LLMS FOR TIME SERIES PREDICTION: A TWO-STAGE MULTIMODAL AP PROACH

Anonymous authors

 Paper under double-blind review

#### Abstract

We introduce **Text-Informed Time Series Prediction (TITSP)**, an innovative multimodal framework that integrates textual knowledge with temporal dynamics using Large Language Models (LLMs). TITSP employs a two-stage process that bridges numerical data with rich contextual information for enhanced forecasting accuracy and interpretability. In the first stage, we present **AutoPrompter**, which captures temporal dependencies from time series data and aligns them with semantically meaningful text embeddings. In the second stage, these aligned embeddings are refined by incorporating task-specific textual instructions through LLM. We evaluate TITSP on several multimodal time series prediction tasks, demonstrating substantial improvements over state-of-the-art baselines. Quantitative results reveal significant gains in predictive performance, while qualitative analyses show that textual context enhances interpretability and actionable insights. Our findings indicate that integrating multimodal inputs not only improves prediction accuracy but also fosters more intuitive, user-centered forecasting.

#### 1 INTRODUCTION

Time series prediction is critical in fields such as finance, healthcare, and climate science, where timely and accurate forecasts drive informed decision-making. Traditional time series forecasting pipelines typically involve three key stages: data preprocessing, model selection, and performance evaluation.



Figure 1: The ideal complete process of Time-series Prediction

However, methods like ARIMA (Shumway et al., 2017), while foundational, often struggle to capture the complex non-linear patterns and long-range dependencies present in real-world datasets.
With the rise of deep learning, models such as Long Short-Term Memory (LSTM) networks (Siami-Namini et al., 2019) and Convolutional Neural Networks (CNNs) (Wang et al., 2019) have shown significant improvements in modeling these complexities. However, they remain constrained by their reliance on numerical data alone, which limits their ability to integrate external contextual information—such as expert insights or macroeconomic events—that could enhance forecast accuracy and interpretability. This shortcoming becomes especially problematic in chaotic or highly volatile systems, such as financial markets or patient health monitoring. Recent advancements in Transformer-based architectures Ilbert et al.; Zhou et al. (2021); Wen et al. (2022) and Large Language Models (LLMs) (Gruver et al., 2024; Jin et al., 2023) have enhanced the ability to capture long-range dependencies in time series data. Yet, these models also face limitations when applied in isolation, often missing out on crucial domain-specific insights provided through other modalities, such as text. In practical scenarios, domain experts may wish to provide instructions or insights that can guide the forecasting process, but existing models do not easily accommodate such interactions.

061 To address these challenges, we propose **Text-Informed Time Series Prediction** (TITSP), a novel 062 two-stage framework that combines the strengths of deep learning models for time series prediction 063 with the contextual richness of domain-specific textual inputs. TITSP first captures temporal depen-064 dencies using numerical time series data and then refines these predictions by incorporating taskspecific textual instructions through a Large Language Model. This integration allows the model 065 to generate more accurate, context-aware, and interpretable predictions. Our extensive experiments 066 demonstrate that TITSP significantly outperforms state-of-the-art models, particularly in scenarios 067 where expert input is essential. 068



Figure 2: Interactive Prediction of GBP/USD Exchange Rate during the 2008 Financial Crisis compared with other methods.

Figure 2 demonstrates the effectiveness of TITSP in predicting the GBP/USD exchange rate during the 2008 financial crisis. It highlights how expert textual inputs can enhance predictions in volatile environments. Additionally, Appendix J provides details on why a finance expert might offer insights into the logarithmic decay observed on October 11, 2008.

The remainder of the paper is structured as follows. In Section 2, we discuss related work, covering existing approaches in time series prediction and multimodal learning. Section 3 formalizes the problem statement and the key challenges we aim to address. In Section 4, we describe our proposed methodology, providing a detailed explanation of the architecture and the rationale behind its design. Section 5 presents the experimental results, showcasing the effectiveness of our approach through evaluations on multiple benchmark datasets.

The code to reproduce the results of this paper is included in the supplementary material and will be made publicly available upon acceptance.

## 2 RELATED WORKS

098

096

099

100

069

071

073

074

075

076

077

078

079 080

081

083

084

085

086

087

2.1 TIME SERIES FORECASTING APPROACHES

Time series forecasting began with classical models like ARIMA (Box et al., 2015), which are effective for linear patterns but struggle with non-linearities common in real-world data (Chatfield, 2000). Feature extraction methods like *tsfresh* (Christ et al., 2018) and machine learning algorithms such as Random Forests (Breiman, 2001) and Support Vector Machines (Vapnik, 1998) improved predictive accuracy by capturing non-linear dependencies.

The adoption of deep learning models, particularly RNNs (Rumelhart et al., 1986), LSTMs (Hochreiter & Schmidhuber, 1997), and Transformers (Vaswani et al., 2017), has further advanced the field by automatically learning complex temporal patterns (Wen et al., 2017). Following the success of

transformer, works including TimeXer (Wang et al., 2024), itransformer (Liu et al., 2023), PatchTST (Nie et al., 2022) and Informer (Zhou et al., 2021) are designed to address time-series prediction problems. Despite these advancements, most approaches rely solely on time series data. Our work
 diverges by integrating time series with textual context via Large Language Models (LLMs), offering a richer, more contextual understanding of the data.

114 2.2 MULTIMODAL LEARNING IN TIME SERIES PREDICTION

Integrating time series data with unstructured text has become a key strategy for enhancing forecast ing accuracy through multimodal approaches that leverage diverse data types for richer context.

118 Advances in Model Architectures and Joint Prediction Zhang et al. (2024) introduced the Dual-119 Adapter model to optimize time series representation by balancing textual and temporal information. 120 Building on this, Liu et al. (2024b) proposed UniTime, a unified model leveraging NLP for enhanced time series forecasting across domains, while Liu et al. (2024a) explored LLMs for dynamical sys-121 tems, showcasing their adaptability. Surveys such as (Liang et al., 2024; Deldari et al., 2022) offer 122 taxonomies of multimodal learning, identifying gaps in self-supervised methods for multimodal and 123 temporal data. The Multi-Modal Forecaster by Kim et al. (2024) jointly predicts time series and 124 text, highlighting the potential of cohesive multimodal frameworks. Similarly, Cheng et al. (2024) 125 introduced a prompt-based multimodal framework focused on classification. Recent advancements, 126 such as (Jia et al., 2024) and (Jin et al., 2023), further integrate time series and text using LLMs, 127 enhancing predictions through textual cues. 128

Existing Gaps and Novel Contributions. While recent works have advanced the integration of time series and text, a key gap remains in instruction-based time series forecasting. Our method addresses this by introducing a novel framework for text instruction-based forecasting, leveraging the interplay between textual instructions and time series data for richer, interactive predictions. Additionally, we present a tailored methodology to benchmark such tasks, enabling applications in "what-if" scenarios and dynamic, user-interactive forecasting.

134 135

136

139

140

141

142

143

144

145

146

147

148

149

150

151

152 153

154

113

2.3 CONTRIBUTIONS OF OUR WORK.

Building on recent advancements in multimodal time series forecasting, our work makes the follow-ing key contributions:

- Data Generation Workflow: We propose a novel data generation workflow that enables the evaluation of instruction-based time series forecasting. This workflow facilitates the testing of model applicability in scenarios that require adherence to specific instructions, a previously underexplored area.
- **Innovative Two-Stage Approach:** We introduce a two-stage model that effectively integrates both temporal data and textual instructions. This design captures contextual information while adhering to instruction-based predictions, leading to enhanced performance.
  - **Comprehensive Ablation Studies:** We conduct ablation studies to validate each step of the proposed approach, demonstrating the importance and rationale behind each component.
  - **Comparative Evaluation:** We benchmark our method against state-of-the-art models for multimodal time series and text forecasting, showing competitive or superior performance across key metrics (both in terms of a new metrics called compliance rate and in terms of MSE).
- These contributions address a significant gap in instruction-based forecasting, offering a structured and interpretable solution.
- 155 156 157

158

**3 PROBLEM DEFINITION** 

Let  $\mathbf{X} = \{x_1, x_2, \dots, x_T\}$  denote a multivariate time series, where  $x_t \in \mathbb{R}^d$  represents the *d*dimensional vector at time step  $t \in \{1, 2, \dots, T\}$ . Traditional time series forecasting aims to predict future values  $\hat{\mathbf{X}} = \{\hat{x}_{T+1}, \hat{x}_{T+2}, \dots, \hat{x}_{T+H}\}$ , where *H* is the forecasting horizon, using only the historical observations  $\mathbf{X}$ . In real-world applications, textual information often provides critical context or specific instructions that can influence the expected trajectory of the time series. Let  $\mathbf{S} = s_1, s_2, \ldots, s_n$  be a sequence of n tokens representing a textual description or instruction, where  $s_i \in \mathcal{V}$  and  $\mathcal{V}$  is a fixed vocabulary. This text can convey information about **future events**, **conditions**, or **hypothetical scenarios** that could affect the future behavior of the time series.

Our objective is to design a multimodal framework that integrates both the time series **X** and the text **S** to produce a more **informed** forecast  $\hat{\mathbf{X}}^{(S)} = \{\hat{x}_{T+1}^{(S)}, \hat{x}_{T+2}^{(S)}, \dots, \hat{x}_{T+H}^{(S)}\}$ , conditioned on both the historical data and the information in **S**. Formally, the forecasting function is defined as:

$$\hat{\mathbf{X}}^{(\mathbf{S})} = f(\mathbf{X}, \mathbf{S}; \theta) \tag{1}$$

where  $f(\cdot, \cdot; \theta)$  is a function parameterized by  $\theta$ , capturing the relationships between the time series and the text. The goal is to learn the parameters  $\theta$  to minimize the forecasting error:

$$\theta^* = \arg\min_{\theta} \mathcal{L}\left(\hat{\mathbf{X}}^{(\mathbf{S})}, \mathbf{X}^{\text{true}}\right)$$
(2)

where  $\mathcal{L}$  is a loss function that measures the discrepancy between the forecasted values  $\hat{\mathbf{X}}^{(S)}$  and the true future values  $\mathbf{X}^{\text{true}} = \{ x_{T+1}^{\text{true}}, x_{T+2}^{\text{true}}, \dots, x_{T+H}^{\text{true}} \}$ .

#### 4 PROPOSED METHODOLOGY

#### 4.1 DATA GENERATION

To train and evaluate our multimodal forecasting model, we generate a synthetic dataset that integrates time series data with textual instructions. The goal is to modify forecasted time series values based on text descriptions, simulating real-world scenarios where instructions influence future predictions.

- 190 The data generation process consists of the following steps:
  - Base Time Series Creation: We generate multivariate time series  $\mathbf{X} = \{x_1, x_2, \dots, x_T\}$  using standard generators (e.g., sine waves, real-world data including financial, electricity data among others). The full description of all the base time series used for the generation is provided in Appendix C.
    - **Textual Instruction Integration:** For each time series, we generate corresponding textual instructions **S**, such as *"increase," "decrease," or "stabilize,"*, ..., which provide guidance on how to modify future values.
    - Forecast Modification Based on Text: Forecasted values  $\hat{\mathbf{X}}^{(S)} = { \hat{x}_{T+1}^{(S)}, \dots, \hat{x}_{T+H}^{(S)} }$ are adjusted based on the instructions.

As example, for the instruction "trend up", the forecast is modified by applying a linear increment to the last observed value  $x_T$ :

171 172 173

174

179

180 181 182

183

185

191

192

193

194

196

197

199

200 201

$$\hat{x}_{T+h}^{(\mathbf{S})} = x_T + A \times h, \text{ where } h \in \{1, 2, \dots, H\}$$
 (3)

In this expression, A is a constant controlling the upward trend. In the data generation, A evolves with time while being positive to add complexity to the generative model, as seen in Figure 3, where some non-linear transformations may be observed for linear growth. Further details on the instructions and their corresponding mathematical modifications are provided in Appendix 6.

To demonstrate the realism of the synthetic dataset, we present several example scenarios where textual instructions are used to modify the base time series. Each scenario shows the base time series and its modified forecasted values, highlighting how the text influences the trajectory of future predictions. Figure 3 shows three different scenarios: "*Linear Growth*", "*Exponential Downward*", and "*Logarithmic Decay*". These examples illustrate the diversity of the dataset and how textual instructions are effectively incorporated into the time series predictions. More figures are listed in Appendix C.



Figure 3: Comparison of different orders for the data generation: (a) Linear Growth, (b) Exponential Decay, and (c) logarithm growth. The vertical line indicates the start of the forecast.

#### 4.2 PROPOSED ARCHITECTURE

225

226

227 228

229 230

231

232

233

234

249

253

254

256

257

258

259

260

261

262

264

265

267

268

Our proposed framework employs a two-stage approach to integrate time series data and textual information for predictive tasks. This multimodal architecture leverages both unsupervised and supervised learning paradigms to align time series data with corresponding semantic information derived from textual context, improving interpretability and predictive performance (see Figure 4 for an architectural overview).



Figure 4: Model architecture

# Stage 1: AutoPrompter: Enhancing VQ-VAE for Unsupervised Alignment of Time Series and Pre-trained Language Codebook

In the first stage, we propose a novel mechanism, AutoPrompter, which serves as a bridge, translating time series data into a compressed, semantically meaningful text embedding space. By introducing cross-attention and a pre-trained language codebook to the classical VQ-VAE (Vector Quantised-Variational AutoEncoder; (Van Den Oord et al., 2017)), AutoPrompter quantizes the time series space, allowing it to align effectively with text embeddings. This translation helps capture the underlying semantic patterns within the time series data, facilitating meaningful connections with textual information. The key components are as follows:

- Time Series Embedding: The time series data, X, is processed through a Convolutional Neural Network (CNN), producing a high-dimensional temporal embedding, z<sub>ts</sub>, that captures key temporal patterns.
- 2. Text Embedding with Pre-trained Codebook and Linear Compression: A pre-trained language model, equipped with a fixed vector quantization codebook, provides discrete text embeddings,  $z_{txt}$ . These embeddings are subsequently passed through a linear compression layer to obtain a compressed semantic representation,  $z_{txt.compressed}$ . This compression reduces the dimensionality of the text embeddings, facilitating efficient alignment with the time series embeddings while preserving essential semantic information.
- 3. Cross-Attention Alignment: A cross-attention mechanism aligns the time series embedding  $z_{ts}$  with the compressed text embedding  $z_{txt.compressed}$ . In this setup,  $z_{ts}$  acts as the

270 query, and z<sub>txt\_compressed</sub> serves as the key-value pair. This interaction facilitates the integra-271 tion of temporal and semantic information, producing an aligned multimodal representa-272 tion,  $\mathbf{z}_{aligned}$ . 4. Self-Supervised Learning: Our model is trained using self-supervised learning to align 274 time series embeddings with meaningful text embeddings while preserving key character-275 istics of the original time series data. The training minimizes two key losses: 276 1. Reconstruction Loss ( $\mathcal{L}_{recon}$ ) 2. Vector Quantization Loss ( $\mathcal{L}_{vq}$ ) 277 • Vector Quantization Loss ( $\mathcal{L}_{vq}$ ): The Vector Quantization Loss ensures that the encoder outputs are quantized by aligning them with the nearest vectors from a fixed, 279 pre-trained language codebook. This encourages the model to represent time series data using a discrete set of embeddings from the codebook, facilitating alignment be-281 tween the input data and text-based embeddings. Unlike traditional VQ-VAE, where the codebook is updated during training, our model employs a fixed codebook. However, because we apply a linear compression layer to 284 the text embeddings, we use two distinct terms to maintain effective alignment: - Quantization Latent Loss ( $\mathcal{L}_{vq,latent}$ ) ensures that the quantized embeddings (obtained from the codebook) closely match the encoder's output. 287 **Commitment Loss** ( $\mathcal{L}_{commit}$ ) encourages the encoder's output to remain close to the quantized embeddings. This ensures the model consistently uses the closest 289 embedding in the codebook, preventing the encoder from drifting too far from the 290 quantized space. 291 These losses are defined as:  $\mathcal{L}_{\text{vq\_latent}} = \|\mathbf{z}_{\text{aligned}} - \text{sg}(\mathbf{z}_{\text{encoder}})\|_2^2$ 293  $\mathcal{L}_{\text{commit}} = \|\mathbf{z}_{\text{encoder}} - \text{sg}(\mathbf{z}_{\text{aligned}})\|_2^2$ 295  $\mathcal{L}_{vq} = \mathcal{L}_{vq\_latent} + \mathcal{L}_{commit}$ 296 297 where: 298  $- \mathbf{z}_{encoder}$  is the output of the encoder. 299  $- \mathbf{z}_{aligned}$  refers to the quantized embeddings obtained from the vector quantizer. 300  $- sg(\cdot)$  is the stop-gradient operator, which ensures gradients don't flow back to the 301 quantizer during optimization. 302 By jointly minimizing these two components, we ensure that the encoder's outputs 303 align well with the codebook, while keeping the representations compact and consistent. 305 • **Reconstruction Loss** ( $\mathcal{L}_{recon}$ ): The Reconstruction Loss measures how well the re-306 constructed time series  $\hat{\mathbf{x}}$ , generated by the decoder, matches the original input  $\mathbf{x}$ . 307 This loss ensures that essential information from the original time series is preserved: 308  $\mathcal{L}_{\text{recon}} = \|\mathbf{x} - \hat{\mathbf{x}}\|_{1}$ 310 • Total Loss: The final loss function is a weighted sum of the Vector Quantization Loss 311 and the Reconstruction Loss: 312  $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{vq}} + \lambda_{\text{recon}} \cdot \mathcal{L}_{\text{recon}}$ 313 314 where  $\lambda_{recon}$  is a hyperparameter that balances the importance of the reconstruction 315 task relative to the vector quantization. We further investigate the sensitivity of this 316 parameter in Appendix H. 317 318 Stage 2: Supervised Multimodal Fusion for Prediction In this stage, we perform supervised time 319 series prediction, conditioned on textual information. This leverages the embeddings learned in 320 Stage 1, combined with a Large Language Model (LLM), to improve prediction accuracy. 321 1. Text Embedding Extraction: A pre-trained LLM extracts embeddings from textual inputs 322 (e.g., "orders" or instructions). These embeddings capture the semantic context of the text, enriching the model's understanding for time series prediction.

2. **Multimodal Fusion**: Time series embedding works as query to get text embeddings combined using a cross-attention mechanism, which dynamically identifies relevant parts of both data sources. And then the embeddings are compressed to half of original dimension to further extract the essential information. This is followed by an adaptive Hadamard product, where time series embeddings *A*, text embeddings *B*, and a learnable weight matrix *W* interact element-wise to form the joint representation:

$$C_{ij} = A_{ij} \cdot W_{ij} \cdot B_i$$

3. **Supervised Training Objective**: The fused representation is processed by the LLM, whose output is passed into a CNN-based decoder (serving as the AutoPrompter decoder to translate from word embeddings to time-series again). The model is trained to minimize the mean absolute error (MAE) between predicted and actual values:

$$L_{\text{MAE}} = \frac{1}{N} \sum_{i=1}^{N} |y_i - \hat{y}_i|$$

**Key Innovation**: The core innovation of our framework is the semantic alignment of time series data with textual context through self-supervised learning. This approach enhances both interpretability and predictive accuracy by integrating multimodal data sources, making it particularly effective for applications where numerical signals and natural language descriptions are interdependent.

5 EXPERIMENTS

#### 5.1 EXPERIMENTS SETUP

In this section, we present the experimental framework used to evaluate the performance of our
 proposed models, AutoPrompter and TITSP. The experiments are designed to assess key aspects
 such as unsupervised alignment, zero-shot generalization, and overall performance metrics.

For computational efficiency, we use **Qwen2-1.5B** (Yang et al., 2024) as our backbone language model, although other language models could also be used.

- **Stage 1**: To enhance both performance and resource efficiency, we employ the pre-trained codebook from Qwen, which is linearly compressed to 512 dimensions. This dimensionality was chosen based on empirical evidence, as illustrated in Figure 20 in the Appendix, which demonstrates an optimal balance between these factors.
  - **Stage 2**: The parameters from Stage 1 are kept fixed, and we perform end-to-end training using our designed loss function.
- 362 5.2 AUTOPROMPTER: UNSUPERVISED ALIGNMENT EVALUATION

We evaluate the AutoPrompter using key metrics: reconstruction loss, zero-shot performance,
 and latent space visualization. Our approach is compared against several baseline architectures to
 highlight the contributions of the model's components.

- **Baselines**: The following baseline models are introduced for comparison:
  - No Cross-Attention: This baseline excludes the cross-attention mechanism between time series and text embeddings. Instead, it directly employs a traditional Vector Quantized Variational Autoencoder (VQ-VAE) to map the time series embeddings to the nearest or most similar embeddings in the fixed pre-trained language codebook.
    - **Reduced Codebook Size**: Employs a smaller compressed codebook to evaluate its impact on performance. The codebook is compressed into 50.
- 376 Zero-shot Performance: Table 1 summarizes the performance of AutoPrompter and also its zero 377 shot generalization performance when training on ETTh1 and testing on various datasets. The Re construction loss will demonstrate how well the information of time-series is represented using

385

388 389

390

391

392

393

394

396 397

408

409

415

417

418

419

420 421

422

423

424

425

378	Dataset	AutoPrompter	AutoPrompter (Zero-shot)	No Cross-Attention	No Cross-Attention (zero-shot)	Reduced Codebook
379	ETTh1	0.018	0.018	0.104	0.104	1.025
	National Illness	0.029	0.031	0.138	0.173	1.141
380	Traffic	0.027	0.056	0.273	0.475	2.131
	Exchange Rate	0.005	0.015	0.820	1.090	1.053
381	Weather	0.007	0.008	0.540	0.950	1.101
382	ETTm1	0.019	0.030	0.147	0.102	1.512
001	ETTh2	0.011	0.015	0.129	0.145	1.357
383	Electricity	0.035	0.058	0.255	0.410	2.092
201	ETTm2	0.019	0.026	0.122	0.192	0.947
304	Lorenz Time	0.010	0.011	0.132	0.212	0.133

Table 1: Zero-shot reconstruction loss comparison across different architectures. The bestperforming value for each dataset is in **bold**, and the second-best is in *italics*.

compressed language codebook. The results clearly demonstrate that the AutoPrompter achieves superior performance across all datasets.

**Latent Space Visualization (t-SNE)**: To examine the alignment of time series and text embeddings, we visualize the latent space using t-SNE. Figure 5 shows that AutoPrompter with the compressed codebook results in well-clustered embeddings. In comparison, a trainable randomly initialized codebook produces scattered, unstructured embeddings after training.



Figure 5: t-SNE Visualization: with pre-trained Codebook (Right) vs. with trainable Randomly initialized Codebook (Left)

410 Conclusion: The evaluation results demonstrate that the AutoPrompter significantly outperforms
 411 the baseline models in zero-shot tasks. The t-SNE visualization emphasizes the importance of the
 412 codebook in generating semantically structured embeddings, while the codebook size analysis suggests that an optimal size enhances performance, underscoring the critical role of the codebook in
 414 our framework.

416 5.3 OVERALL PERFORMANCE EVALUATION OF TITSP

To provide a thorough assessment of the TITSP model's performance, we focus on three critical metrics: **Order Compliance Rate**, **Zero-shot Ability**, and the performance of TITSP for long sequence for the instruction.

- Order Compliance Rate: To provide a thorough assessment of the TITSP model's performance, we introduce the *Compliance Rate* metric to quantitatively evaluate how well the model adheres to specified actions guiding time series predictions. This metric reflects the proportion of time steps where the model's predictions align with the expected trends or patterns dictated by the actions.
- The Compliance Rate R is defined as:  $R = \sum_{t=1}^{T_1} C_t$ , where  $C_t$  is an indicator variable such that  $C_t = 1$  if the prediction at time step t adheres to the prescribed action, and  $C_t = 0$ otherwise. Here,  $T_1$  represents the total number of time steps considered.
- A higher compliance rate indicates better adherence to specified actions, showcasing the model's ability to incorporate desired behaviors into its predictions. For detailed definitions and calculations, refer to Appendix E. In this study, we compare our method with Time-LLM (Jin et al., 2023), which uses natural language as input with well-designed prompts

to aid predictions. Additionally, we introduce two other baselines that use only LLMs with a prompt concatenating both instructions and time series, as well as UniTime (Liu et al., 2024b) and GPT4MTS(Jia et al., 2024). We adapt GPT4MTS to Qwen4MTS with the same language model we use so that we can fairly compare. Further insights into the Time-LLM prompt are provided in Appendix A. We also provide the detail implementation of Qwen4MTS and UniTime(Qwen) in Appendix I. Pure LLM results are also provided in Appendix I. We did not compare with Chronos (Ansari et al., 2024) and UniTS (Gao et al., 2024) as they only input time-series.

Instruction	TITSP		Time-LLM		Qwen4MTS		UniTime (Qwen)		Llama-3.1-8B	
Metric	CR MSE CR MSE CR MSE		CR	MSE	CR	MSE				
Linear Growth and Linear Decay	0.83	1.15	0.38	3.45	0.69	<u>1.90</u>	0.54	2.73	0.32	4.95
Linear Growth and Linear Decay	0.79	1.17	0.49	2.85	0.79	<u>1.34</u>	0.57	2.28	0.41	2.80
Linear Trend Up	<u>0.90</u>	1.03	0.63	1.71	0.76	<u>1.08</u>	0.63	1.65	0.91	1.15
Linear Trend Down	0.87	0.88	0.64	1.55	0.71	1.36	0.51	1.59	0.85	0.92
Exponential Growth	0.89	1.33	0.58	2.59	0.63	2.07	0.60	2.38	0.58	2.35
Exponential Decay	0.84	1.25	0.56	2.26	0.67	2.10	0.69	2.05	0.46	2.39
Keep Stable	0.98	0.35	0.76	0.76	0.93	0.48	0.83	0.62	0.95	0.33
Decrease Amplitude	0.90	0.91	0.85	1.04	0.90	0.84	0.79	1.09	0.52	1.89
Increase Amplitude	0.94	0.94	0.79	1.20	0.89	0.96	0.81	1.03	0.75	1.35
Logarithmic Growth	0.77	1.65	0.49	2.31	0.79	1.55	0.60	1.73	0.55	1.94
Logarithmic Decay	0.83	1.68	0.48	2.19	0.81	<u>1.69</u>	0.67	2.04	0.63	2.60

Table 2: Comparison of Compliance Rate (CR) and MSE for TITSP, Time-LLM, Qwen4MTS, UniTime, and Llama-3.1-8B across various instructed actions with highlighted best (in **bold**) and second-best (underline) results.



Figure 6: comparison of Time-LLM and TITSP. Our method succeed in learning the dependency of prediction and instruction, (red: Time-LLM, blue: TITSP)

• Zero-shot Ability: The TITSP model demonstrates strong zero-shot generalization, effectively adapting to variations of action instructions without retraining. Across a range of similar expressions, it achieves high compliance rates, showcasing its robust ability to understand and follow instructions with minimal performance degradation. The detailed performance for selected instructions is presented in Table 3 in Appendix, where the model maintains near-perfect compliance, indicating superior adaptability across different contexts.

To further highlight the TITSP model's zero-shot generalization ability, we present time series predictions for various instructions. These examples in Figure 7 demonstrate the model's capacity to adapt to new, unseen patterns, reinforcing its robustness in generalizing across different instructions.

Training Instruction	Test Instruction	Compliance Rate	MSE
	Linear Upward	0.81	1.27
Linear Trend Up	Linear Goes Up	0.89	0.93
	Linear Growth	0.80	1.13
Linear Growth and Decay	Linear Up and Down	0.71	1.93
Linear Decay and Growth	Linear Down and Up	0.73	1.51
Exponential Growth	Exponential Upward	0.71	1.93
Exponential Growin	Exponential Goes Up	0.82	1.53

Table 3: Zero-shot performance of TITSP on selected test instructions.



Figure 7: Samples of Zero-shot about TITSP, which illustrate that TITSP has strong zero-shot ability, which is also an advantage of LLMs.

• **Keyword Extraction Capability**: The TITSP model exhibits remarkable performance in extracting relevant keywords from long sequences. The attention mechanism effectively identifies critical components in the input data, enhancing the model's ability to discern patterns and relationships in the time series. For visual insights, refer to the attention matrix provided in Appendix G.

518 Conclusion: The results underscore the efficacy of the TITSP model, particularly in its order compliance, zero-shot learning ability, and adeptness at keyword extraction. These strengths position
520 TITSP as a powerful tool for time series prediction, capable of seamlessly integrating and adhering
521 to textual instructions while maintaining high accuracy and interpretability.

To conclude this experimental section, we address the ongoing debate regarding the relevance of
large language models (LLMs) in time-series prediction. While some researchers question their
importance (Tan et al., 2024), others highlight their potential in sequence modeling (Liu et al.,
2024a). Our work contributes to this discussion by systematically replacing LLMs with traditional
architectures, such as transformers and multilayer perceptrons (MLPs), and evaluating the resulting
performance as seen in Appendix D. Additionally, we conduct extensive ablation studies, detailed
in Appendix D, to assess the significance of specific model components..

### 6 CONCLUDING REMARKS

In this paper, we introduced *Text-Informed Time Series Prediction* (TITSP), a novel framework that enhances time series forecasting by integrating domain-specific textual information. Extensive experiments across diverse datasets demonstrate that TITSP significantly outperforms traditional and existing multimodal approaches, improving both predictive accuracy and interpretability. Notably, TITSP exhibits robust zero-shot generalization, enabling effective deployment across various domains without extensive retraining. Our findings underscore the potential of multimodal methodologies to transform time series modeling, offering more accurate and versatile solutions for real-world applications. By bridging numerical data and contextual textual information, TITSP promises substantial impacts in fields requiring precise forecasting and decision-making.

#### 540 REFERENCES 541

547

553

560

561

567

568

569

- Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, 542 Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, et al. 543 Chronos: Learning the language of time series. arXiv preprint arXiv:2403.07815, 2024. 544
- George EP Box, Gwilym M Jenkins, Gregory C Reinsel, and Greta M Ljung. Time series analysis: 546 forecasting and control. John Wiley & Sons, 2015.
- 548 Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- 549 Chris Chatfield. Time-series forecasting. Statistics in Medicine, 19(2):205-207, 2000. 550
- 551 Mingyue Cheng, Yiheng Chen, Qi Liu, Zhiding Liu, and Yucong Luo. Advancing time series clas-552 sification with multimodal language modeling. arXiv preprint arXiv:2403.12371, 2024.
- Michael Christ, Mark Ploesch, and Uwe Römer. Time series feature extraction based on scalable 554 hypothesis tests. In Proceedings of the 21st International Conference on Data Mining, pp. 945-555 950. IEEE, 2018. 556
- Shohreh Deldari, Hao Xue, Aaqib Saeed, Jiayuan He, Daniel V Smith, and Flora D Salim. Beyond 558 just vision: A review on self-supervised representation learning on multimodal and temporal data. 559 arXiv preprint arXiv:2206.02353, 2022.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. 562 arXiv preprint arXiv:2407.21783, 2024. 563
- Shanghua Gao, Teddy Koker, Owen Queen, Thomas Hartvigsen, Theodoros Tsiligkaridis, and 564 565 Marinka Zitnik. Units: Building a unified time series model. arXiv preprint arXiv:2403.00131, 2024. 566
  - Nate Gruver, Marc Finzi, Shikai Oiu, and Andrew G Wilson. Large language models are zero-shot time series forecasters. Advances in Neural Information Processing Systems, 36, 2024.
- 570 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. Neural Computation, 9(8): 571 1735-1780, 1997.
- Romain Ilbert, Ambroise Odonnat, Vasilii Feofanov, Aladin Virmaux, Giuseppe Paolo, Themis Pal-573 panas, and Ievgen Redko. Samformer: Unlocking the potential of transformers in time series 574 forecasting with sharpness-aware minimization and channel-wise attention. In Forty-first Inter-575 national Conference on Machine Learning. 576
- 577 Furong Jia, Kevin Wang, Yixiang Zheng, Defu Cao, and Yan Liu. Gpt4mts: Prompt-based large language model for multimodal time-series forecasting. In Proceedings of the AAAI Conference 578 on Artificial Intelligence, volume 38, pp. 23343-23351, 2024. 579
- 580 Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y Zhang, Xiaoming Shi, Pin-Yu Chen, Yux-581 uan Liang, Yuan-Fang Li, Shirui Pan, et al. Time-llm: Time series forecasting by reprogramming 582 large language models. arXiv preprint arXiv:2310.01728, 2023. 583
- 584 Kai Kim, Howard Tsai, Rajat Sen, Abhimanyu Das, Zihao Zhou, Abhishek Tanpure, Mathew Luo, and Rose Yu. Multi-modal forecaster: Jointly predicting time series and textual data. arXiv 585 preprint arXiv:2411.06735, 2024. URL https://arxiv.org/abs/2411.06735. Li-586 cense: CC BY 4.0. 587
- 588 Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. Foundations & trends in multimodal 589 machine learning: Principles, challenges, and open questions. ACM Computing Surveys, 56(10): 590 1-42, 2024.591
- Toni JB Liu, Nicolas Boullé, Raphaël Sarfati, and Christopher J Earls. Llms learn governing 592 principles of dynamical systems, revealing an in-context neural scaling law. arXiv preprint arXiv:2402.00795, 2024a.

- Xu Liu, Junfeng Hu, Yuan Li, Shizhe Diao, Yuxuan Liang, Bryan Hooi, and Roger Zimmermann.
   Unitime: A language-empowered unified model for cross-domain time series forecasting. In
   *Proceedings of the ACM on Web Conference 2024*, pp. 4095–4106, 2024b.
- Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long.
  itransformer: Inverted transformers are effective for time series forecasting. *arXiv preprint arXiv:2310.06625*, 2023.
- Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64
   words: Long-term forecasting with transformers. *arXiv preprint arXiv:2211.14730*, 2022.
- David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back propagating errors. *Nature*, 323(6088):533–536, 1986.
- Robert H Shumway, David S Stoffer, Robert H Shumway, and David S Stoffer. Arima models. *Time series analysis and its applications: with R examples*, pp. 75–163, 2017.
- Sima Siami-Namini, Neda Tavakoli, and Akbar Siami Namin. The performance of lstm and bilstm in forecasting time series. In *2019 IEEE International conference on big data (Big Data)*, pp. 3285–3292. IEEE, 2019.
- Mingtian Tan, Mike A Merrill, Vinayak Gupta, Tim Althoff, and Thomas Hartvigsen. Are language
   models actually useful for time series forecasting? *arXiv preprint arXiv:2406.16964*, 2024.
- The New York Times. Whiplash ends a roller coaster week, 2008. URL https://www.nytimes.com/2008/10/11/business/11markets.html.
- Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. Advances in
   *neural information processing systems*, 30, 2017.
- <sup>619</sup> Vladimir N Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- Ashish Vaswani, Noam Shard, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
   Kaiser, Andrew Ku, N. Killar, Y. Cheng, et al. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Kang Wang, Kenli Li, Liqian Zhou, Yikun Hu, Zhongyao Cheng, Jing Liu, and Cen Chen. Multiple convolutional neural networks for multivariate time series prediction. *Neurocomputing*, 360:107–119, 2019.
- Yuxuan Wang, Haixu Wu, Jiaxiang Dong, Guo Qin, Haoran Zhang, Yong Liu, Yunzhong Qiu, Jian min Wang, and Mingsheng Long. Timexer: Empowering transformers for time series forecasting
   with exogenous variables. *arXiv preprint arXiv:2402.19072*, 2024.
  - Liang Wen, Qian Wang, and Qiang Yang. Multi-task learning for time series forecasting: A survey. *arXiv preprint arXiv:1712.07450*, 2017.
  - Qingsong Wen, Tian Zhou, Chaoli Zhang, Weiqi Chen, Ziqing Ma, Junchi Yan, and Liang Sun. Transformers in time series: A survey. *arXiv preprint arXiv:2202.07125*, 2022.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. arXiv preprint arXiv:2407.10671, 2024.
  - Weiqi Zhang, Jiexia Ye, Ziyue Li, Jia Li, and Fugee Tsung. Dualtime: A dual-adapter multimodal language model for time series representation. *arXiv preprint arXiv:2406.06620*, 2024.
- Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.
  Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings* of the AAAI conference on artificial intelligence, volume 35, pp. 11106–11115, 2021.

645

631

632

633

634

635

640

641

608

648 649 650	Ro pe	<b>adm</b> r for c	<b>ap.</b> This appendix provides additional experiments and details omitted from the mair conciseness. It is organized as follows:	ı pa-
651 652		•	Section A: Explores the main disadvantage of Time LLM in handling textual context, tivating the present work.	mo-
653 654		•	Section B: Complements the main paper's experiments with additional visualization TITSP's performance and evidence of its zero-shot capability.	is of
655 656		•	Section C: Details the data generation process, emphasizing base time series and transmations to incorporate context and text.	sfor-
658 659		•	Section D: Conducts extensive ablation studies to validate the rationale of each architect element.	tural
660		•	Section E: Provides a formal definition of the compliance rate with descriptive example	es.
661		•	Section F: Presents evidence on Stage 1 and its reconstruction ability.	
662 663 664		•	Sections G and H: Explain how the model handles longer sequences and establish sensity to loss hyperparameters.	itiv-
665 666	Т	ABL	e of Contents	
668 669	А	Lim	itations of Time-LLM: Influence of Textual Prompts	15
670				
671	В	Ove	rall Performance of TITSP	15
673		<b>B</b> .1	Visualization of Instruction-Based Time-Series Predictions	15
674		B.2	Zero-Shot Samples of TITSP	16
675 676		В.3	Visualization of Normal Prompted Time-Series Prediction	17
677 678 670	С	Data	a Generation	18
680	D	Abla	ation Study	21
681		D.1	Importance of W	21
682 683		D.2	Are LLMs Really Important in Time-Series Prediction?	21
685	Е	Con	upliance Rate Definition	21
686		F 1	Notation	21
687		E.1		21
689		E.2		22
690		E.3	Compliance Conditions per Action	22
691		E.4	Calculating the Compliance Rate	23
692 693		E.5	Example Calculation	24
694 695		E.6	Discussion	24
696	F	Rec	onstruction Ability of AutoPrompter	24
697		F.1	Zero-shot Ability	24
699 699		F.2	Visualization of Reconstruction Ability	25
700		E3	Comparison with Other Baselines	25
701		F4	Performance with Different Compressed Codebook Sizes	26
		1.7	renormance with Different Compressed Codebook biles	20

G	Dealing with Long Sequences	26
н	Sensitivity of $\lambda_{recon}$	28
Ι	Detailed Experimental Setup for the Proposed Method	28
J	What Really Happened on October 10, 2008?	31

Incorrect dataset description	0.385	0.40
Random mean, max, min values	0.383	0.403
Opposite trend description	0.383	0.403

Table 4: Performance of Time-LLM under various textual modifications

### A LIMITATIONS OF TIME-LLM: INFLUENCE OF TEXTUAL PROMPTS

In this section, we investigate the limitations of Time-LLM by examining the impact of textual prompts on its performance. We systematically alter the dataset descriptions to assess whether changes in textual context, unrelated to the time series data, affect the model's predictions. The modifications tested include:

- Providing an incorrect dataset description.
- Assigning random mean, max, and min values between 0 and 1.
- Providing a description with an opposite trend.
- To evaluate the model, we use a sequence length of 96 for both input and prediction on the ETTh1 dataset. The results are summarized in Table 4.

These results indicate that altering the textual descriptions had negligible effects on the performance of Time-LLM. The MSE and MAE values remained largely unchanged, suggesting that the model's predictions are primarily driven by the time series data, with minimal influence from the textual context. This observation underscores the limitations of using prompts in this particular setup and motivates the need for further exploration to determine the conditions under which text can meaningfully guide time-series predictions.

### **B** OVERALL PERFORMANCE OF TITSP

This section presents the overall performance of the Text-Informed Time Series Prediction (TITSP) model, including visualizations of instruction-based time-series predictions and zero-shot samples.

B.1 VISUALIZATION OF INSTRUCTION-BASED TIME-SERIES PREDICTIONS

In this subsection, we provide visualizations of the prediction results with and without instructions.
 These visualizations demonstrate the model's ability to adapt to various textual instructions, show-

casing its versatility and accuracy.



#### B.2 ZERO-SHOT SAMPLES OF TITSP

861
862 In this subsection, we present zero-shot samples of the TITSP model. These samples demonstrate
863 the model's ability to generalize to new, unseen instructions without additional training, highlighting its robustness and adaptability.



These visualizations and zero-shot samples provide a comprehensive overview of the TITSP model's capabilities, demonstrating its effectiveness in integrating textual instructions with time-series data.

#### B.3 VISUALIZATION OF NORMAL PROMPTED TIME-SERIES PREDICTION

906

907 908 909

910

915

916

917

In this subsection, we present the full results of the TITSP model using only dataset descriptions and
 task descriptions.

We give here two examples of prompt that are used to show that they may differe from one dataset to another.

• "Dataset description: The Electricity Transformer Temperature (ETT) is a crucial indicator in the electric power long-term deployment. This dataset consists of 2 years data from two separated counties in China. To explore the granularity on the Long sequence time-series forecasting (LSTF) problem, different subsets are created, ETTh1, ETTh2 for 1-hour-level and ETTm1 for 15-minutes-level. Each data point consists of the target value "oil temperature" and 6 power load features. The train/val/test is 12/4/4 months. Task description: forecast the next 96 steps given the previous 96 steps information."

• "Traffic is a collection of hourly data from California Department of Transportation, which describes the road occupancy rates measured by different sensors on San Francisco Bay area freeways. Task description: forecast the next 96 steps given the previous 96 steps information."

To ensure a fair comparison, we adapt GPT4TS and Time-LLM to QWEN2-1.5B (Yang et al., 2024).
The input and output sequence length is set to 96. The results demonstrate that TITSP excels at building the dependency between prompt words and time-series input, fully utilizing the capabilities of large language models.

Models	TI	rsp	GPT	T4TS	Time	LLM	itrans	former	Rlir	near	Patcl	nTST
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1	0.416	0.422	0.428	0.426	0.450	0.445	0.454	0.448	0.446	0.434	0.469	0.455
ETTh2	0.406	0.409	0.405	0.413	0.427	0.430	0.432	0.434	0.422	0.420	0.441	0.437
ETTm1	0.375	0.399	0.390	0.405	0.408	0.419	0.417	0.425	0.406	0.411	0.423	0.426
ETTm2	0.365	0.389	0.373	0.399	0.395	0.414	0.402	0.418	0.390	0.406	0.404	0.418
Weather	0.342	0.382	0.354	0.395	0.382	0.405	0.383	0.407	0.375	0.399	0.387	0.407
Traffic	0.360	0.383	0.357	0.396	0.398	0.414	0.393	0.412	0.391	0.408	0.394	0.412
Electricity	0.357	0.381	0.351	0.392	0.402	0.415	0.392	0.409	0.396	0.409	0.389	0.408
Exchange rate	0.365	0.397	0.349	0.389	0.400	0.410	0.391	0.406	0.398	0.406	0.382	0.402
1st count	1	1	4	5	(	)	(	)	(	)	(	)

Table 5: Comparison of models based on MSE and MAE metrics with highlighted best (red) and second-best (green) results

> These results highlight the superior performance of TITSP in building dependencies between prompt words and time-series input, effectively utilizing the capabilities of large language models. The table presents a comprehensive comparison of various models, with the best and second-best results highlighted in red and green, respectively. This demonstrates the robustness and effectiveness of the TITSP model in time-series prediction tasks.

### C DATA GENERATION

In this section, we provide a detailed definition of the data generation process to ensure high-quality datasets for evaluation. We generate data from various datasets, including ETTh1, ETTh2, ETTm1, ETTm2, weather, traffic, electricity, exchange\_rate.

In this table:

- np.arange(x) generates a sequence of time steps.
- *time* represents the time variable in the time series.
- transition time is set to 0.5 for uniform experiments.
- A, B, and C are constants defining the amplitude and the slopes. In our case, we set  $A \sim \mathcal{U}(0, 0.5)$ , B and  $C \sim \mathcal{U}(0.01, 0.015)$ .

The detailed equations for *Linear Trend Up*, *Linear Trend Down*, *Linear Growth and Linear Decay*,and *Linear Decay and Linear Growth* are listed as follows:

 $\begin{cases}
\text{Original slope } m = \text{linear regression(time, batch_y_i)} \\
\text{New slope } m' = -m + \delta, \quad \delta \sim \mathcal{U}(0.01, 0.015) \\
\text{New sequence } y' = \text{batch_y_i} + m' \times \text{time}
\end{cases}$ (4)

Action	Description	Mathematic Func- tion	Generated Dataset
Linear Trend Up	Linear increase over	see Equation 4	weather, ex-
-	time	-	change_rate
Linear Trend Down	Linear decrease over	see Equation 5	weather, ex-
	time		change_rate
Exponential Growth	Exponential increase	prediction $\times \exp(B \times$	weather, ex-
	over time	np.arange(x))	change_rate, electric-
			ity
Exponential Decay	Exponential decrease	prediction ×	weather, ex-
	over time	$\exp(-B \times$	change_rate, electric-
		np.arange(x))	ity
Logarithmic Growth	Logarithmic growth	prediction $+ C \times$	weather, ex-
	over time	$\log(1 + \text{np.arange}(x))$	change_rate, electric-
			ity
Logarithmic Decay	Logarithmic decay	prediction $-C \times$	weather, ex-
	over time	$\log(1 + \text{np.arange}(x))$	change_rate, electric-
			ity
Keep Stable	Constant value	last input point	All
Linear Growth and	Linear increase fol-	see Equation 6	weather, ex-
Linear Decay	lowed by decrease		change_rate
Linear Decay and	Linear decrease fol-	see Equation 7	weather, ex-
Linear Growth	lowed by increase		change_rate
Increase Amplitude	Scale up predictions	prediction $\times (1 + A)$	ETTh1, ETTh2,
	by a factor		ETTm1, ETTm2,
			traffic
Decrease Amplitude	Scale down predic-	prediction $\times (1 - A)$	ETTh1, ETTh2,
	tions by a factor		ETTm1, ETTm2,
			traffic

Table 6: Summary of actions, their descriptions, corresponding mathematical functions, and generated datasets.

 $\begin{cases} \text{Original slope } m = \text{linear regression(time, batch_y_i)} \\ \text{New slope } m' = -m - \delta, \quad \delta \sim \mathcal{U}(0.01, 0.015) \\ \text{New sequence } y' = \text{batch_y_i} + m' \times \text{time} \end{cases}$ (5)

$$\begin{array}{l} \text{1016} \\ \text{1017} \\ \text{1018} \\ \text{1018} \\ \text{1019} \end{array} \quad \text{trend} = \begin{cases} -\text{initial\_slope} \times \text{time}(t), & t < \text{transition\_time} \\ -\text{initial\_slope} \times \text{time}(\text{transition\_time}) + \text{increase\_slope} \times (t - \text{transition\_time}), & t \geq \text{transition\_time} \\ t \geq \text{transition\_time} \\ (7) \end{cases}$$

1020 For this part, we provide some examples of data generation:

1022 C.0.1 LORENZ TIME SERIES

In addition to the generation of synthetic datasets for general time series tasks, we also generate
 data based on the Lorenz system, which is commonly used to model chaotic dynamics. The Lorenz system is governed by the following set of ordinary differential equations (ODEs):





1080 To numerically solve this system, we use methods such as the fourth-order Runge-Kutta method, 1081 initializing the system with a set of initial conditions  $(x_0, y_0, z_0)$ . The time series for each of the 1082 state variables (x(t), y(t)), and z(t) exhibits chaotic behavior, characterized by high sensitivity to the initial conditions.

1084 The generated Lorenz time series data provide a rich dataset for testing models on chaotic systems, as it challenges the model's ability to predict complex, non-linear temporal dependencies. The chaotic 1086 nature of the Lorenz system makes it an excellent benchmark for evaluating time series forecasting 1087 models.

1088 1089 1090

1099

1101

**ABLATION STUDY** D

1091 In this section, we conduct an ablation study to understand the importance of various components in 1092 the TITSP model. Specifically, we investigate the significance of the weight matrix W and the role 1093 of Large Language Models (LLMs) in time-series prediction. 1094

1095 D.1 IMPORTANCE OF W1096

To evaluate the importance of the weight matrix W, we directly use the following equation:

 $C_{ij} = A_{ij} \cdot B_{ij}$ 

1100 and compare the performance with the baseline model. The results are presented in Section D.2.

1102 D.2 ARE LLMS REALLY IMPORTANT IN TIME-SERIES PREDICTION? 1103

In this part, we follow the settings of (Tan et al., 2024) to replace LLMs with a Multi-Layer Percep-1104 tron (MLP) and compare the performance. The results are summarized in Table 7. 1105

1106									
	Action	TITSP		TITSP (without W)		TITSP (without LI	.M, with MLP)	Time-LLM (Jin et al., 2023)	
1107	Metric	Compliance Rate	MSE	Compliance Rate	MSE	Compliance Rate	MSE	Compliance Rate	MSE
1100	Linear Trend Up	0.90	1.03	0.79	1.19	0.59	2.83	0.63	1.71
1108	Linear Trend Down	0.87	0.88	0.80	0.89	0.60	2.23	0.64	1.55
1100	Exponential Growth	0.89	1.33	0.81	1.57	0.52	3.19	0.58	2.59
1109	Exponential Decay	0.84	1.25	0.79	1.42	0.44	3.25	0.56	2.26
1110	Keep Stable	0.98	0.35	0.82	0.73	0.79	0.70	0.83	0.76
	Decrease Amplitude	0.90	0.91	0.79	1.30	0.67	1.51	0.85	1.04
1111	Increase Amplitude	0.94	0.94	0.68	2.05	0.59	2.02	0.79	1.20
	Logarithmic Growth	0.77	1.65	0.69	1.69	0.53	2.58	0.49	2.31
1112	Logarithmic Decay	0.83	1.68	0.71	1.74	0.54	2.36	0.48	2.19
4440	Linear Growth and Linear Decay 1	0.83	1.15	0.63	1.89	0.54	3.09	0.38	3.45
1113	Linear Growth and Linear Decay 2	0.79	1.17	0.65	1.97	0.64	2.28	0.49	2.85

Table 7: Comparison of Compliance Rate and MSE for TITSP (with various configurations) and Time-LLM across various instructed actions.

1116 1117

1114

1115

1118 In summary, both LLMs and the weight matrix W play important roles in the task of time-series 1119 prediction. Specifically, W helps in dynamically building the complex dependency of input instruc-1120 tions and giving weights to the input embedding. LLMs play a crucial role in modeling complex 1121 patterns, which are essential for accurate time-series prediction.

1122 1123

1124

1129

1131

#### E **COMPLIANCE RATE DEFINITION**

1125 In time-series prediction tasks that are guided by specific actions or directives (as outlined in Ta-1126 ble 6), it is essential to quantitatively evaluate how well the model adheres to these prescribed 1127 behaviors. We introduce the **Compliance Rate** as a metric that measures the proportion of cases 1128 where the model's predictions align with the expected trends or patterns dictated by the actions.

- 1130 E.1 NOTATION
- We define the following terms: 1132
- 1133
- $\hat{y}_t$ : The model's predicted value at time t with the specified action.

1134 1135	• $y_t$ : The model's predicted value at time t without the specified action (i.e., the baseline prediction)
1136	$T_{\rm T}$ The total number of time stars in the time series
1137	• 1: The total number of time steps in the time series.
1138	• $1\{\cdot\}$ : The indicator function, which equals 1 if the condition inside is true, and 0 otherwise.
1139	• $\delta$ : A small positive threshold indicating the minimum acceptable rate of change (slope).
1140	• $\epsilon$ : A small threshold value representing acceptable fluctuation for the "Keep Stable" action.
1141 1142	• $t_{\text{transition}}$ : The transition time point for actions involving a change in trend.
1143 1144	E.2 GENERAL DEFINITION
1145 1146	The compliance rate $C$ is computed as:
1147	Ar.
1148	$\sum_{i=1}^{N} 1$ {Compliance Condition for series <i>i</i> is satisfied}
1149	$C - \frac{i=1}{2} $ (11)
1150	C = N (II)
1151 1152	where $N$ is the total number of time series (or segments) being evaluated.
1153	
1154	E.3 COMPLIANCE CONDITIONS PER ACTION
1155	The compliance conditions avaluate the avauall turn de of the model's predictions using statistical
1156	methods such as linear regression, which accounts for inherent fluctuations in the data
1157	methods such as finear regression, which accounts for inferent nucluations in the data.
1158	1. Linear Trend Up:
1159	• Model Fitting: Fit the following models to $\hat{u}$ :
1160	• Model Fitting. Fit the following models to $y_t$ .
1161	(a) Linear Model: $y_t = mt + c$ (b) Exponential Model: $\hat{x}_t = x_t^{bt}$
1162	(b) Exponential Model: $y_t = ae^{-t}$
1163	(c) Logarithmic Model. $y_t = u \ln(t) + c$
1164	• Goodness-of-Fit: Compute $R^2$ for each model.
1165	• <b>Compliance Condition</b> : The linear model has the highest $R^{\mu}$ and the slope <i>m</i> satisfies:
1166	$m > \delta \tag{12}$
1167	$m \ge 0 \tag{12}$
1168	2. Linear Trend Down: Similar to Linear Trend Up, but the compliance condition is $m \leq \frac{1}{2}$
1169	-0.
1170	3. Exponential Growth:
11/1	• <b>Model Fitting</b> : Fit the following models to $\hat{y}_t$ :
1172	(a) Exponential Model: $\hat{y}_t = ae^{bt}$
1173	(b) Linear Model: $\hat{y}_t = mt + c$
1175	(c) Logarithmic Model: $\hat{y}_t = a \ln(t) + c$
1176	• <b>Goodness-of-Fit</b> : Compute $R^2$ for each model.
1177	• Compliance Condition: The exponential model has the highest $R^2$ and the rate pa-
1178	rameter b satisfies:
1179	$b \ge \delta \tag{13}$
1180	4. <b>Exponential Decay</b> : Similar to Exponential Growth, but the compliance condition is $b \leq b$
1181	$-\delta$ .
1182	5 Logarithmic Growth
1183	
1184	• Model Fitting: Fit the following models to $\hat{y}_t$ :
1185	(a) Logarithmic Model: $\hat{y}_t = a \ln(t) + c$
1186	(b) Linear Model: $\hat{y}_t = mt + c$
1187	(c) Exponential Model: $\dot{y}_t = ae^{ut}$
	• <b>Goodness-of-Fit</b> : Compute $R^2$ for each model.

1188 • Compliance Condition: The logarithmic model has the highest  $R^2$  and the coefficient 1189 *a* satisfies: 1190  $a\geq \delta$ (14)1191 6. Logarithmic Decay: Similar to Logarithmic Growth, but the compliance condition is a < b1192  $-\delta$ 1193 1194 7. Keep Stable: The compliance condition is: 1195  $\sigma \leq \epsilon$ (15)1196 1197 where  $\sigma$  is the standard deviation of  $\hat{y}_t$ . 1198 8. Linear Growth and Linear Decay: 1199 • Segment Division: Divide the time series into two segments at t<sub>transition</sub>: - First segment:  $t \in [1, t_{\text{transition}} - 1]$ 1201 1202 - Second segment:  $t \in [t_{\text{transition}}, T]$ 1203 • Model Fitting for Each Segment: Fit linear, exponential, and logarithmic models to each segment. 1205 • Compliance Condition: The linear model has the highest  $R^2$  in both segments, with 1206 slopes satisfying: 1207  $\begin{cases} m_1 \ge \delta & \text{(First Segment)} \\ m_2 \le -\delta & \text{(Second Segment)} \end{cases}$ (16)1208 1209 1210 9. Linear Decay and Linear Growth: Same as above, but with conditions: 1211  $\begin{cases} m_1 \leq -\delta & \text{(First Segment)} \\ m_2 \geq \delta & \text{(Second Segment)} \end{cases}$ 1212 (17)1213 1214 10. Increase Amplitude: 1215 1216 • **Compliance Condition** (Evaluated per time step): For each time step t: 1217  $|\hat{y}_t - y_t \times (1+A)| \le \epsilon$ (18)1218 1219 11. Decrease Amplitude: 1220 • **Compliance Condition** (Evaluated per time step): For each time step t: 1221 1222  $|\hat{y}_t - y_t \times (1 - A)| \le \epsilon$ (19)1223 1224 E.4 CALCULATING THE COMPLIANCE RATE 1225 1226 The compliance rate C is calculated based on the action: 1227 1228 1. For actions evaluated per action: 1229 (a) Evaluate the compliance condition for each series. 1230 (b) Aggregate and normalize the results: 1231 1232  $C = \left(\frac{\text{Number of Series Satisfying Compliance}}{N}\right) \times 100\%$ (20)1233 2. For actions evaluated per time step: 1236 (a) Evaluate compliance at each time step. 1237 (b) Aggregate and normalize across all time steps:  $C = \left(\frac{\text{Number of Time Steps Satisfying Compliance}}{N_{\text{total}}}\right) \times 100\%$ 1239 (21)1240 1241 where  $N_{\text{total}} = N \times T$ .

## 1242 E.5 EXAMPLE CALCULATION

Suppose we evaluate 100 time series with the Linear Trend Up action. If 85 of these series have a slope m such that  $m \ge \delta$ , the compliance rate is:

 $C = \left(\frac{85}{100}\right) \times 100\% = 85\%$  (22)

#### E.6 DISCUSSION

• Threshold  $\delta$ : The threshold  $\delta$  is chosen based on domain knowledge or acceptable performance criteria. In our case, we set  $\delta = 0.01$  based on the acceptable slope range of (0.01, 0.015).

• **Threshold**  $\epsilon$ : The threshold  $\epsilon$  for "Keep Stable" actions is set at 0.001.

#### F RECONSTRUCTION ABILITY OF AUTOPROMPTER

1264 F.1 ZERO-SHOT ABILITY

In this section, we present the zero-shot generalization performance of the proposed AutoPrompter
 method. The model's ability to follow different trend instructions is measured in terms of the Compliance Rate and Mean Squared Error (MSE). Table 3 provides a detailed performance matrix, showing how well the model adheres to various trends without any fine-tuning.

Perfor	Compliance Rate	MSE	
	Linear Upward	0.81	1.27
Linear Trend Up	Linear Goes Up	0.89	0.93
	Linear Growth	0.80	1.13
	Linear Downward	0.83	1.03
Linear Trend Down	Linear Goes Down	0.88	0.83
	Linear Decay	0.74	1.63
Exponential Crowth	Exponential Upward	0.71	1.93
Exponential Growth	Exponential Goes Up	0.82	1.53
Exponential Decay	Exponential Downward	0.80	1.43
Exponential Decay	Exponential Goes Down	0.73	1.89
Logarithmic Crowth	Logarithmic Upward	0.70	1.91
Logaritinine Growth	Logarithmic Goes Up	0.73	2.13
Logorithmia Docov	Logarithmic Goes Down	0.79	1.88
Logar timine Decay	Logarithmic Downward	0.83	1.75
Linear Growth and Decay	Linear Goes Up and Linear Goes	0.71	1.93
	Down		
Linear Decay and Growth	Linear Goes Down and Linear	0.73	1.51
	Goes Up		
Increase Amplitude	Raise Amplitude	0.65	2.12
Decrease Amplitude	Reduce Amplitude	0.63	2.38

Table 8: Zero-shot performance of AutoPrompter across various instructed actions.

Table 9 and Table 10 show the average reconstruction losses for different datasets when evaluated with input lengths of 512 and 2048, respectively.

1296	Dataset	Average Reconstruction Loss	Dataset	Average Reconstruction Loss
1297	ETTh1	0.028	National Illness	0.039
1298	Traffic	0.047	Exchange Rate	0.005
1299	Weather	0.007	ETTm1	0.019
1300	ETTh2	0.021	Electricity	0.035
1301	ETTm2	0.019	Lorenz Time Series	0.010

Table 9: Reconstruction losses for various datasets: Zero-shot with input length = 512.

Dataset	Average Reconstruction Loss	Dataset	Average Reconstruction Loss
ETTh1	0.028	National Illness	0.039
Traffic	0.047	Exchange Rate	0.005
Weather	0.007	ETTm1	0.019
ETTh2	0.021	Electricity	0.035
ETTm2	0.019	Lorenz Time Series	0.010

Table 10: Reconstruction losses for various datasets: Zero-shot with input length = 2048.

#### F.2 VISUALIZATION OF RECONSTRUCTION ABILITY

Figures 14 and 18 depict the model's reconstruction ability across various datasets for input lengths of 512 and 2048, respectively.



Figure 11: Electricity Dataset (512)



Figure 14

Figure 12: ETTh1 Dataset (512)



Figure 13: ETTh2 Dataset (512)



To further evaluate the performance of AutoPrompter, we compare it against other baseline methods.
 Figure 19 illustrates how AutoPrompter performs in terms of reconstruction loss when compared to various baselines on the ETTh1 dataset.





#### G DEALING WITH LONG SEQUENCES

1401 In this section, we discuss the influence of different lengths of prompts. In practical use, users may input long sequences, but only a few words may be crucial for instructing time-series prediction. In 1402 our TITSP model, the attention query is the aligned embedding of the time-series, and the key and 1403 value are the embeddings of the prompt.







Figure 22: MSE and Compliance rate performance with different input sequence lengths

Our model can effectively extract essential information when dealing with long input sequences, performing well in this zero-shot situation. However, as the figure shows, the performance worsens as the input sequence length increases.

1483 H SENSITIVITY OF  $\lambda_{\text{RECON}}$ 

1485 In this section, we discuss the sensitivity of  $\lambda_{recon}$ . As described in Section 4.2, there are two parts 1486 to the loss function, and  $\lambda_{recon}$  is set to balance them. We choose  $\lambda_{recon}$  to be 0.1, 1, 5, 10, and 20 to 1487 evaluate the performance of the data reconstruction loss on the ETTh1 dataset. The following figure 1488 visualizes the results, and we finally choose 10 as the value of  $\lambda_{recon}$ .



Figure 23: Reconstruction and zero-shot ability with different  $\lambda_{\text{recon}}$ 

1511 I DETAILED EXPERIMENTAL SETUP FOR THE PROPOSED METHOD

1512 To ensure the quality of data generation, we utilize different datasets for different instructions. Ta-1513 ble 6 outlines which dataset is used to generate new data based on the corresponding instructions. 1514

Additionally, Section C provides details about how we select the hyperparameters and functions. 1515 Here, we summarize the process based on the original dataset.

• The sequence length is set to 192, with the first 96 points used as the input x.

• For the prediction target y, the values are adjusted according to the specific instructions.

- 1516
- 1517 1518
- 1519
- 1520 1521

1522

Table 11: Datasets with different instructions and the number of training examples.

1523	Dataset	Instruction	Count
1524		increase amplitude	
1525	ETTH1	decrease amplitude	14307
1526		keep stable	
1527		increase amplitude	
1528	ETTH2	decrease amplitude	14307
1529		keep stable	
1530		linear trend up	
1531		linear trend down	
1532		exponential growth	
1502		exponential decay	
1000	Weather	logarithmic growth	52603
1534		logarithmic decay	
1535		keep stable	
1536		linear growth and linear decay	
1537		linear decay and linear growth	
1538		linear trend up	
1539		linear trend down	
1540		exponential growth	
1541		exponential decay	
1542	Exchange rate	logarithmic growth	7207
15/13		logarithmic decay	
1545		keep stable	
1544		linear growth and linear decay	
1545		linear decay and linear growth	
1546		linear trend up	
1547		linear trend down	
1548		exponential growth	
1549		exponential decay	
1550	Electricity	logarithmic growth	26210
1551		logarithmic decay	
1552		keep stable	
1553		linear growth and linear decay	
1554		linear decay and linear growth	
1555		increase amplitude	
1556	ETTm1	decrease amplitude	57507
1000		keep stable	
1557		increase amplitude	
1558	ETTm2	decrease amplitude	57507
1559		keep stable	
1560		increase amplitude	
1561	Traffic	decrease amplitude	16476
1562		keep stable	

1563

1564

The detailed function of each instruction is described in Section C. We select the instruction for each 1565 dataset based on the original features of the dataset.

## 1566 DETAILS ABOUT THE CONVOLUTION LAYERS AND PATCHING

#### 1568 PATCHING PROCESS 1569

1570 The input time series are divided into smaller sequences of length 96, similar to the approach used1571 in PatchTST (Nie et al., 2022).

#### 1573 ENCODER ARCHITECTURE

#### The encoder in **AutoPrompter** consists of:

- Three 1-dimensional convolutional layers.
- One residual connection applied to the first layer.
- ReLU activation function after the second and third convolutional layers.
- A linear layer projects features to a dimensionality of 1536, matching the language model embedding size.

The encoder structure is summarized in Table 12. Layers with a stride of 2 compress the input sequence by a factor of 4, meaning every four points are combined into one embedding in the compressed semantic space.

#### Table 12: Encoder Configuration

Layer	Input Channels	<b>Output Channels</b>	Kernel Size	Stride	Padding
Convolution Layer 1	1	128	4	2	1
Convolution Layer 2	128	256	4	2	1
Convolution Layer 3	256	256	3	1	1

1592 1593

1594

1595

1597

1598 1599

1601 1602

1603

1609

1615

1616

1572

1574

1575 1576

1577

1578

1579 1580

1581

1586

1587

#### DECODER ARCHITECTURE

#### 1596 The decoder consists of:

- Three convolutional layers.
- One residual connection.

The decoder structure is shown in Table 13.

#### Table 13: Decoder Configuration

Layer	Input Channels	<b>Output Channels</b>	Kernel Size	Stride	Padding
Convolution Layer 1	1536	256	3	1	1
Convolution Layer 2	256	128	4	2	1
Convolution Layer 3	128	1	4	2	1

# 1610 DECODER CONNECTED TO LLM OUTPUT

The decoder connecting to the LLM output consists of a single 1-dimensional convolutional layer
(Table 14). Since the original time series is compressed by a factor of 4, there are 24 embeddings of
size 1536. The input channel is the product of 24 and 1536 after flattening.

#### Table 14: Decoder-LLM Configuration

18	Layer	Input Channels	<b>Output Channels</b>	Kernel Size	Stride	Padding
619	Convolution Layer 1	36864	96	3	1	1

## 1620 HYPERPARAMETERS USED

1622<br/>1623The hyperparameters used in the experiments are summarized in Table 15. Early stopping is applied<br/>if no improvement is observed on the test set for more than 10 epochs. For the electricity and traffic<br/>datasets, the number of epochs is reduced to 5 due to faster convergence.

11yper par ameter	value
Learning Rate Batch Size Epochs Optimizer	5.00E-04 16 20 (default), 5 (electricity/traffic) Adam
IMPLEMENTATION DETAILS FOR Llama-8B-instruct	QWEN4MTS, UNITIME AND
QWEN4MTS (FROM GPT4MTS)	
Following GPT4MTS (Jia et al., 2024):	
• Word Embedding, Position Embable.	edding, Add & Norm, and Output Linear Layer are tr
• The embedding size is set to 153	6, consistent with Qwen.
• The learning rate is set to 5e-5.	
UniTime	
Following UniTime (Liu et al., 2024b):	
• A binary indicator is used to gen	erate the mask.
• Time-series embeddings are cond	catenated with sentence embeddings.
• The hidden dimension is set to 1:	536 (matching Qwen), with n_embd set to 1536.
• The mask rate is 0.5, and the lear	ming rate is set to $1e-4$ .
• The lightweight transformer is re	placed with Qwen LLM.
• All experiments are trained for 1	0 epochs.
Llama-3.1-8B-instruct	
Prompt Design for Llama-3.1-8B-instruct	(Dubey et al., 2024):
The prompt include task description, ins example.	truction and specific input number. The following is
• "[Task Description]: forecast the	e next 96 steps given the previous 96 steps information
• "[Instruction] the prediction sho	ould follow 'Linear Growth'"
• "[Input Number]: 0.173, 0.125,	

The sharp drop in the GBP/USD exchange rate on October 11, 2008, can be traced to events on
 October 10, 2008, when global financial markets experienced a massive sell-off, exacerbating the effects of the ongoing financial crisis.

On October 10, 2008, the Dow Jones Industrial Average plunged 679 points (or about 7.3%) by the end of the trading day, after dipping as much as 800 points during intra-day trading. European markets saw even steeper declines, with the FTSE 100 (UK's benchmark index) dropping around 8.9% and Germany's DAX falling 7%. In Asia, Japan's Nikkei 225 dropped 9.6%, adding to the already severe financial instability.

1679
1680
1681
1681
1682
1682
1682
1684
1685
1685
1685
1686
1686
1686
1686
1686
1687
1688
1688
1689
1689
1680
1680
1681
1681
1681
1682
1682
1682
1683
1684
1684
1685
1685
1686
1686
1687
1687
1688
1688
1688
1689
1689
1689
1680
1680
1681
1681
1682
1682
1682
1684
1684
1685
1685
1685
1686
1686
1687
1687
1687
1688
1688
1688
1688
1689
1689
1689
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680
1680</l

In addition to the stock market crashes, the growing fear of a deep global recession weighed heavily on the British pound, especially because the UK economy was seen as particularly vulnerable due to its dependence on the financial sector. This combination of stock market turmoil, concerns over the UK's banking system, and the safe-haven demand for the U.S. dollar caused the sharp decline in the GBP/USD exchange rate on October 11, 2008 (Times, 2008).