

TS-REASONER: ALIGNING TIME SERIES FOUNDATION MODELS WITH LLM REASONING

Anonymous authors

Paper under double-blind review

ABSTRACT

Time series reasoning is crucial to decision-making in diverse domains, including finance, energy usage, traffic, weather, and scientific discovery. While existing time series foundation models (TSFMs) can capture low-level dynamic patterns and provide accurate forecasting, further analysis usually requires additional background knowledge and sophisticated reasoning, which are lacking in most TSFMs but can be achieved through large language models (LLMs). On the other hand, without expensive post-training, LLMs often struggle with the numerical understanding of time series data. Although it is intuitive to integrate the two types of models, developing effective training recipes that align the two modalities for reasoning tasks is still an open challenge. To this end, we propose TS-REASONER that aligns the latent representations of TSFMs with the textual inputs of LLMs for downstream understanding/reasoning tasks. Specifically, we propose a simple yet effective method to curate diverse, synthetic pairs of time series and textual captions for alignment training. We then develop a two-stage training recipe that applies instruction finetuning after the alignment pretraining. Unlike existing works that train an LLM to take time series as inputs, we leverage a pretrained TSFM and freeze it during training. Extensive experiments on several benchmarks demonstrate that TS-REASONER not only outperforms a wide range of prevailing LLMs, Vision Language Models (VLMs), and Time Series LLMs, but also achieves this with remarkable data efficiency, e.g., using less than half the training data.

1 INTRODUCTION

Time series analysis has long been fundamental to various real-world applications in finance, energy, weather, traffic, and other domains (Prakarsha & Sharma, 2022; Xu et al., 2023; Nie et al., 2024). Its ability to model dynamics and predict future states based on historical data makes it an indispensable tool for informed decision-making and strategic planning. While numerical attributes form the bedrock of time series analysis, human decision-making is often complemented by rich prior knowledge and qualitative contextual information, including news articles, social media trends, and expert assessments. This gap prevents analytical models from achieving a deeper, more contextualized understanding of the events and dynamics driving the numerical data. By enabling machines to understand both contextual information and numerical time series patterns, we can empower them as automated systems that assist humans in gaining deeper insights into complex phenomena.

Recent advances in Time Series Foundation Models (TSFMs) have significantly enhanced the understanding of time series data through large-scale pretraining. These models are capable of generalizing across a wide variety of time series tasks and domains. Although TSFMs (Goswami et al., 2024; Das et al., 2024) demonstrate strong modeling capabilities, most are pre-trained exclusively on unimodal numerical time series and cannot therefore comprehend or integrate textual information. On the other hand, large Language Models (LLMs) and Vision Language Models (VLMs) can take texts and images as input context, and have demonstrated remarkable reasoning and problem-solving abilities across various tasks (Wei et al., 2022; Yao et al., 2023; Hao et al., 2023; Yu et al., 2024), sparking interest in transferring their capabilities to time series analysis. Some studies (Gruver et al., 2023; Liu et al., 2024c; Jia et al., 2024) transform numerical time series into string form and perform time series forecasting on LLMs by prompting them with the strings. However, despite their strong reasoning abilities, LLMs struggle to capture temporal dependencies due to their inherent lack of temporal understanding (Fons et al., 2024; Merrill et al.,

2024) and limited ability to interpret numerical values. These limitations hinder their understanding of time series data. As shown in the figure 1, TSFM and LLM have complementary strengths; the former specializes in temporal understanding, while the latter excels at text understanding.

To combine the complementary strengths of TSFMs and LLMs while overcoming their respective limitations, we propose TS-REASONER, a Time Series Large Language Model (TSLLM) designed to enhance time series reasoning by aligning a TSFM with an LLM. Specifically, we first employ the TSFM to extract rich temporal representations from numerical time series data. To effectively incorporate this temporal information into the LLM, TS-REASONER introduces a TS-to-Text adapter, which projects the TSFM-extracted temporal features into the LLM’s input embedding space. This enables seamless integration of the TSFM’s temporal understanding with the LLM’s powerful linguistic and reasoning capabilities. Our training framework consists of two stages: pretraining and fine-tuning. In the pretraining stage, we finetune TS-REASONER to produce textual captions of input time series and achieve a fundamental alignment. To this end, we propose a simple yet effective prompting strategy to curate high-quality captions for diverse time series data using advanced LLMs/VLMs. In the fine-tuning stage, we further enhance the model’s reasoning abilities through instruction tuning, ensuring robust performance in downstream tasks.

Our work makes unique contributions to a recent line of research combining TSFMs and LLMs. First, our formulation sets up the connection between LLMs and TSFMs, facilitating time series reasoning through the integration of rich contextual information and LLM reasoning. Second, we address a critical data bottleneck by a simple yet effective time series captioning method, which diversifies the training data for aligning LLMs and TSFMs. Finally, we offer new empirical insights into the strengths and limitations of existing approaches.

We evaluate the understanding and reasoning capabilities of our approach on two standard benchmarks: TimeSeriesExam (Cai et al., 2024a) and MTBench (Chen et al., 2025). TS-REASONER significantly outperforms a wide range of baseline models, including LLMs, VLMs, and the TSLLMs, as shown in Figure 2. Finally, comprehensive analyses, including extensive ablation studies, validate the effectiveness of our key designs and establish the superiority of TS-REASONER in generalization performance, scalability, and training data efficiency.

2 RELATED WORK

LLMs for Time Series. LLMs have recently garnered significant interest in time series analysis. Traditional time series forecasting relies on statistical models (RB, 1990) or data-driven neural networks (Liu et al., 2021; Lim et al., 2021; Wu et al., 2021; Zhou et al., 2022; Li et al., 2023b; 2024b) for tasks like weather and stock prediction. Recent efforts explore LLMs for this task, with some designing prompts to elicit forecasting abilities (Cao et al., 2023; Chuang* et al., 2024). Others focus on enabling LLMs to understand time series data by converting it into textual sequences or aligning its embeddings with language model embeddings via prompting or semantic information (Jin et al., 2023; Sun et al., 2023; Pan et al., 2024). In addition, multimodal vision-based LLMs are being

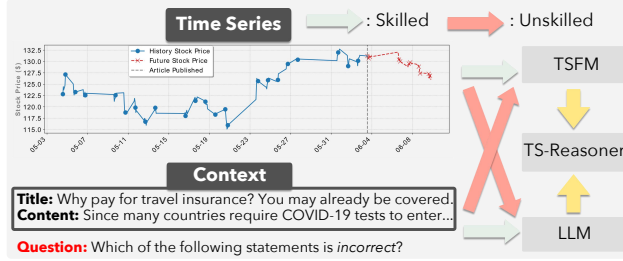


Figure 1: Time series forecasting vs. reasoning. The time series reasoning task requires both contextual reasoning (e.g., news) by LLMs and numerical understanding by TSFM.

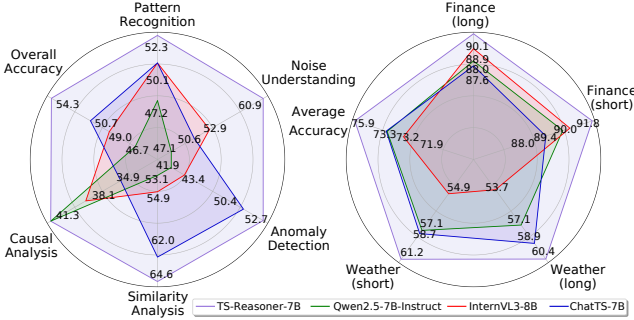


Figure 2: Results on time series understanding and reasoning benchmarks. TS-REASONER demonstrates a consistent advantage over the prevailing LLMs, VLMs, and TSLLMs.

investigated for time series prediction (Chen et al., 2024c; Zhong et al., 2025). Though LLMs exhibit non-trivial performance on some forecasting tasks, Merrill et al. (Merrill et al., 2024) indicate that LLMs struggle to reason time series. To tackle this challenge, several works (Chow et al., 2024; Zhang et al., 2025a; Xie et al., 2024) propose to enable LLMs to understand the time series with context. TS-REASONER lies in this direction, distinguishing itself by employing a pre-trained Time Series Foundation Model to ground the LLM’s reasoning in robust temporal features.

Modality Alignment. Modality alignment methods are widely studied in the multimodal domain (Li et al., 2022; Lai et al., 2024; Li et al., 2023a; Liu et al., 2024b). Inspired by the success of multimodal alignment, recent works treat time series as another modality and align it to the LLM (Xie et al., 2024; Zhang et al., 2025a). Though they achieve a certain degree of time series understanding, they focus on narrow domains (e.g., electricity) and tasks (e.g., time series understanding), and train time series encoders from scratch. In contrast, we adapt the successful training paradigm in VLMs, identify and address the key challenges (e.g., Integration of characteristics of time series into LLMs, and the shortage of time series-text pairs) faced in applying this paradigm to the unique modality of time series, exploring pre-trained time series foundation models to exploit rich time series knowledge.

Time Series Foundation Models. Recent advancements in pre-training methods are significantly contributing to the development of foundation models for time series analysis. Early efforts, such as TST (Zerveas et al., 2021) and PatchTST (Nie et al., 2022), applied BERT-like masked pretraining techniques, focusing on point-level and patch-level masking, respectively. A separate line of work, exemplified by models like TimesFM (Das et al., 2024), Timer (Liu et al., 2024e), TTMS (Ekambaram et al., 2024), Chronos (Ansari et al., 2024), and Time-MoE (Shi et al., 2024), Moirai (Liu et al., 2024d), TimesBERT (Zhang et al., 2025b), and Sundial (Liu et al., 2025) demonstrates the advantages of large-scale pre-training for improving forecasting performance. Exploring diverse pre-training objectives, MOMENT (Goswami et al., 2024) leverages a T5 encoder to achieve strong downstream multi-task capabilities. ChronoSteer (Wang et al., 2025a) also explores the alignment between TSFMs and LLMs, yet it leverages the LLM’s revisions to enhance TSMF’s forecasting capability.

3 TS-REASONER FOR TEMPORAL REASONING

As illustrated in Figure 3, TS-REASONER is composed of (1) a pretrained TSFM that encodes normalized, non-overlapping patches of input time series into compact embeddings; (2) a pretrained LLM, and (3) a TS-to-Text adapter that projects the TSFM’s output embedding to the input space of the LLM. The LLM concatenates the sequence of projected time series features with the sequence of embeddings for input text tokens, with the former demarcated by special tokens “<ts><ts/>”. The training of TS-REASONER consists of two stages: (1) a pretraining stage to align time series features from the TSFM with the LLM, using time series caption data synthesized by an advanced LLM/VLM, and (2) an instruction tuning stage to enhance complex reasoning capabilities on downstream tasks.

3.1 MODEL ARCHITECTURE

Given a natural language context \mathcal{X} and a corresponding set of time series $\mathcal{S} = \{\mathcal{T}_0, \mathcal{T}_1, \dots, \mathcal{T}_K\}$, we first project both into a shared embedding space. Specifically, for each time series $\mathcal{T}_i \in \mathbb{R}^{L_i}$, where L_i is the length of the series, we first apply instance normalization to standardize its distribution to zero mean and unit variance. This preprocessing step ensures that the model is robust to shifts and scales in the input data. Subsequently, we partition the normalized time series into a sequence of non-overlapping patches, each of a fixed length P . This patching strategy yields a sequence of $N_i = \lfloor L_i/P \rfloor$ patches, transforming the time series into a tensor $\mathcal{T}_i^p \in \mathbb{R}^{N_i \times P}$. These patches are then encoded using the TSFM, which acts as our time series feature extractor. The TSFM processes the sequence of patches and produces a sequence of embedding vectors:

$$\mathcal{Z}_i^T = \text{TSFM}(\mathcal{T}_i^p) \in \mathbb{R}^{N_i \times d_{ts}}, \quad (1)$$

where d_{ts} denotes the dimension of the time series embeddings. Concurrently, the natural language context \mathcal{X} is tokenized and fed into the pre-trained LLM’s embedding layer. This process converts the textual input into a sequence of contextualized token embeddings:

$$\mathcal{Z}^L = \text{LLM}_{\text{embed}}(\mathcal{X}) \in \mathbb{R}^{M \times d_{\text{text}}}, \quad (2)$$

where M is the number of tokens in the instruction, and d_{text} is the dimensionality of the LLM’s hidden states. To align the dimension and semantics of embeddings between LLM and TSFM, we

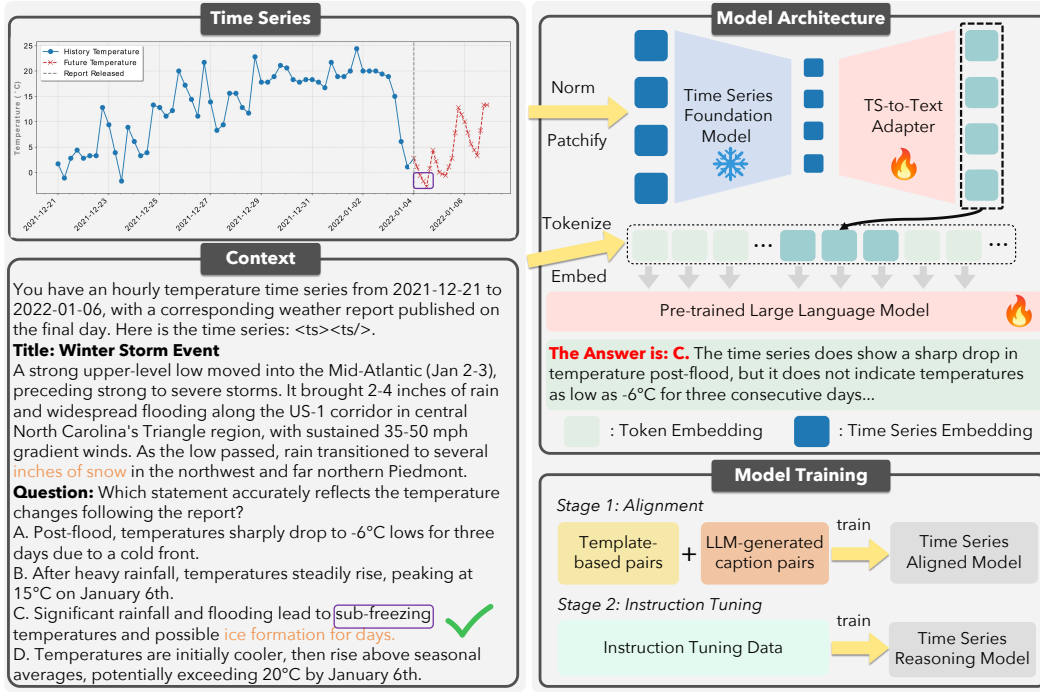


Figure 3: Overview of TS-REASONER architecture and training pipeline. To perform reasoning, a time series is first encoded by a pretrained Time Series Foundation Model (TSFM). Its output features are then projected into the LLM’s input embedding space by a trainable TS-to-Text Adapter and subsequently processed by the LLM. The model is trained in two stages: (1) a pretraining stage that aligns the TSFM outputs with the LLM inputs using both template-based (code-synthesized) and LLM-generated captions, as described in §3.2, and (2) an instruction-tuning stage to improve complex reasoning capabilities.

use a multilayer perceptron (MLP) as a TS-to-Text Adapter to transform the time series embedding into the text embedding space:

$$\mathcal{H}_i^T = \text{MLP}(Z_i^T) \in \mathbb{R}^{N_i \times d_{\text{text}}}, \quad (3)$$

To form a unified input sequence for the LLM that accommodates multiple time series, we structure the natural language instruction \mathcal{X} to include K indicators, $\{K \cdot \langle \text{ts} \rangle \langle \text{ts} \rangle\}$. The i -th placeholder $\langle \text{ts} \rangle \langle \text{ts} \rangle$ marks the insertion point for the corresponding i -th time series \mathcal{T}_i .

Let $\{\mathcal{H}_i^T \in \mathbb{R}^{N_i \times d_{\text{text}}}\}_{i=1}^K$ be the set of projected time series embeddings. The final input sequence H is constructed by sequentially inserting the embedding to each $\langle \text{ts} \rangle \langle \text{ts} \rangle$ with its corresponding time series embedding sequence \mathcal{H}_i^T . This substitution process results in a composite sequence where language and time series representations are interleaved. The total length of this fused sequence is $M + \sum_{i=1}^K N_i$. The final tensor fed to the LLM’s transformer layers is therefore: $\mathcal{H} \in \mathbb{R}^{(M + \sum_{i=1}^K N_i) \times d_{\text{text}}}$. This strategy enables the LLM to process multiple, arbitrarily placed time series within a single, coherent context and capture complex inter-series and text-series dependencies. After the combination, the input embedding H is fed to LLM to produce the final prediction \mathcal{Y} .

3.2 TRAINING RECIPE

Our training process consists of two sequential stages: the first stage aligns time series data with the LLM to establish a foundational understanding of temporal-textual relationships, while the second stage refines the LLM’s reasoning capabilities to interpret and analyze these aligned representations. Throughout both stages, we keep the parameters of the TSFM frozen to preserve its pretrained temporal knowledge, while allowing the LLM’s parameters to remain trainable, ensuring adaptive learning without compromising the integrity of the encoded time series features.

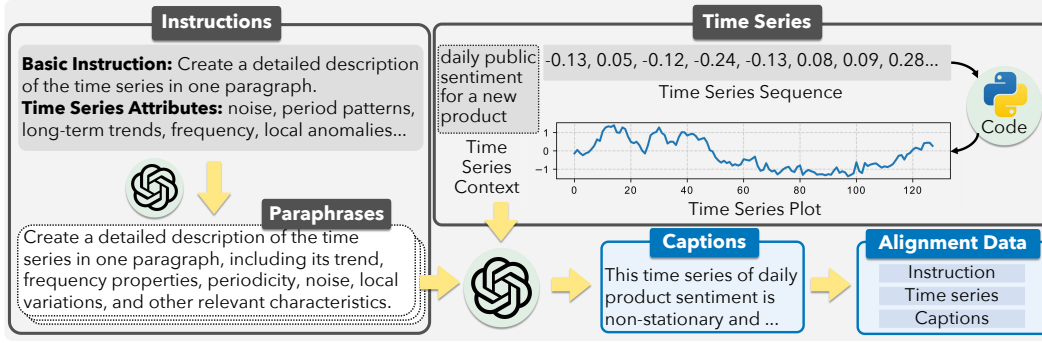


Figure 4: Workflow for our attribute-aware caption synthesis, designed to curate training data for alignment in stage 1. It enriches basic instructions with key attributes and generates diverse paraphrases, yielding the high-fidelity captions to train TS-REASONER effectively.

Stage 1: Pre-training for Language-Timeseries Alignment. In this stage, our primary objective is to align temporal data with textual information. We initially leverage synthesized data from (Xie et al., 2024), which provides predefined templates to describe time series attributes. However, while this template-based data offers accurate numerical information, its focus on specific time series patterns limits diversity, and the caption structure is monotonous. This lack of diversity can lead to overfitting to the templates, encouraging the model to learn shallow patterns and resulting in poor generalization ability (Dong et al., 2025; Choi et al., 2024). To alleviate this problem, we draw inspiration from captioning techniques in multimodal LLMs (Chen et al., 2024a). We synthesize comprehensive captions using advanced LLMs (e.g., GPT-4.1) to enrich our alignment data. Specifically, we collect time series from two sources: (Merrill et al., 2024), which includes contextual information, and synthetic data from Chronos (Ansari et al., 2024), which provides pure numerical time series.

Attribute-aware Captioning. Caption generation has been extensively investigated in visual domains (Cheng et al., 2023; 2025; Chen et al., 2024b), playing a crucial role in multimodal alignment. However, time series captioning remains largely underexplored, presenting a significant impediment to achieving comprehensive alignment. To address this gap, we introduce a straightforward approach for generating scalable time series captions, as shown in Figure 4.

Given a time series \mathcal{T} with a temporal context \mathcal{C} , we begin by defining a fundamental captioning instruction, denoted as $\mathcal{I}_{\text{base}}$. To facilitate enhanced comprehension by LLMs, we transform the time series into an image plot via Python code, $I_{TS} = \Phi(\mathcal{T})$. As evidenced in Table 1 (Section 4.1), presenting the time series as an image to advanced LLMs (e.g., GPT-4.1) demonstrates a substantial advantage in understanding compared to providing it as a raw numerical series.

To enrich the generated captions, we first identify a set of G pertinent attributes of the time series, denoted as $\{a_1, a_2, \dots, a_G\}$ (e.g., trend, frequency, periodicity, noise, local variations). These attributes are then incorporated into the basic instruction, yielding an augmented instruction $\mathcal{I}' = \mathcal{I}_{\text{base}} \cup \{a_1, a_2, \dots, a_G\}$. To further promote caption diversity, we leverage the LLM to paraphrase \mathcal{I}' into R distinct instructions, forming a candidate set of prompts $\mathcal{P} = \{\mathcal{I}'_1, \mathcal{I}'_2, \dots, \mathcal{I}'_R\}$. For each time series \mathcal{T} , a single prompt \mathcal{I}'' is uniformly sampled from this set. The final caption is then generated conditioned on the sampled prompt and the time series visualization:

$$\text{Caption} = \text{LLM}(\mathcal{I}'', I_{TS}), \quad (4)$$

where $\mathcal{I}'' \sim \mathcal{U}(\mathcal{P})$. The prompts are shown in the Figure 10 in Appendix E. We randomly sample 10K time series from each of two distinct sources: the Chronos synthetic dataset (Ansari et al., 2024), which contains purely numerical time series, and a dataset of text-attributed time series from Merrill et al. (Merrill et al., 2024), which provides contextual backgrounds. The construction of data offers two benefits: (1) Pure time series data enables the model to build a foundational understanding of temporal patterns by focusing solely on the intrinsic characteristics of the data. (2) Context-augmented time series enhances domain-specific comprehension by linking numerical trends to real-world scenarios, thereby improving the model’s ability to generalize across diverse applications.

Stage 2: Instruction Finetuning for Time Series Reasoning. To elevate the model’s capabilities from foundational understanding to complex reasoning, we employ an instruction fine-tuning stage

based on the instruction tuning dataset (Xie et al., 2024), which encompasses a wide range of Q&As and instruction-following tasks. This training facilitates TS-REASONER with two critical abilities: the fidelity to adhere to complex instructions and structured response formats, and the capacity for nuanced, context-driven reasoning on time series-specific queries.

4 EXPERIMENTS

Datasets To assess the capabilities of TS-REASONER, we conduct comparative experiments against various baselines on benchmarks tailored for time series reasoning. Our evaluation incorporates TimeSeriesExam (Cai et al., 2024a), a comprehensive multiple-choice question answering dataset. TimeSeriesExam is specifically engineered to systematically evaluate a model’s time series understanding and reasoning abilities across several key aspects: Pattern Recognition (PR), which addresses identifying trends, cycles, and stationarity; Noise Understanding (NU), focused on recognizing noise types such as white noise and random walks; Anomaly Detection (AD), for detecting unusual patterns; Similarity Analysis (SA), which involves comparing the shape and distribution of two time series; and Causality Analysis (CA), assessing the recognition of Granger Causality between time series. Furthermore, we evaluate on MTBench (Chen et al., 2025), a large-scale benchmark for evaluating time series reasoning in the real-world financial and weather domains, featuring questions that span both short-term (7-day) and long-term (14-day) temporal horizons.

Baselines and Evaluation Metrics We compare our method against three types of baselines: closed-source LLMs / VLMs, open-source LLMs / VLMs, and TSLLMs. Specifically, for closed-source models, we include GPT-4o, GPT-4.1 (gpt, 2024; Achiam et al., 2023), Claude-Sonnet-3.7 (The), and DeepSeek-Chat (Liu et al., 2024a). For open source LLM, we evaluated LLama-3.1-8B-Instruct (Grattafiori et al., 2024), Qwen-2.5-7B-Instruct (Yang et al., 2024), GLM-4-9B-Chat (GLM et al., 2024), InternLM3-8B-Instruct (Cai et al., 2024b), and Ministral-8B-Instruct (Jiang et al., 2024). Time series are transformed into textual sequences of numbers for LLMs. For open-source VLM models, we compare Qwen2.5-VL-7B (Bai et al., 2025), Phi-4-Multimodal-Instruct (Abouelenin et al., 2025), Llama3-LLaVA-Next-8B (Li et al., 2024a), InternVL3-8B (Zhu et al., 2025), and MiniCPM-V-2.6 (Yao et al., 2024). Time series are transformed into plots via code for VLMs. For TSLLMs models, we compare with ChatTime-7B (Wang et al., 2025b), ChatTS-14B (Xie et al., 2024), and we use the official training data and code to fine-tune a 7B model for a fair comparison. As all benchmarks are multiple-choice Q&As, we use accuracy as the evaluation metric. See more implementation details in Appendix A.

4.1 MAIN RESULTS

Table 1 presents the performance of all models on the two benchmarks. The best results are bolded, and the second-best results are underlined. Based on the results, we have the following key observations:

(i) TS-REASONER achieves superior overall performance on all benchmarks among models of the same size. Specifically, TS-REASONER demonstrates superior performance, surpassing the best-performing LLM by 7.60% overall, the best VLM by 5.25% overall, and the TSLLM by 3.54% overall on the TimeSeriesExam benchmark. Compared to the backbone model, TS-REASONER improves on our backbone LLM performance by a substantial 16.29%. TS-REASONER also excels the best baseline on MTBench by around 2%. In addition, TS-REASONER performs even competently with ChatTS-14B, which has a larger base model. The notable improvement demonstrates the effectiveness of our model in various time series reasoning scenarios by introducing the temporal information of TSFM for the LLM.

(ii) TS-REASONER delivers consistent gains in most time series understanding and reasoning subtasks. In particular, it surpasses the second best baseline with absolute improvements of 2.16% in *Pattern Recognition*, 8.05% in *Noise Understanding*, 2.33% in *Anomaly Detection*, and 2.65% in *Similarity Analysis*, while also yielding around 2% improvements on both financial and weather reasoning tasks. These gains stem from two key factors: (1) aligning time series with text during training substantially strengthens TS-REASONER’s understanding of time series patterns; and (2) this improved understanding further enhances its ability to reason over time series in context when combined with textual information.

Table 1: Performance of LLMs, VLMs, TSLLMs, and proprietary models on time series understanding and reasoning benchmarks. Our baselines also include ChatTS-14B, which uses a larger base model.

Model	TimeSeriesExam (Cai et al., 2024a)						MTBench (Chen et al., 2025)			
	PR	NU	AD	SA	CA	OA	Finance (long)	Finance (short)	Weather (long)	Weather (short)
<i>Proprietary models</i>										
DeepSeek-Chat	65.23	55.17	52.71	63.71	42.86	59.89	89.15	90.02	59.75	58.76
Claude-Sonnet-3.7	62.26	55.17	48.06	72.57	50.79	59.63	84.11	88.56	51.24	47.91
GPT-4o	59.03	55.17	53.49	62.83	31.75	55.96	84.30	82.69	48.07	48.22
GPT-4o (vision)	67.12	62.07	62.79	64.60	26.98	62.12	84.11	80.65	46.43	48.53
GPT-4.1 (vision)	69.81	68.97	68.22	75.22	41.27	67.89	93.41	91.45	56.04	55.35
<i>Open-source Large Language Models</i>										
Llama-3.1-8B-Instruct	37.73	37.93	30.23	36.28	28.57	35.52	63.37	35.52	40.25	40.00
Qwen2.5-7B-Instruct	47.17	47.13	41.86	53.10	41.27	46.66	87.98	89.41	57.14	58.44
GLM-4-9B-chat	41.78	39.08	37.21	47.79	38.09	41.28	71.31	77.19	50.27	50.85
InternLM3-8B-Instruct	43.93	51.72	26.35	52.21	34.92	42.33	71.70	71.08	45.05	46.67
Ministral-8B-Instruct	43.13	37.93	39.53	44.25	36.51	41.55	46.32	50.71	39.15	40.93
<i>Open-source Vision Language Models</i>										
Qwen2.5-VL-7B-Instruct	25.34	32.18	19.38	42.48	12.70	26.61	81.98	86.35	52.06	46.82
Phi-4-Multimodal-Instruct	36.39	34.48	30.23	38.94	14.28	33.68	70.35	74.54	48.35	49.77
Llama3-LLaVA-Next-8B	31.27	35.63	29.46	30.09	38.09	31.85	52.14	51.50	47.53	47.29
InternVL3-8B	50.13	52.87	43.41	54.87	38.09	49.01	88.95	90.00	53.71	54.88
MiniCPM-V-2.6	29.11	39.08	27.13	51.33	31.75	33.42	81.78	83.09	48.63	45.12
<i>Time Series Large Language Models</i>										
ChatTime-7B	42.85	49.42	35.65	44.24	34.92	41.94	25.97	28.10	47.80	42.79
ChatTS-7B	50.13	50.57	50.38	61.95	34.92	50.72	87.60	88.01	58.92	58.75
ChatTS-14B*	59.30	54.02	51.16	62.83	41.27	56.36	89.22	91.22	59.61	59.22
TS-REASONER-7B (ours)	52.29	60.92	52.71	64.60	41.27	54.26	90.12	91.85	60.44	61.24
Δ Over Best 7B	+2.16	+8.05	+2.33	+2.65	+0.00	+3.54	+1.17	+1.85	+1.52	+2.49

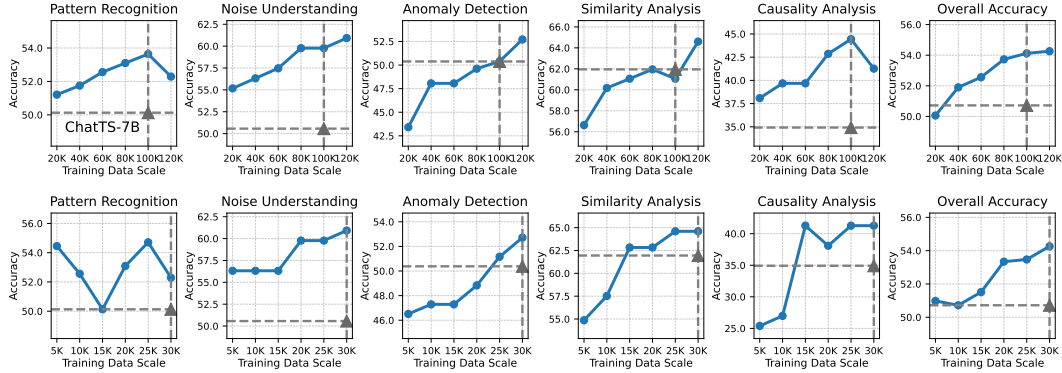


Figure 5: Data scaling and efficiency of TS-REASONER. The top (bottom) row illustrates how the performance of TS-REASONER varies when increasing the training data for alignment (instruction tuning). The columns correspond to sub-tasks in TimeSeriesExam. ChatTS-7B (Xie et al., 2024) is included for reference, denoted by the gray triangle.

4.2 ANALYSIS OF DATA SCALING AND EFFICIENCY

Figure 5 presents our data scaling analysis on the TimeSeriesExam benchmark. TS-REASONER demonstrates remarkable data efficiency compared to the ChatTS-7B baseline. For the alignment stage, TS-REASONER achieves superior overall accuracy using just 60K samples, less than half the data required by the baseline. This efficiency is even more stark in the instruction tuning stage, where 10K samples suffice to outperform ChatTS-7B. This significant reduction in data dependency stems from our pre-trained TSFM and effective alignment, which equips the LLM with a robust temporal foundation. Consequently, TS-REASONER develops advanced reasoning capabilities with a substantially smaller amount of data, marking a key advantage for practical deployments where labeled data is scarce.

4.3 CHOICES OF CAPTIONING MODEL FOR ALIGNMENT

The quality of the generated captions is a critical factor in the efficacy of our time-series-language alignment. To validate this, we conducted an experiment where we trained TS-REASONER using three distinct sets of captioning data, each generated by a model with varying capabilities:

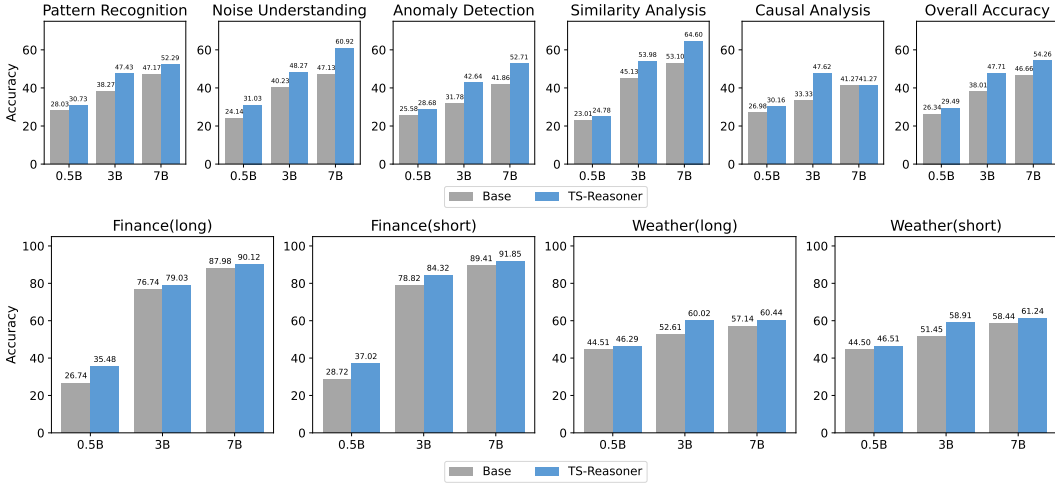


Figure 6: Performance of TS-REASONER and their associated LLM backbones (Qwen2.5 series). The top row and bottom row report the performance on TimeSeriesExam and MTBench, respectively.

the state-of-the-art GPT-4.1, and two VLMs, InternVL3-8B and Qwen2.5-VL-7B-Instruct. As illustrated in Figure 7, the results demonstrate that the performance of TS-REASONER is directly correlated with the fidelity of the captioning model. A distinct performance hierarchy emerges across both benchmarks: the model trained on GPT-4.1 captions consistently outperforms the one trained on InternVL3-8B captions, which in turn surpasses the one trained on Qwen2.5-VL-7B-Instruct captions. The higher performance gain from GPT-4.1 is attributed to its advanced capability in time series understanding. It is not surprising that the captions generated by InternVL3-8B achieve higher performance than Qwen2.5-VL-7B-Instruct, as its better time series understanding capability is shown in Table 1.

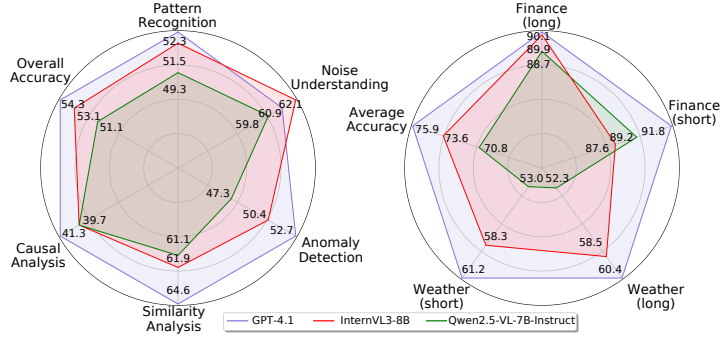


Figure 7: Comparison of multimodal LLMs used to generate time series captions for training TS-REASONER. **Left:** performance on TimeSeriesExam. **Right:** Performance on MTBench.

4.4 CHOICES OF TSFM AND LLM IN TS-REASONER

Different choices of TSFMs. To investigate the performance of TS-REASONER with different TSFMs, we replaced TimesFM (200M) with MOMENT-1-base (200M), a TSFM of the same size, and re-evaluated its performance on the TimeSeriesExam benchmark. The results presented in Table 2 reveal a substantial performance degradation when using MOMENT, with overall accuracy falling from 54.26% to 45.74%. The result is expected because TimesFM outperforms Moment on various time series forecasting benchmarks (Shi et al., 2024; Mulayim et al., 2024). This suggests that TimesFM provides better time series representations, enabling TS-REASONER to better understand and reason for the time series.

Different choices of LLMs. To investigate the scalability and robustness of our approach with different LLM backbones, we evaluate TS-REASONER against across three distinct sizes of the

Table 2: Comparison of TS-REASONER using different TSFMs on the TimeSeriesExam benchmark.

Model	PR	NU	AD	SA	CA	OA
MOMENT	46.90	47.13	41.86	54.87	28.57	45.74
TimesFM	52.29	60.92	52.71	64.60	41.27	54.26

Table 3: Ablation study results of different components in TS-REASONER.

Model	TimeSeriesExam (Cai et al., 2024a)						MTBench (Chen et al., 2025)			
	PR	NU	AD	SA	CA	OA	Finance (long)	Finance (short)	Weather (long)	Weather (short)
TS-REASONER-7B (ours)	52.29	60.92	52.71	64.60	41.27	54.26	90.12	91.85	60.44	61.24
<i>Ablation on Training Data</i>										
- LLM-caption	51.21	56.32	52.71	56.54	36.51	51.25	88.67	89.40	58.24	59.69
- Attributes	52.02	57.47	48.83	62.83	39.68	52.69	89.71	89.20	57.28	59.07
<i>Ablation on Training Stages</i>										
- Stage 1	47.98	54.02	37.98	57.52	30.16	46.92	80.24	83.71	52.88	55.34
- Stage 2	33.42	28.73	13.95	25.67	1.59	25.81	88.07	86.76	56.86	58.60
<i>Ablation on Model Architecture</i>										
- TSFM	51.48	52.87	51.16	63.71	38.09	51.76	89.43	89.70	58.65	60.62

Qwen-2.5-Instruct backbone: 0.5B, 3B, and 7B. The results, shown in Figure 6, confirm that TS-REASONER is both highly effective and robustly scalable. We observe a clear positive scaling law for both TS-REASONER and a baseline. More importantly, TS-REASONER maintains a consistent and significant lead across all models, with Overall Accuracy improvements of +3.15% (29.49% vs. 26.34%), +9.70% (47.71% vs. 38.01%), and +7.61% (54.26% vs. 46.65%) for the 0.5B, 3B, and 7B models, respectively. This demonstrates that our approach performs robustly across different LLM backbones for complex time series reasoning.

4.5 ABLATION STUDIES

To further demonstrate the effectiveness of TS-REASONER, we conduct ablation studies to analyze the impact of individual components. Table 3 summarizes our component-wise ablations from both training and model architecture perspectives:

(1) **Attribute-aware captioning is critical for robust language-timeseries alignment.** (1) Removing our attribute-aware captioning data entirely degrades overall accuracy by 3.01% on TimeSeriesExam and 2% on MTBench. (2) Removing the attributes from the captioning instructions still results in a performance drop of 1.57% and 2.09%, respectively. These results confirm that fine-grained details are vital for learning nuanced temporal patterns. The quality of these captions is confirmed through a quantitative analysis in Appendix C. We provide a qualitative case study in Appendix D to illustrate how attribute-rich captions provide crucial details for model comprehension.

(2) **Absence of any training stage significantly harms the performance.** When removing stage 1 and training only with instruction tuning data, the performance on both benchmarks drops to a large extent due to the weak time series understanding ability. Lack of stage 2 leads to a performance drop by 28.45% on the TimeSeriesExam benchmark and 3.34% on MTBench. The significant performance gap is attributed to the weak ability to understand time series instructions. We observe that removing stage 1 (alignment) leads to a larger drop on MTBench, while removing stage 2 (instruction tuning) causes a larger drop on TimeSeriesExam. This difference comes from the task characteristics of the benchmarks: Removing Stage 1 hurts MTBench more because its tasks require reasoning across both time series and textual news, a skill entirely dependent on the cross-modal alignment learned in Stage 1. In contrast, removing Stage 2 impacts TimeSeriesExam more severely because it directly tests the model’s ability to follow specific analytical commands, which is precisely the skill taught in Stage 2.

(3) **Pretrained TSFM is crucial for effective time series feature extraction.** We remove the pretrained TSFM and repurpose the TS-to-Text adapter to directly project time series patches into the LLM’s embedding space. As shown in Table 3, this modification leads to a performance decrease of 2.50% on the TimeSeriesExam benchmark and 1.31% on MTBench. This result underscores the importance of the TSFM as a powerful temporal feature extractor.

5 CONCLUSION

We introduce TS-REASONER, a framework that advances the ability of LLMs to understand and reason about time series via bridging with the TSFM. To mitigate the intrinsic semantic gap, we further developed an attribute-aware captioning method that enriches time-series alignment data, fostering a more robust alignment. Extensive experiments demonstrate that TS-REASONER substantially outperforms a wide range of baselines on time series understanding and reasoning benchmarks.

ETHICS AND REPRODUCIBILITY STATEMENTS

(1) **Ethics:** Our work aims to improve the time series understanding and reasoning ability, and the experiments conducted in this paper adopt open-source data only for research purposes. It is far from exceeding the understanding of humanity, which does not anticipate any ethical concerns with this work.

(3) **Reproducibility:** Sections 3 and 4 describe our methods and experiments. Further experiment details and results are available in Appendix C. Finally, we include our code repository in the supplemental materials.

REFERENCES

- The claude 3 model family: Opus, sonnet, haiku. URL <https://api.semanticscholar.org/CorpusID:268232499>.
- Openai gpt-4o., 2024. URL <https://openai.com/index/hello-gpt-4o/>.
- Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkinson, Hany Awadalla, Nguyen Bach, Jianmin Bao, Alon Benham, Martin Cai, Vishrav Chaudhary, Congcong Chen, et al. Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras. *arXiv preprint arXiv:2503.01743*, 2025.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, et al. Chronos: Learning the language of time series. *arXiv preprint arXiv:2403.07815*, 2024.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- Yves Bestgen. Measuring lexical diversity in texts: The twofold length problem. *Language Learning*, 74(3):638–671, 2024.
- Yifu Cai, Arjun Choudhry, Mononito Goswami, and Artur Dubrawski. Timeseriesexam: A time series understanding exam. *arXiv preprint arXiv:2410.14752*, 2024a.
- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, et al. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*, 2024b.
- Defu Cao, Furong Jia, Sercan O Arik, Tomas Pfister, Yixiang Zheng, Wen Ye, and Yan Liu. Tempo: Prompt-based generative pre-trained transformer for time series forecasting. *arXiv preprint arXiv:2310.04948*, 2023.
- Jialin Chen, Aosong Feng, Ziyu Zhao, Juan Garza, Gaukhar Nurbek, Cheng Qin, Ali Maatouk, Leandros Tassioulas, Yifeng Gao, and Rex Ying. Mtbench: A multimodal time series benchmark for temporal reasoning and question answering. *arXiv preprint arXiv:2503.16858*, 2025.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. In *European Conference on Computer Vision*, pp. 370–387. Springer, 2024a.
- Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Zhenyu Tang, Li Yuan, et al. Sharegpt4video: Improving video understanding and generation with better captions. *Advances in Neural Information Processing Systems*, 37:19472–19495, 2024b.
- Mouxian Chen, Lefei Shen, Zhuo Li, Xiaoyun Joy Wang, Jianling Sun, and Chenghao Liu. Visions: Visual masked autoencoders are free-lunch zero-shot time series forecasters. *arXiv preprint arXiv:2408.17253*, 2024c.

- Kanzhi Cheng, Wenpo Song, Zheng Ma, Wenhao Zhu, Zixuan Zhu, and Jianbing Zhang. Beyond generic: Enhancing image captioning with real-world knowledge using vision-language pre-training model. In *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 5038–5047, 2023.
- Kanzhi Cheng, Wenpo Song, Jiaxin Fan, Zheng Ma, Qiushi Sun, Fangzhi Xu, Chenyang Yan, Nuo Chen, Jianbing Zhang, and Jiajun Chen. Caparena: Benchmarking and analyzing detailed image captioning in the llm era. *arXiv preprint arXiv:2503.12329*, 2025.
- Juhwan Choi, Junehyoung Kwon, JungMin Yun, Seunguk Yu, and YoungBin Kim. Voldoger: Llm-assisted datasets for domain generalization in vision-language tasks. *arXiv preprint arXiv:2407.19795*, 2024.
- Winnie Chow, Lauren Gardiner, Haraldur T Hallgrímsson, Maxwell A Xu, and Shirley You Ren. Towards time series reasoning with llms. *arXiv preprint arXiv:2409.11376*, 2024.
- Yu-Neng Chuang*, Songchen Li*, Jiayi Yuan*, Guanchu Wang*, Kwei-Herng Lai*, Leisheng Yu, Sirui Ding, Chia-Yuan Chang, Qiaoyu Tan, Daochen Zha, and Xia Hu. Understanding different design choices in training large time series models. *arXiv preprint arXiv:2406.14045*, 2024.
- Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for time-series forecasting. In *Forty-first International Conference on Machine Learning*, 2024.
- Hao Dong, Moru Liu, Kaiyang Zhou, Eleni Chatzi, Juho Kannala, Cyrill Stachniss, and Olga Fink. Advances in multimodal adaptation and generalization: From traditional approaches to foundation models. *arXiv preprint arXiv:2501.18592*, 2025.
- Vijay Ekambaram, Arindam Jati, Pankaj Dayama, Sumanta Mukherjee, Nam Nguyen, Wesley M Gifford, Chandra Reddy, and Jayant Kalagnanam. Tiny time mixers (ttms): Fast pre-trained models for enhanced zero/few-shot forecasting of multivariate time series. *Advances in Neural Information Processing Systems*, 37:74147–74181, 2024.
- Elizabeth Fons, Rachneet Kaur, Soham Palande, Zhen Zeng, Tucker Balch, Manuela Veloso, and Svitlana Vyetenko. Evaluating large language models on time series feature understanding: A comprehensive taxonomy and benchmark. *arXiv preprint arXiv:2404.16563*, 2024.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, et al. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024.
- Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski. Moment: A family of open time-series foundation models. In *International Conference on Machine Learning*, 2024.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew G Wilson. Large language models are zero-shot time series forecasters. *Advances in Neural Information Processing Systems*, 36:19622–19635, 2023.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*, 2023.
- 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 on Artificial Intelligence*, volume 38, pp. 23343–23351, 2024.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.

- Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y Zhang, Xiaoming Shi, Pin-Yu Chen, Yuxuan Liang, Yuan-Fang Li, Shirui Pan, et al. Time-llm: Time series forecasting by reprogramming large language models. *arXiv preprint arXiv:2310.01728*, 2023.
- Zhengfeng Lai, Haotian Zhang, Bowen Zhang, Wentao Wu, Haoping Bai, Aleksei Timofeev, Xianzhi Du, Zhe Gan, Jiulong Shan, Chen-Nee Chuah, et al. Veclip: Improving clip training via visual-enriched captions. In *European Conference on Computer Vision*, pp. 111–127. Springer, 2024.
- Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. Llava-next: Stronger llms supercharge multi-modal capabilities in the wild, May 2024a. URL <https://llava-vl.github.io/blog/2024-05-10-llava-next-stronger-llms/>.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pp. 12888–12900. PMLR, 2022.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pp. 19730–19742. PMLR, 2023a.
- Yuxin Li, Wenchao Chen, Bo Chen, Dongsheng Wang, Long Tian, and Mingyuan Zhou. Prototype-oriented unsupervised anomaly detection for multivariate time series. In *International Conference on Machine Learning*, pp. 19407–19424. PMLR, 2023b.
- Yuxin Li, Wenchao Chen, Xinyue Hu, Bo Chen, Mingyuan Zhou, et al. Transformer-modulated diffusion models for probabilistic multivariate time series forecasting. In *The Twelfth International Conference on Learning Representations*, 2024b.
- Bryan Lim, Sercan Ö Arık, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers for interpretable multi-horizon time series forecasting. *International Journal of Forecasting*, 37(4): 1748–1764, 2021.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024a.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024b.
- Haoxin Liu, Zhiyuan Zhao, Jindong Wang, Harshvardhan Kamarthi, and B Aditya Prakash. Lst-prompt: Large language models as zero-shot time series forecasters by long-short-term prompting. *arXiv preprint arXiv:2402.16132*, 2024c.
- Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X Liu, and Schahram Dustdar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and forecasting. In *International Conference on Learning Representations*, 2021.
- Xu Liu, Juncheng Liu, Gerald Woo, Taha Aksu, Yuxuan Liang, Roger Zimmermann, Chenghao Liu, Silvio Savarese, Caiming Xiong, and Doyen Sahoo. Moirai-moe: Empowering time series foundation models with sparse mixture of experts. *arXiv preprint arXiv:2410.10469*, 2024d.
- Yong Liu, Haoran Zhang, Chenyu Li, Xiangdong Huang, Jianmin Wang, and Mingsheng Long. Timer: Generative pre-trained transformers are large time series models. *arXiv preprint arXiv:2402.02368*, 2024e.
- Yong Liu, Guo Qin, Zhiyuan Shi, Zhi Chen, Caiyin Yang, Xiangdong Huang, Jianmin Wang, and Mingsheng Long. Sundial: A family of highly capable time series foundation models. *arXiv preprint arXiv:2502.00816*, 2025.
- Mike A Merrill, Mingtian Tan, Vinayak Gupta, Tom Hartvigsen, and Tim Althoff. Language models still struggle to zero-shot reason about time series. *arXiv preprint arXiv:2404.11757*, 2024.

- Ozan Baris Mulayim, Pengrui Quan, Liying Han, Xiaomin Ouyang, Dezhi Hong, Mario Bergés, and Mani Srivastava. Are time series foundation models ready to revolutionize predictive building analytics? In *Proceedings of the 11th ACM International Conference on Systems for Energy-Efficient Buildings, Cities, and Transportation*, pp. 169–173, 2024.
- 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.
- Yuqi Nie, Yaxuan Kong, Xiaowen Dong, John M Mulvey, H Vincent Poor, Qingsong Wen, and Stefan Zohren. A survey of large language models for financial applications: Progress, prospects and challenges. *arXiv preprint arXiv:2406.11903*, 2024.
- Zijie Pan, Yushan Jiang, Sahil Garg, Anderson Schneider, Yuriy Nevmyvaka, and Dongjin Song. s^2 IP-LLM: Semantic space informed prompt learning with LLM for time series forecasting. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 39135–39153. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/pan24c.html>.
- Kandukuri Ratna Prakarsha and Gaurav Sharma. Time series signal forecasting using artificial neural networks: An application on ecg signal. *Biomedical Signal Processing and Control*, 76:103705, 2022.
- CLEVELAND RB. Stl: A seasonal-trend decomposition procedure based on loess. *J Off Stat*, 6: 3–73, 1990.
- Xiaoming Shi, Shiyu Wang, Yuqi Nie, Dianqi Li, Zhou Ye, Qingsong Wen, and Ming Jin. Time-moe: Billion-scale time series foundation models with mixture of experts. *arXiv preprint arXiv:2409.16040*, 2024.
- Chenxi Sun, Hongyan Li, Yaliang Li, and Shenda Hong. Test: Text prototype aligned embedding to activate llm’s ability for time series. *arXiv preprint arXiv:2308.08241*, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Chengsen Wang, Qi Qi, Zhongwen Rao, Lujia Pan, Jingyu Wang, and Jianxin Liao. Chronosteer: Bridging large language model and time series foundation model via synthetic data. *arXiv preprint arXiv:2505.10083*, 2025a.
- Chengsen Wang, Qi Qi, Jingyu Wang, Haifeng Sun, Zirui Zhuang, Jinming Wu, Lei Zhang, and Jianxin Liao. Chattime: A unified multimodal time series foundation model bridging numerical and textual data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 12694–12702, 2025b.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in Neural Information Processing Systems*, 34:22419–22430, 2021.
- Zhe Xie, Zeyan Li, Xiao He, Longlong Xu, Xidao Wen, Tieying Zhang, Jianjun Chen, Rui Shi, and Dan Pei. Chatts: Aligning time series with llms via synthetic data for enhanced understanding and reasoning. *arXiv preprint arXiv:2412.03104*, 2024.
- Jingwen Xu, Fei Lyu, and Pong C Yuen. Density-aware temporal attentive step-wise diffusion model for medical time series imputation. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, pp. 2836–2845, 2023.

- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.
- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*, 2024.
- Fangxu Yu, Lai Jiang, Haoqiang Kang, Shibo Hao, and Lianhui Qin. Flow of reasoning: Training llms for divergent problem solving with minimal examples. *arXiv preprint arXiv:2406.05673*, 2024.
- George Zerveas, Srideepika Jayaraman, Dhaval Patel, Anuradha Bhamidipaty, and Carsten Eickhoff. A transformer-based framework for multivariate time series representation learning. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pp. 2114–2124, 2021.
- Haochuan Zhang, Chunhua Yang, Jie Han, Liyang Qin, and Xiaoli Wang. Tempogpt: Enhancing temporal reasoning via quantizing embedding. *arXiv preprint arXiv:2501.07335*, 2025a.
- Haoran Zhang, Yong Liu, Yunzhong Qiu, Haixuan Liu, Zhongyi Pei, Jianmin Wang, and Mingsheng Long. Timesbert: A bert-style foundation model for time series understanding. *arXiv preprint arXiv:2502.21245*, 2025b.
- Siru Zhong, Weilin Ruan, Ming Jin, Huan Li, Qingsong Wen, and Yuxuan Liang. Time-vlm: Exploring multimodal vision-language models for augmented time series forecasting. *arXiv preprint arXiv:2502.04395*, 2025.
- Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. *arXiv preprint arXiv:2201.12740*, 2022.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025.
- Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. Texusgen: A benchmarking platform for text generation models. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pp. 1097–1100, 2018.

A IMPLEMENTATION DETAILS

TS-REASONER uses the Qwen-2.5-7B-Instruct as the LLM backbone across all the experiments with an embedding dimension of 5120, and uses the TimesFM-1.0-200M (Das et al., 2024) as our backbone TSFM with an embedding dimension of 1080. All the parameters of the backbone are finetuned during training. The detailed derivation of these time series embeddings from TimesFM can be found in the Appendix B. All training and inference procedures for TS-REASONER were conducted locally on $8 \times \text{L40s}$ GPUs. Comprehensive training parameters are further detailed in Table 4.

Table 4: Training details of TS-REASONER.

	Stage-1	Stage-2
Patch Size	32	32
Dataset	Captions	Instructions
#Samples	120K	30K
TSFM LLM Backbone	TimesFM-1.0-200M Qwen2.5-7B-Instruct	
Trainable Params.	7.3B	7.3B
Batch Size	64	32
Learning Rate:	1×10^{-5}	2×10^{-5}
Epoch	1	2

B TIMESFM FOR TIME SERIES EMBEDDING

Given a time series $\mathcal{T} \in \mathbb{R}^L$, where L is the length of the time series. We first normalize it to have a mean of zero and a variance of one. We then segment \mathcal{T} into consecutive, non-overlapping patches of fixed length P , resulting in a total of $N = \lfloor L/P \rfloor$ patches. This yields a patched time series $\mathcal{T}_p \in \mathbb{R}^{N \times P}$.

Following the approach of (Das et al., 2024), j -th patch \mathcal{T}_p^j is passed through a residual block to project it into the model dimension. This block is implemented as a two-layer MLP with a skip connection, processing each patch independently. The input token for the j -th patch is computed as:

$$\mathcal{E}_p^j = \text{InputResidualBlock}(\mathcal{T}_p^j) + \text{PE}_j, \quad (5)$$

where PE_j is the position encoding for the j -th patch, as defined in the original transformer (Vaswani et al., 2017). These encoded patch representations are then fed into an M -layer stacked Transformer to produce the final sequence of time series features:

$$\mathcal{Z}_T = \text{StackedTransformer}([\mathcal{E}_p^{(0)}, \mathcal{E}_p^{(1)}, \dots, \mathcal{E}_p^{(N)}]), \quad (6)$$

where $\mathcal{Z}_T \in \mathbb{R}^{N \times d_{ts}}$ and d_{ts} denotes the embedding dimension for each time series patch. Refer to more details of TimesFM in (Das et al., 2024).

C CAPTION ANALYSIS

A critical limitation of synthetic datasets is the risk of models learning spurious correlations from similar templates. To mitigate this, our attribute-aware generation process is designed to produce captions that are both lexically diverse. To quantitatively validate the richness of our approach, we compare it against the template-based method. We evaluate both lexical diversity using the Measure of Textual Lexical Diversity (MTLD) (Bestgen, 2024) and Self-BLEU-4 (Zhu et al., 2018) on a random sample of 1K captions from each dataset. The results presented in Table 5 show that our attribute-aware captions achieve an MTLD score of 133.30, a nearly 3 times increase over the template-based score of 42.95. Furthermore, the Self-BLEU-4 score is almost halved from 0.82 to 0.45. This substantial improvement in lexical diversity confirms that our method generates a significantly more expressive and diverse set of captions, crucial for training robust and generalizable models.

Table 5: Comparison of lexical diversity between template-based pairs and LLM-generated pairs.

Metrics	MTLD \uparrow	Self-BLEU-4 \downarrow
Template-based pairs	42.95	0.82
LLM-generated pairs	133.30	0.45

To ensure comprehensive data coverage, we curated time series with context from a wide range of domains. The distribution of these domains is visualized in Figure 8.

D QUALITATIVE ANALYSIS: A CASE STUDY

To qualitatively evaluate the distinct advantages of our approach, we conduct a case study comparing three methods: (1) our proposed attribute-aware captioning, which leverages visual time series plots and explicit attribute guidance; (2) a basic captioning baseline that operates on visual plots but lacks attribute guidance; and (3) LLM prompted with the raw textual (numerical) time series data. Our analysis, illustrated in Figure 9, yields two key insights.

(i) **Attribute-Aware Captions Provide Semantically Richer Descriptions.** A primary limitation of basic captioning is its tendency to produce superficial, chronological narrations of the data. As shown in Figure 9, the captioner describes the series’ movements (e.g., "the value increases, then decreases sharply") but fails to extract deeper, underlying characteristics. While factually correct, this description omits properties crucial for a comprehensive understanding. In contrast, our attribute-aware captioning enriches this chronological account with critical semantic attributes. It not only captures the temporal dynamics but also identifies and articulates the series’ overall trend, periodicity, and noise level. This multifaceted analysis provides a more holistic understanding of the time series, which is essential for TS-REASONER to conduct reasoning on downstream tasks.

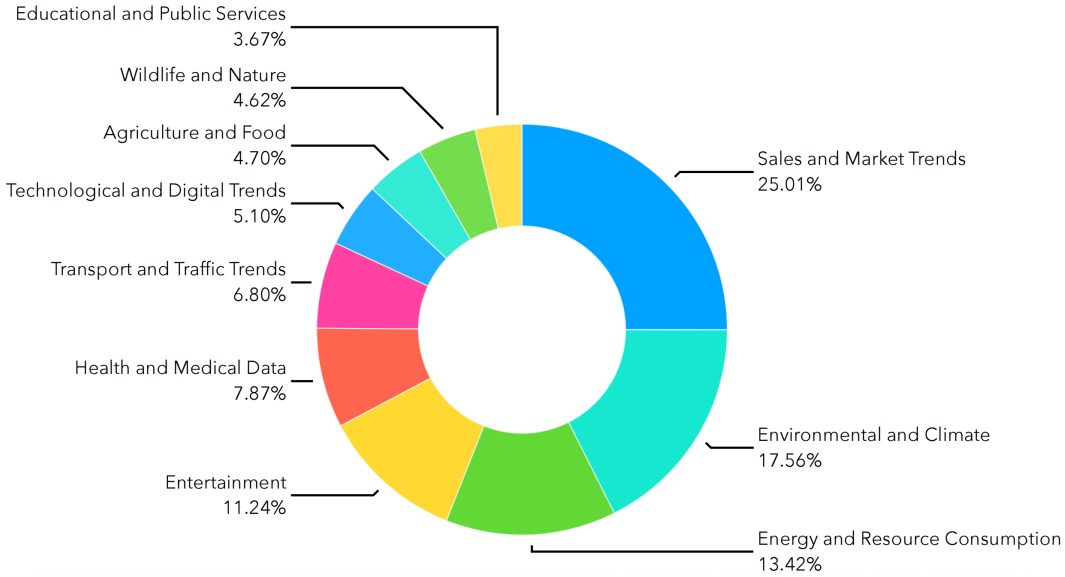


Figure 8: Domain distribution of LLM-generated time series with context.

(ii) **Visual Representation is Crucial for Capturing Global Temporal Patterns.** When comparing our visually-grounded method to an LLM processing raw numerical data, a significant gap emerges in the ability to identify global patterns. The text-based LLM, while capable of discerning local features like high-frequency oscillations or noise within a limited window, consistently fails to recognize the overarching periodicity of the entire series. We hypothesize that this failure stems from the inherent inefficiency of representing long numerical sequences as text. The exceeding length of input may distract the LLM, preventing it from observing the complete pattern. Conversely, a time series plot serves as a highly compressed, holistic representation. It enables the model to perceive the entire sequence as a single input, making global structures like periodicity visually salient and readily identifiable.

E PROMPTS FOR ATTRIBUTE-AWARE CAPTIONING

Below, we detail all the prompts we used. Figure 10 shows our captioning prompts.

F THE USE OF LARGE LANGUAGE MODELS

In this work, LLMs were only used to review and polish small sections of writing

LIMITATION

Our current experimental setup utilizes LLMs with up to 7B parameters, primarily due to existing resource constraints. While we anticipate that our findings will extend to larger models, future work should aim to validate these results with more substantial LLMs. Additionally, exploring alternative architectural designs, such as different encoder architectures, and investigating more effective training strategies for time series text alignment would be valuable avenues for future research. Future research could also explore incorporating RL approaches to elicit reflection behaviors.

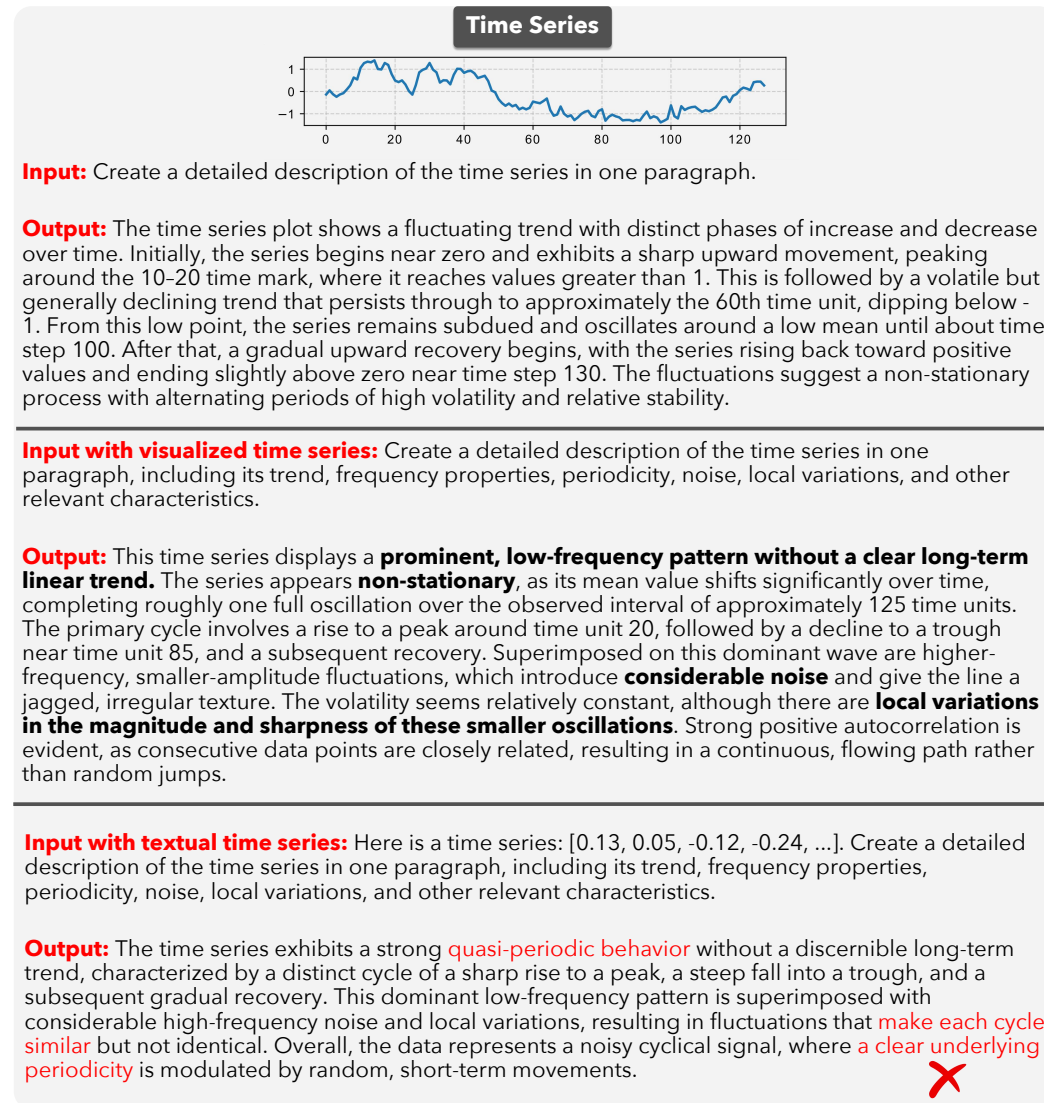


Figure 9: Case study comparing time series captioning methods. Basic visual captioning (row 1) gives a chronological description. Our attribute-aware visual captioning (row 2) provides a richer description including key attributes. Textual time series captioning (row 3) identifies some attributes but fails to capture the periodicity pattern.

Instructions for prompting LLMs to generate time series captions.

1. Write a paragraph that analyzes the time series, covering its local behaviors, noise levels, periodic structures, overall trend, frequency content, and any other characteristics you consider important.
2. Create a detailed description of the time series in one paragraph, including its trend, frequency properties, periodicity, noise, local variations, and other relevant characteristics.
3. Provide a paragraph summarizing the time series characteristics such as noise, periodic patterns, long-term trends, frequency behavior, local anomalies, and any other significant features.
4. Compose a detailed caption describing the frequency characteristics, noise, trends, local variations, periodic structures, and any other meaningful patterns you observe in the time series.
5. Craft a one-paragraph summary of the time series, noting local fluctuations, periodic behavior, frequency features, trend, noise content, and any other insights you find important.
6. Generate a descriptive paragraph detailing the time series' key attributes, including frequency structure, noise patterns, trend direction, local features, periodic elements, and other notable aspects.
7. Give a thorough one-paragraph explanation of the time series, addressing periodicity, noise, frequency components, trend, local variations, and other relevant characteristics.
8. Write a narrative paragraph explaining the time series, focusing on noise, frequency characteristics, periodicity, localized structures, the overall trend, and other important features you identify.
9. Summarize the time series in a paragraph, describing its fluctuations, recurring patterns, noise levels, frequency-domain features, trend direction, and any additional traits you find significant.
10. Develop a paragraph that captures the key features of the time series, such as frequency traits, trend, noise, periodic components, local behaviors, and other characteristics worth noting.
11. Provide a one-paragraph caption analyzing the time series data in terms of noise, trend, periodicity, local features, frequency-related behavior, and any additional characteristics of interest.
12. Create a rich paragraph description of the time series, including its trend, local anomalies, periodic activity, noise artifacts, spectral content, and other important descriptive elements.
13. Write a descriptive paragraph for the time series, highlighting frequency properties, trend behavior, periodic patterns, local structures, noise, and other characteristics you consider relevant.
14. Generate a compact yet thorough paragraph explaining the time series in terms of periodicity, trend movement, noise level, frequency details, local dynamics, and any other key aspects.
15. Construct a one-paragraph analysis of the time series by examining its local variations, noise, trend, periodic elements, frequency spectrum, and other notable features you deem important.
16. Write a summary paragraph that discusses the time series' periodic features, trend behavior, local patterns, noise levels, frequency domain signals, and other characteristics worth mentioning.
17. Create a detailed one-paragraph commentary on the time series that outlines its noise characteristics, periodicity, frequency content, trends, localized behaviors, and other useful insights.
18. Prepare a paragraph-long description of the time series covering its trend, noise, frequency-related traits, local fluctuations, periodic structures, and any additional attributes of note.
19. Offer a one-paragraph interpretation of the time series, highlighting its frequency features, periodic nature, local patterns, noise, trend line, and any other important characteristics you observe.
20. Compose a detailed summary in one paragraph focusing on the time series' periodic behavior, frequency spectrum, localized fluctuations, overall trend, noise, and other relevant descriptive elements.

Figure 10: The list of instructions for attributes-aware time series captioning.