
Position: What Can Large Language Models Tell Us about Time Series Analysis

Ming Jin^{1*} Yifan Zhang^{2*} Wei Chen^{3*} Kexin Zhang⁴ Yuxuan Liang^{†3} Bin Yang⁵
Jindong Wang⁶ Shirui Pan^{†1} Qingsong Wen^{†7}

Abstract

Time series analysis is essential for comprehending the complexities inherent in various real-world systems and applications. Although large language models (LLMs) have recently made significant strides, the development of artificial general intelligence (AGI) equipped with time series analysis capabilities remains in its nascent phase. Most existing time series models heavily rely on domain knowledge and extensive model tuning, predominantly focusing on prediction tasks. In this paper, we argue that current LLMs have the potential to revolutionize time series analysis, thereby promoting efficient decision-making and advancing towards a more universal form of time series analytical intelligence. Such advancement could unlock a wide range of possibilities, including time series modality switching and question answering. We encourage researchers and practitioners to recognize the potential of LLMs in advancing time series analysis and emphasize the need for trust in these related efforts. Furthermore, we detail the seamless integration of time series analysis with existing LLM technologies and outline promising avenues for future research.

1. Introduction

Time series, a fundamental data type for recording dynamic system variable changes, is widely applied across diverse disciplines and applications (Hamilton, 2020; Wen et al., 2022). Its analysis is instrumental in uncovering patterns and relationships over time, thus facilitating the understand-

^{*}Equal contribution ¹Griffith University. ²Chinese Academy of Sciences. ³The Hong Kong University of Science and Technology (Guangzhou). ⁴Zhejiang University. ⁵East China Normal University. ⁶Microsoft Research Asia. ⁷Squirrel AI. Correspondence to: Yuxuan Liang <yuxliang@outlook.com>, Shirui Pan <s.pan@griffith.edu.au>, Qingsong Wen <qingsongedu@gmail.com>.

Proceedings of the 41st International Conference on Machine Learning, Vienna, Austria. PMLR 235, 2024. Copyright 2024 by the author(s).

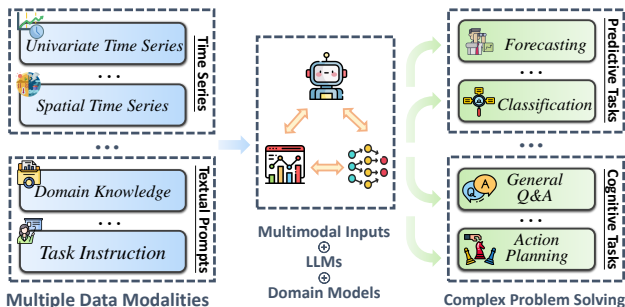


Figure 1: Across a myriad of time series analytical domains, the integration of time series and LLMs demonstrates potential in solving complex real-world problems.

ing of complex real-world systems and supporting informed decision-making. Many real-world dynamic laws, such as financial market fluctuations (Tsay, 2005) and traffic patterns during peak hours (Alghamdi et al., 2019), are fundamentally encapsulated in time series data. In addressing practical scenarios, time series analysis employs methods ranging from traditional statistics (Fuller, 2009) to recent deep learning techniques (Gamboa, 2017; Wen et al., 2021). In the era of sensory artificial intelligence, these domain-specific models efficiently extract meaningful representations for prediction tasks like forecasting and classification. Despite such successes, a notable gap persists between mainstream time series research and the development of artificial general intelligence (AGI) (Bubeck et al., 2023) with time series capabilities to address various problems in a unified manner.

The recent emergence of large language models (LLMs), such as Llama (Touvron et al., 2023a;b) and GPT-4 (Achiam et al., 2023), have swept through and propelled advancements in various interdisciplinary fields (Zhao et al., 2023). Their outstanding zero-shot capabilities (Kojima et al., 2022), along with emerging reasoning and planning abilities (Wang et al., 2024), have garnered increasing attention. However, their focus has primarily been on text sequences. The exploration of extending LLMs’ capabilities to accommodate and process more data modalities, such as images (Zhang et al., 2024a) and graphs (Chen et al., 2024), has begun to receive preliminary attention.

With the integration of LLMs, time series analysis is undergoing significant transformation (Jin et al., 2023b; Liang

et al., 2024). Time series models are conventionally designed for specific tasks, depend heavily on prior domain knowledge and extensive model tuning, lacking assurances of effective updates and validations (Zhou et al., 2023a). Conversely, LLMs hold enormous potential not only to improve prediction performance (Jin et al., 2024) but also to support cross-disciplinary (Yan et al., 2023), interactive (Xue et al., 2023), and interpretative (Gu et al., 2024) analyses. By aligning time series and natural language, large language and foundation/specialistic time series models constitute a new technology paradigm, where the LLM is prompted with both time series and textual instructions. In this paradigm, time series and textual information provide essential contexts, LLMs contribute internal knowledge and reasoning capabilities, and pre-trained time series models offer fundamental pattern recognition assurances. This novel integration is depicted in Figure 1, where the successful amalgamation of these components showcases the potential for a general-purpose, unified system in next-generation time series analysis.

Why This Position Paper? Given the remarkable capabilities emerging in recent research (Jin et al., 2023b; Liang et al., 2024; Zhang et al., 2024e), we believe that the field of time series analysis research is undergoing an exciting transformative moment. Our standpoint is that LLMs can act as the central hub for understanding and advancing the analysis of time series data. Specifically, we present key insights that LLMs can profoundly impact time series analysis in three fundamental ways with their capability boundaries illustrated in Figure 2: **(1) as effective data and model enhancers**, augmenting time series data and existing approaches with enhanced external knowledge and analytical prowess; **(2) as superior predictors**, utilizing their extensive internal knowledge and emerging reasoning abilities to benefit a range of prediction tasks; and **(3) as next-generation agents**, transcending conventional roles to actively engage in and transform time series analysis. We advocate attention to related research and efforts, moving towards more universal intelligent systems for general-purpose time series analysis. To this end, we thoroughly examine relevant literature, present and discuss potential formulations of LLM-centric time series analysis to bridge the gap between the two. We also identify and outline prospective research opportunities and challenges, calling for greater commitment and exploration in this promising interdisciplinary field.

Contributions: The contributions of this work can be summarized in three aspects: **(1) offering new perspectives**. We articulate our stance on LLM-centric time series analysis, outlining the potential synergies between LLMs and time series analytical models. This underscores the need for increased research focus and dedication in this area; **(2) systematic review and categorization**. We meticulously exam-

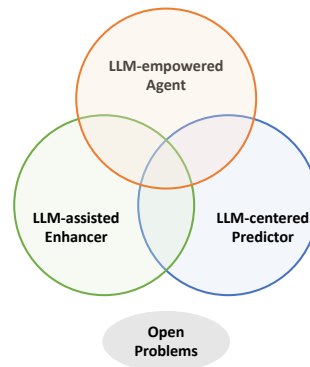


Figure 2: Task-solving capability boundaries on the roles of LLMs for time series analysis: as data/model enhancers, effective predictors, or next-generation agents.

ine existing preliminary work and present a clear roadmap, highlighting three potential integration forms of LLMs and time series analysis; **(3) identifying future opportunities**. We explore and articulate areas that current research has not yet addressed, presenting promising directions for future investigations in this evolving interdisciplinary field.

2. Background

2.1. Time Series Analysis

Data Modality. Time series data, comprising sequential observations over time, can be either regularly or irregularly sampled, with the latter often leading to missing values. This data falls into two main categories: *univariate* and *multivariate*. Univariate time series consist of single scalar observations over time, represented as $X = \{x_1, x_2, \dots, x_T\} \in \mathbb{R}^T$. Multivariate time series, on the other hand, involve N -dimensional vector observations, denoted as $X \in \mathbb{R}^{N \times T}$. In complex real-world systems, multivariate time series often exhibit intricate spatial dependencies in addition to temporal factors. This has led to some recent studies modeling them as graphs (Jin et al., 2023a), also referred to as *spatial time series*. In this approach, a time series is conceptualized as a sequence of graph snapshots, $\mathcal{G} = \{\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_T\}$, with each $G_t = (A_t, X_t)$ representing an attributed graph characterized by an adjacency matrix $A_t \in \mathbb{R}^{N \times N}$ and node features $X_t \in \mathbb{R}^{N \times D}$.

Analytical Tasks. Time series analysis is crucial for deriving insights from data, with recent deep learning advancements spurring a rise in neural network-based methods (Wen et al., 2023). These methods focus on modeling complex inter-temporal and/or inter-variable relationships in time series (Zhang et al., 2024b; Jin et al., 2023b), aiding in tasks like forecasting, classification, anomaly detection, and imputation. Forecasting predicts future values, classification categorizes series by patterns, anomaly detection identifies

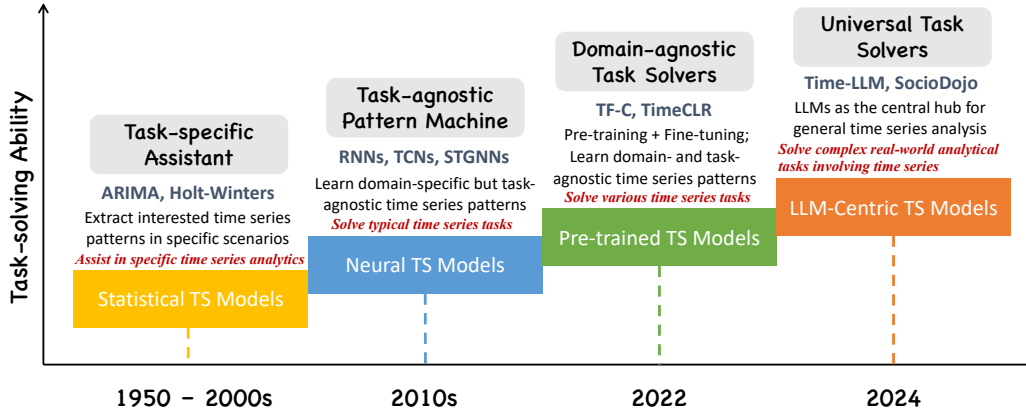


Figure 3: A roadmap of time series analysis delineating four generations of models based on their task-solving capabilities.

anomalous events, and imputation estimates missing data. Emerging research has also shown promise in time series modality switching and question answering (Xue & Salim, 2023; Jin et al., 2024; Yang et al., 2022a). These novel approaches highlight the potential for cross-disciplinary, interactive, and interpretative advancements in time series analytics. Such advancements open a realm of possibilities in practical applications, such as (zero-shot) medical question answering (Yu et al., 2023a; Oh et al., 2024) and intelligent traffic agents (Da et al., 2024; Lai et al., 2023).

2.2. Large Language Models

Basic Concept. Large language models typically refer to transformer-based pre-trained language models (PLMs) with billions or more parameters. The scaling of PLMs, both in terms of model and data size, has been found to enhance model performance across various downstream tasks (Zhao et al., 2023). These models such as GPT-4 (Achiam et al., 2023), PaLM (Chowdhery et al., 2023), and Llama (Touvron et al., 2023a), undergo extensive pre-training on extensive text corpora, enabling them to acquire wide-ranging knowledge and problem-solving capabilities for diverse NLP tasks. Technically, language modeling (LM) is a fundamental pre-training task in LLMs and a key method for advancing machine language intelligence. The primary objective of LM is to model the probability of generating word sequences, encompassing both non-autoregressive and autoregressive language model categories. Autoregressive models, like the GPT series (Bubeck et al., 2023), predict the next token y based on a given context sequence X , trained by maximizing the probability of the token sequence given the context:

$$P(y | X) = \prod_{t=1}^T P(y_t | x_1, x_2, \dots, x_{t-1}), \quad (1)$$

where T represents the sequence length. Through this, the model achieves intelligent compression and language generation in an autoregressive manner.

Emergent Abilities of LLMs. Large language models exhibit emergent abilities that set them apart from traditional neural networks. These abilities, present in large models but not in smaller ones, are a significant aspect of LLMs (Wei et al., 2022a). Three key emergent abilities of LLMs include: (1) *in-context learning (ICL)*, introduced by GPT-3 (Brown et al., 2020), allowing LLMs to generate relevant outputs for new instances using instructions and examples without additional training; (2) *instruction following*, where LLMs, through instruction tuning, excel at novel tasks presented in an instructional format, enhancing their generalization (Victor et al., 2022); (3) *step-by-step reasoning*, where LLMs use strategies like chain-of-thought (CoT) (Wei et al., 2022b) or other prompting strategies (Yao et al., 2024; Besta et al., 2023) to address complex tasks requiring multiple reasoning steps.

2.3. Research Roadmap

Time series analytical model development spans four generations: (1) *statistical models*, (2) *deep neural networks*, (3) *pre-trained models*, and (4) *LLM-centric models*, as shown in Figure 3. This categorization hinges on the evolving task-solving capabilities of each model generation. Traditional analytics relied on statistical models like ARIMA (Shumway et al., 2017) and Holt-Winters (Kalekar et al., 2004), optimized for small-scale data and based on heuristics like stationarity and seasonality (Hamilton, 2020). These models assumed past trends would continue into the future. Deep neural networks, like recurrent and temporal convolution neural networks (Gamboa, 2017), processed larger, complex datasets, capturing non-linear and long-term dependencies without heavy reliance on prior knowledge, thus transforming predictive time series analysis. Recent research like TimeCLR (Yeh et al., 2023) introduced pre-training on diverse, large-scale time series data, allowing fine-tuning for specific tasks with relatively smaller data samples (Jin et al., 2023b), reducing the time and resources required for model

training. This allows for the application of sophisticated models in scenarios where collecting large-scale time series data is challenging. Despite the successes of previous generations, we posit that the emergence of LLMs is set to revolutionize time series analysis, shifting it from predictive to general intelligence. LLM-centric models, processing both language instructions and time series (Jin et al., 2024; Cheng & Chin, 2024), extend capabilities to general question answering, interpretable predictions, and complex reasoning, moving beyond conventional predictive analytics.

3. LLM-assisted Enhancer for Time Series

The role of LLM-assisted enhancers delves into whether LLMs can augment our understanding of time series data and extend the knowledge of existing models. Numerous methods have been devised to address temporal data, but the vast internal knowledge and reasoning capabilities of LLMs may significantly enhance both data understanding and model performance; thus, we intuitively distinguish LLM-assisted enhancers from *data* and *model* perspectives.

3.1. Data-based Enhancer

LLM-assisted enhancers not only enhance data interpretability but also provide supplementary improvements, facilitating a more thorough understanding and effective use of time series data. For interpretability, LLMs offer textual descriptions and summaries, helping to understand patterns and anomalies in time series data. Examples include LLM-MPE (Liang et al., 2023) for human mobility data, SignalGPT (Liu et al., 2023a) for biological signals, and Insight Miner (Zhang et al., 2023c) for trend mining. Additionally, AmicroN (Chatterjee et al., 2023) and SST (Ghosh et al., 2023) use LLMs for detailed sensor and spatial time series analysis. Supplementary enhancements involve integrating diverse data sources, enriching time series data context and improving model robustness, as explored in (Yu et al., 2023b) and (Fatouros et al., 2024) for financial decision-making. Such enhancements help improve domain models' inherent capabilities and make them more robust.

3.2. Model-based Enhancer

Model-based enhancers aim to augment time series models by addressing their limitations in external knowledge and domain-specific contexts. Transferring knowledge from LLMs boosts the performance of domain models in handling complex tasks. Such approaches often employ a dual-tower model, like those in (Qiu et al., 2023b; Li et al., 2024), use frozen LLMs for electrocardiogram (ECG) analysis. Some methods further utilize contrastive learning to achieve certain alignments. For example, IMU2CLIP (Moon et al., 2023) aligns text and video with sensor data, while STLLM (Zhang et al., 2024c) enhances spatial time series

prediction. Another branch utilizes prompting techniques to harness the inferential decision-making capability of LLMs. For instance, TrafficGPT (Zhang et al., 2024d) exemplifies decision analysis, integrating traffic models with LLMs for user-tailored solutions, offering detailed insights to enhance system interpretability.

3.3. Discussion

LLM-assisted enhancers effectively address the inherent sparsity and noise characteristics of time series data while also providing existing time series models with enhanced external knowledge and analytical capabilities. Time series, perhaps more than other modalities, benefits significantly from data enhancements due to its potentially lower information density. For example, comparing raw data volumes, video data inherently contain vastly more information per second than typical time series data such as audio. This disparity underscores the importance of time series data augmentation in many use cases, where LLMs can leverage both internal and external knowledge, along with their inherent reasoning capabilities, to address this challenge. On the other hand, the development of foundation models for time series is still in its infancy, presenting an opportunity for LLMs to extend the knowledge boundaries of existing time series models. Moreover, this technology is plug-and-play, enabling flexible assistance for real-world time series data and model challenges. However, a notable hurdle is that using LLM as an enhancer introduces significant time and cost overheads when dealing with large-scale datasets. In addition, the inherent diversity and range of application scenarios in time series data add layers of complexity to the creation of universally effective LLM-assisted enhancers.

Our position: *LLM-assisted enhancers represent a promising avenue for augmenting time series data and models, meriting further exploration. Future directions should focus on developing efficient, accountable, and universally adaptable plug-and-play solutions that effectively address practical challenges, such as data sparsity and noise, while also considering the time and cost efficiencies for large-scale dataset applications.*

4. LLM-centered Predictor for Time Series

LLM-centered predictors utilize the extensive knowledge within LLMs for diverse time series tasks such as prediction and anomaly detection. Adapting LLMs to time series data involves unique challenges such as differences in data sampling and information completeness. In the following discussion, approaches are categorized into *tuning-based* and *non-tuning-based* methods based on whether access to LLM parameters, primarily focusing on building general or domain-specific time series models.

4.1. Tuning-based Predictor

Tuning-based predictors use accessible LLM parameters, typically involving patching and tokenizing numerical signals and related text data, followed by fine-tuning for time series tasks. Figure 4(a) shows this process: (1) with a $\text{Patching}(\cdot)$ operation (Nie et al., 2022), a time series is chunked to form patch-based tokens \mathcal{X}_{inp} . An additional option is to perform $\text{Tokenizer}(\cdot)$ operation on time series-related text data to form text sequence tokens \mathcal{T}_{inp} ; (2) time series patches (and optional text tokens) are fed into the LLM with accessible parameters; (3) an extra task layer, denoted as $\text{Task}(\cdot)$, is finally introduced to perform different analysis tasks with the instruction prompt P (not shown in the figure). This process is formulated below:

$$\begin{aligned} \text{Pre-processing:} \quad & \mathcal{X}_{inp} = \text{Patching}(\mathcal{X}), \\ & \mathcal{T}_{inp} = \text{Tokenizer}(\mathcal{T}), \quad (2) \\ \text{Analysis:} \quad & \hat{Y} = \text{Task}(f_{LLM}^{\Delta}(\mathcal{X}_{inp}, \mathcal{T}_{inp}, P)), \end{aligned}$$

where \mathcal{X} and \mathcal{T} denote the set of time series samples and related text samples, respectively. These two (the latter is optional) are fed together into LLM f_{LLM}^{Δ} with partial unfreezing or additional adapter layers to predict label \hat{Y} .

Adapting out-of-box LLMs directly to raw time series numerical signals for downstream time series analysis tasks is often counterintuitive due to the inherent modality gap between text and time series data. Nevertheless, OFA (Zhou et al., 2023a) and similar studies found that LLMs, even when frozen, can perform comparably in time series tasks due to the self-attention mechanism’s universality. Others, like GATGPT (Chen et al., 2023b) and ST-LLM (Liu et al., 2024a), applied these findings to spatial-temporal data, while UniTime (Liu et al., 2024b) used manual instructions for domain identification. This allows them to handle time series data with different characteristics and distinguish between different domains.

However, the above methods all require modifications that disrupt the parameters of the original LLMs, potentially leading to catastrophic forgetting. In contrast, another line of work, inspired by this, aims to avoid this by introducing additional lightweight adaptation layers. Time-LLM (Jin et al., 2024) uses text data as a prompt prefix and reprograms input time series into language space, enhancing LLM’s performance in various forecasting scenarios. TEST (Sun et al., 2024) tackles inconsistent embedding spaces by constructing an encoder for time series data, employing alignment contrasts and soft prompts for efficient fine-tuning with frozen LLMs. TEMPO (Cao et al., 2023) combines seasonal and trend decompositions with frozen LLMs, using prompt pooling to address distribution changes in forecasting non-stationary time series.

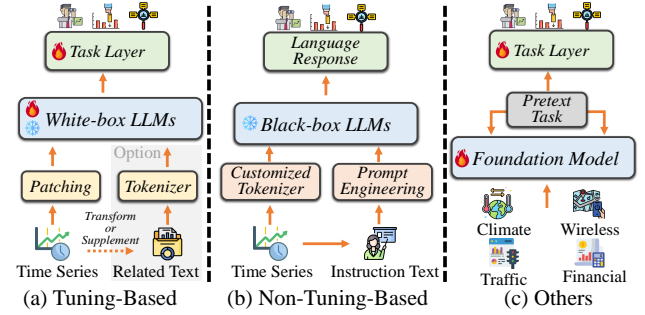


Figure 4: Categories of LLM-centered predictor.

4.2. Non-tuning-based Predictor

Non-tuning-based predictors, suitable for closed-source models, involve preprocessing time series data to fit the input spaces of LLMs. As Figure 4(b) illustrates, this typically involves two steps: (1) preprocessing raw time series, including optional operations such as $\text{Template}(\cdot)$ and $\text{Customized Tokenizer}(\cdot)$; (2) feeding the processed inputs \mathcal{X}_{inp} into the LLM to obtain responses. A $\text{Parse}(\cdot)$ function is then employed to retrieve prediction labels. This process is formulated below:

$$\begin{aligned} \text{Pre-processing:} \quad & \mathcal{X}_{inp} = \text{Template}(\mathcal{X}, P), \\ & \text{or } \mathcal{X}_{inp} = \text{Tokenizer}(\mathcal{X}), \quad (3) \\ \text{Analysis:} \quad & \hat{Y} = \text{Parse}(f_{LLM}^{\Delta}(\mathcal{X}_{inp})), \end{aligned}$$

where P represents the instruction prompt for the current analysis task, and f_{LLM}^{Δ} denotes the black-box LLM model.

(Spathis & Kawsar, 2023) initially noted that LLM tokenizers, not designed for numerical values, separate continuous values and ignore their temporal relationships. They suggested using lightweight embedding layers and prompt engineering as solutions. Following this, LLMTime (Gruber et al., 2023) introduced a novel tokenization approach, converting tokens into flexible continuous values, enabling non-tuned LLMs to match or exceed zero-shot prediction performance in domain-specific models. This success is attributed to LLMs’ ability to represent multimodal distributions. Using in-context learning, evaluations were performed in tasks like sequence transformation and completion. (Mirchandani et al., 2023) suggested that LLMs’ capacity to handle abstract patterns positions them as foundational general pattern machines. This has led to applying LLMs in areas like human mobility mining (Wang et al., 2023b; Zhang et al., 2023d), financial forecasting (Lopez-Lira & Tang, 2023), and health prediction (Kim et al., 2024).

4.3. Others

Beyond the previously discussed methods, another significant approach in temporal analysis involves building foundation models from scratch, as shown in Figure 4(c). This

approach focuses on creating large, scalable models, both generic and domain-specific, aiming to emulate the *scaling law* (Kaplan et al., 2020) of LLMs. More details can be found in (Liang et al., 2024), as this branch of methods is somewhat peripheral to our primary positions in this paper.

4.4. Discussion

LLM-centric predictors have advanced significantly in time series analysis, often outperforming many domain-specific models in few-shot and zero-shot scenarios. Tuning-based methods, with their adjustable parameters, generally show better performance and adaptability to specific domains. However, they are prone to catastrophic forgetting and involve high training costs due to parameter modification. While adapter layers have somewhat alleviated this issue, the challenge of expensive training persists. Conversely, non-tuning methods, offering text-based predictions, depend heavily on manual prompt engineering, and their prediction stability is not always reliable. Additionally, building foundational time series models from scratch involves balancing high development costs against their applicability. Therefore, further refinement is needed to address these challenges in LLM-centric predictors.

Our position: *LLM-centric predictors, though burgeoning in time series analysis, are still in their infancy and warrant deeper consideration. Our position posits a crucial hypothesis that LLMs excel at processing time series tasks. Future advancements may not only build upon but also involve pre-trained time series models. By harnessing the unique capabilities of LLMs, these advancements can further reduce tuning costs and improve prediction stability and reliability.*

5. LLM-empowered Agent for Time Series

As demonstrated in the previous section, tuning-based approaches in time series utilize LLMs as robust model checkpoints, attempting to adjust certain parameters for specific domain applications. However, this approach often sacrifices the interactive capabilities of LLMs and may not fully exploit the benefits offered by LLMs, such as in-context learning or chain-of-thought. On the other hand, non-tuning approaches, integrating time series data into textual formats or developing specialized tokenizers, face limitations due to LLMs’ primary training on linguistic data, hindering their comprehension of complex time series patterns not easily captured in language. Addressing these challenges, there are limited works that directly leverage LLMs as time series agents for general-purpose analysis and problem-solving. We first endeavor to provide an overview of such approaches across various modalities in Appendix B, aiming to delineate strategies for constructing a robust general-purpose

time series analysis agent.

In the subsequent section, we employ prompt engineering techniques to compel LLMs to assist in executing basic time series analytical tasks. Our demonstration reveals that LLMs undeniably possess the potential to function as time series agents. Nevertheless, their proficiency is constrained when it comes to comprehending intricate time series data, leading to the generation of hallucinatory outputs. Ultimately, we identify and discuss promising avenues that can empower us to develop more robust and reliable general-purpose single-agent and multi-agent time series systems.

5.1. Empirical Insights: LLMs as Time Series Analysts

This subsection presents experiments evaluating the LLM’s zero-shot capability as an agent for human interaction and time series data analysis. We utilize the HAR (Anguita et al., 2013) database, derived from recordings of 30 study participants engaged in activities of daily living (ADL) while carrying a waist-mounted smartphone equipped with inertial sensors. The end goal is to classify activities into four categories (*Stand, Sit, Lay, Walk*), with ten instances per class for evaluation. The prompts used for GPT-3.5 are illustrated in Figure 8, and the classification confusion matrix is presented in Figure 5. Our key observations include:



Figure 5: Confusion matrix of HAR classification.

LLM as Effective Analytical Agent. The experiments demonstrate that the LLM serves adeptly as an agent for human interaction and time series data analysis, producing accurate predictions as shown in Figure 8. Notably, all instances with label *Stand* were correctly classified, underscoring the LLMs’ proficiency in zero-shot tasks. The models exhibit a profound understanding of common-sense behaviors, encompassing various labels in time series classification, anomaly detection, and skillful application of data augmentation (Figure 9).

Interpretability and Truthfulness. This single-agent system prioritizes high interpretability and truthfulness, allowing users to inquire about the reasons behind their decisions with confidence. The intrinsic classification reasoning is articulated in natural language, fostering a user-friendly interaction.

Limitations in Understanding Complex Patterns. Despite their capabilities, current LLMs show limitations in comprehending complex time series patterns. When faced with complex queries, they may initially refuse to provide

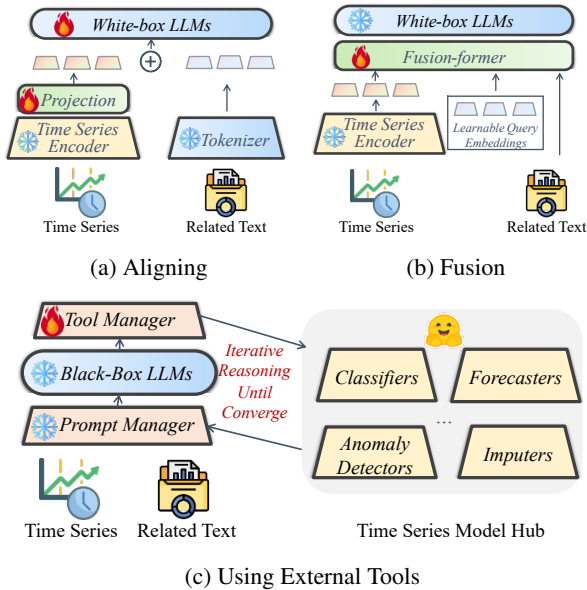


Figure 6: Different directions for incorporating time series knowledge to LLMs.

answers, citing the lack of access to detailed information about the underlying classification algorithm.

Bias and Task Preferences. LLMs display a bias towards the training language distributions, exhibiting a strong preference for specific tasks. In Figure 8, instances of *Lay* are consistently misclassified as *Sit* and *Stand*, with better performance observed for *Sit* and *Stand*.

Hallucination. LLMs are susceptible to hallucination problem, generating reasonable but false answers. For instance, in Figure 9, augmented data is merely a copy of given instances, although the model knows how to apply data augmentation: *These instances continue the hourly trend of oil temperature and power load features, maintaining the structure and characteristics of the provided dataset.* Subsequent inquiries into the misclassification in Figure 5, particularly regarding why LLMs classify *Lay* instances as *Sit* and *Stand*, elicit seemingly plausible justifications (see Table 1). However, these justifications expose the model’s inclination to fabricate explanations.

5.2. Key Lessons for Advancing Time Series Agents

General-purpose single-agent (and multi-agent) systems remain a missing piece in the puzzle of modern time series analysis. In light of the empirical insights from earlier experiments, it is apparent that LLMs, when serving as advanced single-agent time series systems, exhibit notable limitations when dealing with questions about data distribution and specific features. Their responses often show a reliance on requesting additional information or highlight an inability to provide accurate justifications without access to the

underlying model or specific data details.

To surmount such limitations and develop practical time series agents built upon LLMs, it becomes paramount to seamlessly integrate time series knowledge into LLMs. Drawing inspiration from studies that have successfully injected domain-specific knowledge into LLMs (Wang et al., 2023a; Liu et al., 2023b; Wu et al., 2023; Schick et al., 2024), we propose several research directions. These include innovative methods to enhance LLM-based single-agent systems’ proficiency in time series analysis by endowing them with a deep understanding of temporal patterns and relevant contextual information.

- **Aligning Time Series Features with Language Model Representations** (Figure 6a). Explicitly aligning time series features with pre-trained language model representations can potentially enhance the model’s understanding of temporal patterns. This alignment may involve mapping specific features to the corresponding linguistic elements within the model.
- **Fusing Text Embeddings and Time Series Features** (Figure 6b). Exploring the fusion of text embeddings and time series features in a format optimized for LLMs is a promising avenue. This fusion aims to create a representation that leverages the strengths of LLMs in natural language processing while accommodating the intricacies of time series data.
- **Teaching LLMs to Utilize External Tools** (Figure 6c). The goal here is to instruct the LLM to identify the appropriate pre-trained time series models or analytical tools from an external “toolbox” and guide their usage based on user queries. The time series knowledge resides within this external model hub, while the LLM assumes the role of a high-level agent, responsible for orchestrating their utilization and facilitating interaction with users.

Differentiating from approaches like model repurposing or fine-tuning on specific tasks, the focus of future research should be on harnessing the inherent zero-shot capabilities of LLMs for general pattern manipulation. Establishing a framework that facilitates seamless interaction between users and LLM agents for solving general time series problems through in-context learning is an exciting direction.

5.3. Exploring Alternative Research Avenues

Addressing the urgent and crucial need to enhance the capabilities of time series agents built upon LLMs, we recognize that incorporating time series knowledge is a pivotal direction. Concurrently, mitigating risks associated with such agents is equally paramount. In this regard, we pinpoint key challenges and suggest potential directions to boost both the reliability and effectiveness of our time series agents.

Hallucination, a recurring challenge in various foundational models (Zhou et al., 2023b; Rawte et al., 2023; Li et al., 2023), is a significant concern in deploying LLM-based agent systems for time series analysis, as our experiments have shown. Addressing this issue typically involves two methods: identifying reliable prompts (Vu et al., 2023; Madaan et al., 2024) and fine-tuning models with dependable instruction datasets (Tian et al., 2023; Zhang et al., 2023a). However, these approaches require substantial human effort, posing scalability and efficiency challenges. Some initiatives integrate domain-specific knowledge into ICL prompts (Da et al., 2023; Yang et al., 2022b) and construct instruction datasets for specific domains (Liu et al., 2023b; Ge et al., 2023), but the best formats for instructions or prompts for effective time series analysis are still unclear. Developing guidelines for crafting impactful instructions in time series analysis is a promising area for future research.

Multi-agent time series systems is another promising direction composed of multiple interacting time series agents. These systems can leverage the strengths of individual agents, each specialized in different aspects of time series analysis, to provide more comprehensive and accurate results (Cheng et al., 2024). For instance, one agent might excel in identifying patterns and trends, while another focuses on anomaly detection or forecasting. The collaboration among these agents can lead to more robust solutions, particularly in complex scenarios where single-agent time series systems fall short.

Ongoing concerns about **aligning time series agents with human preferences** (Lee et al., 2023), such as generating helpful and harmless content (Bai et al., 2022), highlight the need for more robust and trustworthy agents. Additionally, the internet’s constant evolution, adding petabytes of new data daily (Wenzek et al., 2020), accentuates the importance of handling concept drift in time series data (Tsymbal, 2004), where future data may differ from past patterns. Addressing this challenge requires **enabling agents to continually acquire new knowledge** (Garg et al., 2023) or adopting lifelong learning without costly retraining.

Our Position: *LLMs hold promise as agent systems for various time series applications, yet they encounter challenges such as occasional inaccuracies and hallucination. Enhancing their reliability requires effective instruction guidelines and domain-specific knowledge integration. Aligning with human preferences and adapting to evolving time series are crucial for maximizing their capabilities and minimizing risks. Our vision is to develop robust LLM-empowered agents capable of handling time series complexities, including exploring multi-agent systems for improved performance and reliability.*

6. Further Discussion

Our perspectives initiate ongoing discussion. Acknowledging diverse views and potential curiosities regarding LLM-centric time series analysis, we objectively examine several alternate viewpoints:

Accountability and Transparency. LLMs remain somewhat enigmatic, raising fundamental questions about their capabilities, mechanisms, and efficiency levels, especially in recent studies like PromptCast (Xue & Salim, 2023). We advocate for understanding underlying mechanisms (Gruver et al., 2023) and establishing transparent development and evaluation frameworks, including consistent model reporting and clear explanations of internal processes and outputs (Liao & Vaughan, 2024).

Privacy and Security. LLM-centric time series analysis poses significant privacy and security challenges due to the sensitivity of industrial time series data. LLMs are known to sometimes memorize segments of their training data, which may include private information (Peris et al., 2023). Measures against threats like data leakage and misuse are crucial, along with ethical guidelines and regulatory frameworks to ensure responsible and secure application (Zhuo et al., 2023).

Environmental and Computational Costs. Critics highlight the environmental and computational costs of LLM-centric time series analysis, suggesting optimization opportunities in LLM development and exploring more efficient alignment and inference strategies, especially for handling tokenized high-precision numerical data.

7. Conclusion

This paper aims to draw the attention of researchers and practitioners to the potential of LLMs in advancing time series analysis and to underscore the importance of trust in these endeavors. Our key position is that LLMs can serve as the central hub for understanding and advancing time series analysis, steering towards more universal intelligent systems for general-purpose analysis, whether as enhancers, predictors, or agents. To substantiate our positions, we have reviewed relevant literature, exploring and debating possible directions towards LLM-centric time series analysis to bridge existing gaps.

Our objective is to amplify the awareness of this area within the research community and pinpoint avenues for future investigations. While our positions may attract both agreement and dissent, the primary purpose of this paper is to spark discussion on this interdisciplinary topic. If it serves to shift the discourse within the community, it will have achieved its intended objective.

Acknowledgements

This material is based on research partially sponsored by the CSIRO – National Science Foundation (US) AI Research Collaboration Program. S. Pan was supported in part by the Australian Research Council (ARC) under grants FT210100097 and DP240101547. Y. Liang was supported by the Guangzhou-HKUST(GZ) Joint Funding Program (No. 2024A03J0620).

Impact Statement

This position paper aims to reshape perspectives within the time series analysis community by exploring the untapped potential of LLMs. We advocate a shift towards integrating LLMs with time series analysis, proposing a future where decision-making and analytical intelligence are significantly enhanced through this synergy. While our work primarily contributes to academic discourse and research directions, it also touches upon potential societal impacts, particularly in decision-making processes across various industries. Ethically, the responsible and transparent use of LLMs in time series analysis is emphasized, highlighting the need for trust and understanding in their capabilities. While we foresee no immediate societal consequences requiring specific emphasis, we acknowledge the importance of ongoing ethical considerations and the potential for future societal impacts as this interdisciplinary field evolves.

References

- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Alghamdi, T., Elgazzar, K., Bayoumi, M., Sharaf, T., and Shah, S. Forecasting traffic congestion using arima modeling. In *2019 15th international wireless communications & mobile computing conference (IWCMC)*, pp. 1227–1232. IEEE, 2019.
- Anguita, D., Ghio, A., Oneto, L., Parra, X., Reyes-Ortiz, J. L., et al. A public domain dataset for human activity recognition using smartphones. In *Esann*, volume 3, pp. 3, 2013.
- Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Gianinazzi, L., Gajda, J., Lehmann, T., Podstawski, M., Niewiadomski, H., Nyczyk, P., et al. Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687*, 2023.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y., Lundberg, S., et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- Cao, D., Jia, F., Arik, S. O., Pfister, T., Zheng, Y., Ye, W., and Liu, Y. Tempo: Prompt-based generative pre-trained transformer for time series forecasting. In *The Twelfth International Conference on Learning Representations*, 2023.
- Chang, C., Peng, W.-C., and Chen, T.-F. Llm4ts: Two-stage fine-tuning for time-series forecasting with pre-trained llms. *arXiv preprint arXiv:2308.08469*, 2023.
- Chatterjee, S., Mitra, B., and Chakraborty, S. Amicron: A framework for generating annotations for human activity recognition with granular micro-activities. *arXiv preprint arXiv:2306.13149*, 2023.
- Chen, S., Long, G., Shen, T., and Jiang, J. Prompt federated learning for weather forecasting: toward foundation models on meteorological data. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, pp. 3532–3540, 2023a.
- Chen, Y., Wang, X., and Xu, G. Gatgpt: A pre-trained large language model with graph attention network for spatiotemporal imputation. *arXiv preprint arXiv:2311.14332*, 2023b.
- Chen, Z., Mao, H., Li, H., Jin, W., Wen, H., Wei, X., Wang, S., Yin, D., Fan, W., Liu, H., et al. Exploring the potential of large language models (llms) in learning on graphs. *ACM SIGKDD Explorations Newsletter*, 25(2): 42–61, 2024.
- Cheng, J. and Chin, P. Sociodojo: Building lifelong analytical agents with real-world text and time series. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=s9z0HzWJJP>.
- Cheng, Y., Zhang, C., Zhang, Z., Meng, X., Hong, S., Li, W., Wang, Z., Wang, Z., Yin, F., Zhao, J., et al. Exploring large language model based intelligent agents: Definitions, methods, and prospects. *arXiv preprint arXiv:2401.03428*, 2024.

- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- Da, L., Gao, M., Mei, H., and Wei, H. Llm powered sim-to-real transfer for traffic signal control. *arXiv preprint arXiv:2308.14284*, 2023.
- Da, L., Liou, K., Chen, T., Zhou, X., Luo, X., Yang, Y., and Wei, H. Open-ti: Open traffic intelligence with augmented language model. *International Journal of Machine Learning and Cybernetics*, pp. 1–26, 2024.
- Das, A., Kong, W., Sen, R., and Zhou, Y. A decoder-only foundation model for time-series forecasting. In *the 41st International Conference on Machine Learning*, 2024.
- Ekambaram, V., Jati, A., Nguyen, N. H., Dayama, P., Reddy, C., Gifford, W. M., and Kalagnanam, J. Ttms: Fast multi-level tiny time mixers for improved zero-shot and few-shot forecasting of multivariate time series. *arXiv preprint arXiv:2401.03955*, 2024.
- Fatouros, G., Metaxas, K., Soldatos, J., and Kyriazis, D. Can large language models beat wall street? unveiling the potential of ai in stock selection. *arXiv preprint arXiv:2401.03737*, 2024.
- Fuller, W. A. *Introduction to statistical time series*. John Wiley & Sons, 2009.
- Gamboa, J. C. B. Deep learning for time-series analysis. *arXiv preprint arXiv:1701.01887*, 2017.
- Garg, S., Farajtabar, M., Pouransari, H., Vemulapalli, R., Mehta, S., Tuzel, O., Shankar, V., and Faghri, F. Tic-clip: Continual training of clip models. In *NeurIPS 2023 Workshop on Distribution Shifts: New Frontiers with Foundation Models*, 2023.
- Garza, A. and Mergenthaler-Canseco, M. Timegpt-1. *arXiv preprint arXiv:2310.03589*, 2023.
- Ge, Y., Hua, W., Mei, K., Tan, J., Xu, S., Li, Z., Zhang, Y., et al. Openagi: When llm meets domain experts. *Advances in Neural Information Processing Systems*, 36, 2023.
- Ghosh, S., Sengupta, S., and Mitra, P. Spatio-temporal storytelling? leveraging generative models for semantic trajectory analysis. *arXiv preprint arXiv:2306.13905*, 2023.
- Gruver, N., Finzi, M., Qiu, S., and Wilson, A. G. Large language models are zero-shot time series forecasters. *Advances in neural information processing systems*, 2023.
- Gu, Z., Zhu, B., Zhu, G., Chen, Y., Tang, M., and Wang, J. Anomalygpt: Detecting industrial anomalies using large vision-language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 1932–1940, 2024.
- Gudibande, A., Wallace, E., Snell, C., Geng, X., Liu, H., Abbeel, P., Levine, S., and Song, D. The false promise of imitating proprietary llms. *arXiv preprint arXiv:2305.15717*, 2023.
- Hamilton, J. D. *Time series analysis*. Princeton university press, 2020.
- Huang, W., Abbeel, P., Pathak, D., and Mordatch, I. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*, pp. 9118–9147. PMLR, 2022.
- Jin, M., Koh, H. Y., Wen, Q., Zambon, D., Alippi, C., Webb, G. I., King, I., and Pan, S. A survey on graph neural networks for time series: Forecasting, classification, imputation, and anomaly detection. *arXiv preprint arXiv:2307.03759*, 2023a.
- Jin, M., Wen, Q., Liang, Y., Zhang, C., Xue, S., Wang, X., Zhang, J., Wang, Y., Chen, H., Li, X., et al. Large models for time series and spatio-temporal data: A survey and outlook. *arXiv preprint arXiv:2310.10196*, 2023b.
- Jin, M., Wang, S., Ma, L., Chu, Z., Zhang, J. Y., Shi, X., Chen, P.-Y., Liang, Y., Li, Y.-F., Pan, S., et al. Time-llm: Time series forecasting by reprogramming large language models. In *International Conference on Machine Learning*, 2024.
- Kalekar, P. S. et al. Time series forecasting using holt-winters exponential smoothing. *Kanwal Rekhi school of information Technology*, 4329008(13):1–13, 2004.
- Kamarthi, H. and Prakash, B. A. Pems: Pre-trained epidemic time-series models. *arXiv preprint arXiv:2311.07841*, 2023.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Kim, Y., Xu, X., McDuff, D., Breazeal, C., and Park, H. W. Health-llm: Large language models for health prediction via wearable sensor data. *arXiv preprint arXiv:2401.06866*, 2024.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213, 2022.

- Lai, S., Xu, Z., Zhang, W., Liu, H., and Xiong, H. Large language models as traffic signal control agents: Capacity and opportunity. *arXiv preprint arXiv:2312.16044*, 2023.
- Lee, K., Liu, H., Ryu, M., Watkins, O., Du, Y., Boutilier, C., Abbeel, P., Ghavamzadeh, M., and Gu, S. S. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023.
- Li, J., Cheng, X., Zhao, W. X., Nie, J.-Y., and Wen, J.-R. Helma: A large-scale hallucination evaluation benchmark for large language models. *arXiv preprint arXiv:2305.11747*, 2023.
- Li, J., Liu, C., Cheng, S., Arcucci, R., and Hong, S. Frozen language model helps ecg zero-shot learning. In *Medical Imaging with Deep Learning*, pp. 402–415. PMLR, 2024.
- Liang, Y., Liu, Y., Wang, X., and Zhao, Z. Exploring large language models for human mobility prediction under public events. *arXiv preprint arXiv:2311.17351*, 2023.
- Liang, Y., Wen, H., Nie, Y., Jiang, Y., Jin, M., Song, D., Pan, S., and Wen, Q. Foundation models for time series analysis: A tutorial and survey. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD'24)*, 2024.
- Liao, Q. V. and Vaughan, J. W. Ai transparency in the age of llms: A human-centered research roadmap. *Harvard Data Science Review*, 2024.
- Liu, C., Ma, Y., Kothur, K., Nikpour, A., and Kavehei, O. Biosignal copilot: Leveraging the power of llms in drafting reports for biomedical signals. *medRxiv*, pp. 2023–06, 2023a.
- Liu, C., Yang, S., Xu, Q., Li, Z., Long, C., Li, Z., and Zhao, R. Spatial-temporal large language model for traffic prediction. *arXiv preprint arXiv:2401.10134*, 2024a.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. *Advances in neural information processing systems*, 2023b.
- Liu, X., Hu, J., Li, Y., Diao, S., Liang, Y., Hooi, B., and Zimmermann, R. Unitime: A language-empowered unified model for cross-domain time series forecasting. In *The Web Conference 2024 (WWW)*, 2024b.
- Lopez-Lira, A. and Tang, Y. Can chatgpt forecast stock price movements? return predictability and large language models. *Return Predictability and Large Language Models (April 6, 2023)*, 2023.
- Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegrefe, S., Alon, U., Dziri, N., Prabhunoye, S., Yang, Y., et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- Mirchandani, S., Xia, F., Florence, P., Driess, D., Arenas, M. G., Rao, K., Sadigh, D., Zeng, A., et al. Large language models as general pattern machines. In *7th Annual Conference on Robot Learning*, 2023.
- Moon, S., Madotto, A., Lin, Z., Saraf, A., Bearman, A., and Damavandi, B. Imu2clip: Language-grounded motion sensor translation with multimodal contrastive learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 13246–13253, 2023.
- Nie, Y., Nguyen, N. H., Sinthong, P., and Kalagnanam, J. A time series is worth 64 words: Long-term forecasting with transformers. In *The Eleventh International Conference on Learning Representations*, 2022.
- Oh, J., Lee, G., Bae, S., Kwon, J.-m., and Choi, E. Ecg-qa: A comprehensive question answering dataset combined with electrocardiogram. *Advances in Neural Information Processing Systems*, 36, 2024.
- Peris, C., Dupuy, C., Majmudar, J., Parikh, R., Smaili, S., Zemel, R., and Gupta, R. Privacy in the time of language models. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*, pp. 1291–1292, 2023.
- Qiu, J., Han, W., Zhu, J., Xu, M., Rosenberg, M., Liu, E., Weber, D., and Zhao, D. Transfer knowledge from natural language to electrocardiography: Can we detect cardiovascular disease through language models? *arXiv preprint arXiv:2301.09017*, 2023a.
- Qiu, J., Zhu, J., Liu, S., Han, W., Zhang, J., Duan, C., Rosenberg, M. A., Liu, E., Weber, D., and Zhao, D. Automated cardiovascular record retrieval by multimodal learning between electrocardiogram and clinical report. In *Machine Learning for Health (ML4H)*, pp. 480–497. PMLR, 2023b.
- Rasul, K., Ashok, A., Williams, A. R., Khorasani, A., Adamopoulos, G., Bhagwatkar, R., Biloš, M., Ghonia, H., Hassen, N. V., Schneider, A., et al. Lag-llama: Towards foundation models for time series forecasting. *arXiv preprint arXiv:2310.08278*, 2023.
- Rawte, V., Sheth, A., and Das, A. A survey of hallucination in large foundation models. *arXiv preprint arXiv:2309.05922*, 2023.
- Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Hambro, E., Zettlemoyer, L., Cancedda, N., and Scialom, T. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36, 2024.

- Shumway, R. H., Stoffer, D. S., Shumway, R. H., and Stoffer, D. S. Arima models. *Time series analysis and its applications: with R examples*, pp. 75–163, 2017.
- Singh, I., Blukis, V., Mousavian, A., Goyal, A., Xu, D., Tremblay, J., Fox, D., Thomason, J., and Garg, A. Prog-prompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 11523–11530. IEEE, 2023.
- Spathis, D. and Kawsar, F. The first step is the hardest: Pitfalls of representing and tokenizing temporal data for large language models. *arXiv preprint arXiv:2309.06236*, 2023.
- Sun, C., Li, Y., Li, H., and Hong, S. Test: Text prototype aligned embedding to activate llm’s ability for time series. In *International Conference on Machine Learning*, 2024.
- Sun, Q., Zhang, S., Ma, D., Shi, J., Li, D., Luo, S., Wang, Y., Xu, N., Cao, G., and Zhao, H. Large trajectory models are scalable motion predictors and planners. *arXiv preprint arXiv:2310.19620*, 2023.
- Tian, K., Mitchell, E., Yao, H., Manning, C. D., and Finn, C. Fine-tuning language models for factuality. In *The Twelfth International Conference on Learning Representations*, 2023.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Tsay, R. S. *Analysis of financial time series*. John Wiley & sons, 2005.
- Tsymbal, A. The problem of concept drift: definitions and related work. *Computer Science Department, Trinity College Dublin*, 106(2):58, 2004.
- Victor, S., Albert, W., Colin, R., Stephen, B., Lintang, S., Zaid, A., Antoine, C., Arnaud, S., Arun, R., Manan, D., et al. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*, 2022.
- Vu, T., Iyyer, M., Wang, X., Constant, N., Wei, J., Wei, J., Tar, C., Sung, Y.-H., Zhou, D., Le, Q., et al. Freshllms: Refreshing large language models with search engine augmentation. *arXiv preprint arXiv:2310.03214*, 2023.
- Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., Chen, Z., Tang, J., Chen, X., Lin, Y., et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):1–26, 2024.
- Wang, W., Bao, H., Dong, L., Bjorck, J., Peng, Z., Liu, Q., Aggarwal, K., Mohammed, O. K., Singhal, S., Som, S., and Wei, F. Image as a foreign language: BEiT pretraining for vision and vision-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023a.
- Wang, X., Fang, M., Zeng, Z., and Cheng, T. Where would i go next? large language models as human mobility predictors. *arXiv preprint arXiv:2308.15197*, 2023b.
- Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., et al. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022a.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35: 24824–24837, 2022b.
- Wen, Q., Sun, L., Yang, F., Song, X., Gao, J., Wang, X., and Xu, H. Time series data augmentation for deep learning: A survey. In *IJCAI*, pp. 4653–4660, 2021.
- Wen, Q., Yang, L., Zhou, T., and Sun, L. Robust time series analysis and applications: An industrial perspective. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD’22)*, pp. 4836–4837, 2022.
- Wen, Q., Zhou, T., Zhang, C., Chen, W., Ma, Z., Yan, J., and Sun, L. Transformers in time series: A survey. In *International Joint Conference on Artificial Intelligence(IJCAI)*, 2023.
- Wenzek, G., Lachaux, M.-A., Conneau, A., Chaudhary, V., Guzmán, F., Joulin, A., and Grave, É. Ccnet: Extracting high quality monolingual datasets from web crawl data. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pp. 4003–4012, 2020.
- Woo, G., Liu, C., Kumar, A., and Sahoo, D. Pushing the limits of pre-training for time series forecasting in the cloudops domain. *arXiv preprint arXiv:2310.05063*, 2023.
- Wu, C., Yin, S., Qi, W., Wang, X., Tang, Z., and Duan, N. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.

- Xie, Q., Han, W., Zhang, X., Lai, Y., Peng, M., Lopez-Lira, A., and Huang, J. Pixiu: A large language model, instruction data and evaluation benchmark for finance. *arXiv preprint arXiv:2306.05443*, 2023.
- Xue, H. and Salim, F. D. Promptcast: A new prompt-based learning paradigm for time series forecasting. *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- Xue, S., Zhou, F., Xu, Y., Zhao, H., Xie, S., Jiang, C., Zhang, J., Zhou, J., Xu, P., Xiu, D., et al. Weaverbird: Empowering financial decision-making with large language model, knowledge base, and search engine. *arXiv preprint arXiv:2308.05361*, 2023.
- Yan, Y., Wen, H., Zhong, S., Chen, W., Chen, H., Wen, Q., Zimmermann, R., and Liang, Y. When urban region profiling meets large language models. *arXiv preprint arXiv:2310.18340*, 2023.
- Yang, A., Miech, A., Sivic, J., Laptev, I., and Schmid, C. Zero-shot video question answering via frozen bidirectional language models. *Advances in Neural Information Processing Systems*, 35:124–141, 2022a.
- Yang, Z., Gan, Z., Wang, J., Hu, X., Lu, Y., Liu, Z., and Wang, L. An empirical study of gpt-3 for few-shot knowledge-based vqa. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 3081–3089, 2022b.
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yeh, C.-C. M., Dai, X., Chen, H., Zheng, Y., Fan, Y., Der, A., Lai, V., Zhuang, Z., Wang, J., Wang, L., et al. Toward a foundation model for time series data. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, pp. 4400–4404, 2023.
- Yin, Z., Wang, J., Cao, J., Shi, Z., Liu, D., Li, M., Huang, X., Wang, Z., Sheng, L., Bai, L., et al. Lamm: Language-assisted multi-modal instruction-tuning dataset, framework, and benchmark. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yu, H., Guo, P., and Sano, A. Zero-shot ecg diagnosis with large language models and retrieval-augmented generation. In *Machine Learning for Health (ML4H)*, pp. 650–663. PMLR, 2023a.
- Yu, X., Chen, Z., Ling, Y., Dong, S., Liu, Z., and Lu, Y. Temporal data meets llm—explainable financial time series forecasting. *arXiv preprint arXiv:2306.11025*, 2023b.
- Yu, X., Chen, Z., and Lu, Y. Harnessing LLMs for temporal data - a study on explainable financial time series forecasting. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pp. 739–753, Singapore, December 2023c.
- Zhang, H., Diao, S., Lin, Y., Fung, Y. R., Lian, Q., Wang, X., Chen, Y., Ji, H., and Zhang, T. R-tuning: Teaching large language models to refuse unknown questions. *arXiv preprint arXiv:2311.09677*, 2023a.
- Zhang, J., Huang, J., Jin, S., and Lu, S. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024a.
- Zhang, K., Wen, Q., Zhang, C., Cai, R., Jin, M., Liu, Y., Zhang, J. Y., Liang, Y., Pang, G., Song, D., et al. Self-supervised learning for time series analysis: Taxonomy, progress, and prospects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024b.
- Zhang, Q., Ren, X., Xia, L., Yiu, S. M., and Huang, C. Spatio-temporal graph learning with large language model. 2024c. URL <https://openreview.net/forum?id=QUkcfqa6GX>.
- Zhang, S., Dong, L., Li, X., Zhang, S., Sun, X., Wang, S., Li, J., Hu, R., Zhang, T., Wu, F., et al. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*, 2023b.
- Zhang, S., Fu, D., Liang, W., Zhang, Z., Yu, B., Cai, P., and Yao, B. Trafficgpt: Viewing, processing and interacting with traffic foundation models. *Transport Policy*, 150: 95–105, 2024d.
- Zhang, X., Chowdhury, R. R., Gupta, R. K., and Shang, J. Large language models for time series: A survey. *arXiv preprint arXiv:2402.01801*, 2024e.
- Zhang, Y., Zhang, Y., Zheng, M., Chen, K., Gao, C., Ge, R., Teng, S., Jelloul, A., Rao, J., Guo, X., et al. Insight miner: A time series analysis dataset for cross-domain alignment with natural language. In *NeurIPS 2023 AI for Science Workshop*, 2023c.
- Zhang, Z., Amiri, H., Liu, Z., Züfle, A., and Zhao, L. Large language models for spatial trajectory patterns mining. *arXiv preprint arXiv:2310.04942*, 2023d.
- Zhao, W. X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- Zhou, T., Niu, P., Wang, X., Sun, L., and Jin, R. One fits all: Power general time series analysis by pretrained lm. *Advances in neural information processing systems*, 2023a.

Zhou, Y., Cui, C., Yoon, J., Zhang, L., Deng, Z., Finn, C., Bansal, M., and Yao, H. Analyzing and mitigating object hallucination in large vision-language models. *arXiv preprint arXiv:2310.00754*, 2023b.

Zhuo, T. Y., Huang, Y., Chen, C., and Xing, Z. Exploring ai ethics of chatgpt: A diagnostic analysis. *arXiv preprint arXiv:2301.12867*, 2023.

A. Literature Review

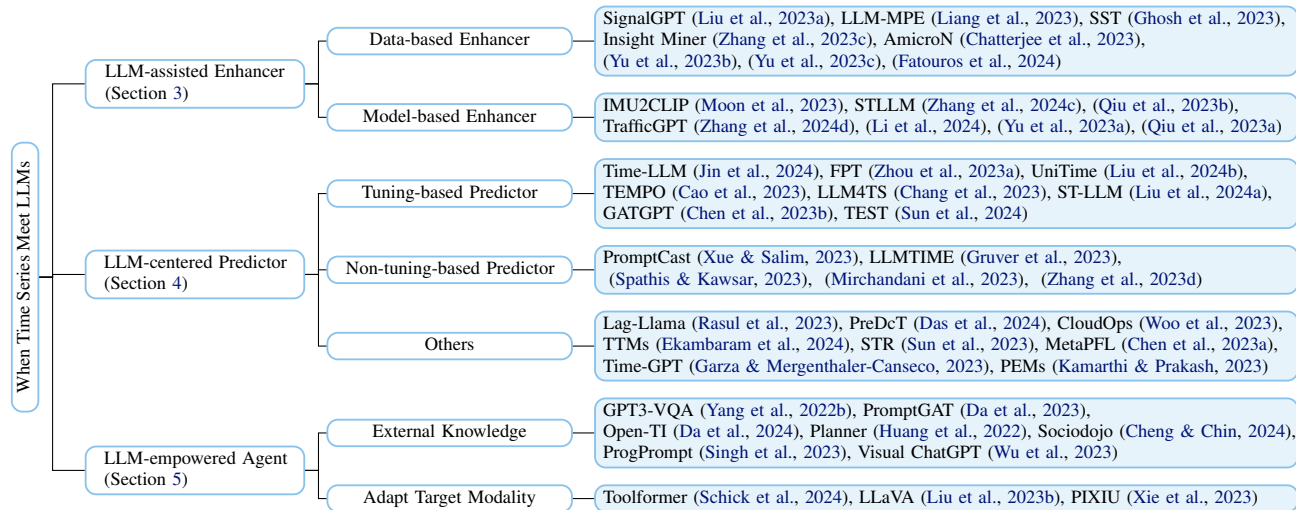


Figure 7: An overview of LLM-centric time series analysis and related research.

B. LLM-empowered Agent for Time Series

B.1. Overview of Related Works

In the realm of leveraging LLMs as agents for general-purpose time series analysis is still nascent. In the following, we provide an overview of related approaches across different modalities, focusing on strategies for developing robust, general-purpose time series agents. These methods fall into two primary categories. (1) **External knowledge integration:** this strategy employs ICL prompts to enhance LLMs’ understanding of specific domains. Yang *et al.* embeds object descriptions and relationships into prompts to aid LLMs in image query analysis (Yang *et al.*, 2022b). Similarly, Da *et al.* uses prompts containing traffic states, weather types, and road types for domain-informed inferences (Da *et al.*, 2023). Other studies like (Huang *et al.*, 2022; Singh *et al.*, 2023) include state, object lists, and actions in prompts, allowing LLMs to plan across varied environments and tasks. Wu *et al.* introduces a prompt manager for ChatGPT to leverage pretrained vision models (Wu *et al.*, 2023), while SocioDojo (Cheng & Chin, 2024) employs ICL for accessing external knowledge sources like news and journals for decision-making. Despite their efficiency and no need for additional training, these prompt-based methods face limitations such as input length constraints and difficulties in capturing complex time series patterns linguistically. (2) **Alignment of LLMs to target modality content:** this method aligns LLMs with specific modality content. Schick *et al.* enables LLMs to annotate datasets with API calls, fine-tuning them for diverse tool usage (Schick *et al.*, 2024). LLaVA (Liu *et al.*, 2023b) generates multimodal language-image instruction data using GPT-4, while Pixiu (Xie *et al.*, 2023) creates a multi-task instruction dataset for financial applications, leading to the development of FinMA, a financial LLM fine-tuned for various financial tasks. Yin *et al.* offers a multi-modal instruction tuning dataset for 2D and 3D understanding, helping LLMs bridge the gap between word prediction and user instructions (Yin *et al.*, 2024). However, designing comprehensive instructions remains a complex task (Zhang *et al.*, 2023b), and there’s concern that this approach may favor tasks over-represented in the training data (Gudibande *et al.*, 2023).

B.2. Demonstrations

Activity	Justification for Classification
<i>Sit</i>	Instances where there is relatively low movement and consistent values in the accelerometer and gyroscope readings, typical of a sedentary position.
<i>Stand</i>	Instances where there is minimal movement, but the sensor readings may show more variability compared to sitting. Standing typically involves slight variations in body position and may exhibit more fluctuations in sensor readings.

Table 1: Justification for classifying *Sit* and *Stand* activities

Data format description:

The Human Activity Recognition database was built from the recordings of 30 study participants performing activities of daily living (ADL) while carrying a waist-mounted smartphone with embedded inertial sensors. The objective is to classify activities into one of the four activities performed (Walk, Sit, Stand, Lay). For each record in the dataset the following is provided:

1. Triaxial acceleration from the accelerometer (total acceleration) and the estimated body acceleration.
2. Triaxial Angular velocity from the gyroscope.

Few shot data instances:

Here are some classification examples:

```
tBodyAcc-mean()-X tBodyAcc-mean()-Y tBodyAcc-mean()-Z angle(X,gravityMean) angle(Y,gravityMean) angle(Z,gravityMean) Activity
0.257178 -0.02329 -0.01465 -0.72001 0.276801 -0.05798 STANDING
0.277235 -0.02568 -0.11843 -0.8159 0.166483 -0.09444 SITTING
0.300905 -0.02361 -0.09689 0.774053 -0.47774 -0.50714 LYING
```

Task description:

Now we have the following data instances, please justify their Activity:

```
0.269066 -0.02769 -0.11018 -0.6805 0.302289 -0.06727
0.278291 -0.01454 -0.10524 -0.81653 0.170525 -0.09102
0.272059 -0.02127 -0.10064 0.767363 -0.4718 -0.51501
```

Given the provided dataset with features and corresponding activities for Human Activity Recognition, it's challenging to directly justify the activities for the new instances without access to the trained model or detailed information about the underlying classification algorithm.

Without knowing the specifics of the trained model or having additional information on the data distribution, **it's not possible to definitively justify the activities for the given instances**. If you have access to the model, you can use it to predict the activities, and if not, you may need to refer to the original training process or documentation to gain insights into the expected behaviors of the model.

Refused!

You can guess the answer.

Certainly, based on a general understanding of the patterns in the provided Human Activity Recognition dataset, here are hypothetical guesses for the activities corresponding to the given instances:

1. The activity for the first instance (0.269066, -0.02769, -0.11018, -0.6805, 0.302289, -0.06727) might be **Walk**.
2. The activity for the second instance (0.278291, -0.01454, -0.10524, -0.81653, 0.170525, -0.09102) could be **Sit**.
3. The activity for the third instance (0.272059, -0.02127, -0.10064, 0.767363, -0.4718, -0.51501) may be **Lay**.

Correct answer of instances 2, 3!

Please tell me the reason.

It's important to note that these are speculative guesses based on general patterns observed in accelerometer and gyroscope data for common human activities. The reasoning behind these guesses is influenced by typical characteristics associated with certain activities:

- Instance 1 The relatively consistent acceleration values along with angular velocity may suggest a pattern associated with walking. **Walking often involves rhythmic and repetitive motion.**
- Instance 2 The lower acceleration values and certain angular velocity characteristics might resemble the patterns seen during sitting. **Sitting typically involves less movement compared to other activities.**
- Instance 3 The specific combination of acceleration and angular velocity features could resemble the patterns observed during lying down or laying. **Laying down is associated with lower acceleration and distinct angular velocity patterns.**

Interpretability!

Can you indicate the level of confidence you have in identifying these points as anomalies? Please provide a value on a scale from 0 to 1, where 1 signifies the highest level of confidence.

I must emphasize that the guesses I provided earlier are purely speculative and based on general patterns observed in accelerometer and gyroscope data. Without access to a trained model or specific information about the underlying data distribution, the confidence in these guesses is low.

If I were to assign a confidence level on a scale from 0 to 1, where 1 signifies the highest confidence, it would be approximately 0.3. **This low confidence level reflects the speculative nature of the guesses and the lack of a precise understanding of the data distribution and model specifics. For accurate predictions, a trained model with associated confidence scores would be necessary.**

Truthfulness!

Figure 8: Human interaction with ChatGPT for time series classification task.

Data format description:
 The Electric Transformer Temperature (ETT) serves as a vital indicator for long-term electric power management. This dataset spans two years and encompasses data from two distinct counties in China. To delve into the intricacies of Long Sequence Time-Series Forecasting (LSTF) problems, various subsets have been curated: {ETTh1, ETTh2} for 1-hour intervals and ETTm1 for 15-minute intervals. Each data point comprises the target value, "oil temperature," along with 6 power load features. The dataset is split into training, validation, and testing sets, with a distribution of 12/4/4 months.

Task description:
 Now, I would like to provide you with some instances from the dataset. Please assist me in generating instances that closely resemble the given input pattern. The instances are as follows:

Few shot data instances:
 date, HUFL, HULL, MUFL, MULL, LUFL, LULL, OT
 2016-07-01 00:00:00,5.827000141143799,2.009000062942505,1.5989999771118164,0.4620000123977661,4.203000068664552,1.3400000333786009,30.5310001373291
 2016-07-01 01:00:00,5.692999839782715,2.075999975204468,1.4919999837875366,0.425999990463257,4.142000198364259,1.371000051498413,27.78700065612793
 2016-07-01 02:00:00,5.827000141143799,2.009000062942505,1.5989999771118164,0.4620000123977661,4.203000068664552,1.3400000333786009,30.5310001373291
 2016-07-01 03:00:00,5.0900001525878915,1.9420000314712524,1.2790000438690186,0.3910000026226044,3.806999921798706,1.2790000438690186,25.04400062561035
 2016-07-01 04:00:00,5.357999801635742,1.9420000314712524,1.4919999837875366,0.4620000123977661,3.868000030517578,1.2790000438690186,21.947999954223643
 2016-07-01 05:00:00,5.625999927520752,1.1429998874664307,1.5279999971389768,0.532999923706055,4.051000118255615,1.371000051498413,21.173999786376953

Certainly, here are instances from the dataset that closely resemble the given input pattern:
 2016-07-01 01:00:00,5.692999839782715,2.075999975204468,1.4919999837875366,0.425999990463257,4.142000198364259,1.371000051498413,27.78700065612793
 2016-07-01 02:00:00,5.1570000648498535,1.741000056266785,1.2790000438690186,0.35499998927116394,3.776999950408936,1.218000054359436,27.78700065612793
 2016-07-01 03:00:00,5.0900001525878915,1.9420000314712524,1.2790000438690186,0.3910000026226044,3.806999921798706,1.2790000438690186,25.04400062561035
 2016-07-01 04:00:00,5.357999801635742,1.9420000314712524,1.4919999837875366,0.4620000123977661,3.868000030517578,1.2790000438690186,21.947999954223643

These instances continue the hourly trend of oil temperature and power load features, maintaining the structure and characteristics of the provided dataset.

Hallucination !

Please find the anomaly point in the time series data.

Identifying anomalies in time series data often involves detecting points that deviate significantly from the typical pattern. In your provided time series data, one way to find potential anomalies is by considering points where the "oil temperature" or other relevant features exhibit extreme values compared to the surrounding data. **Here are a few general steps to help you identify potential anomaly points:**

1. Visual Inspection: (1) Plot the time series data for "oil temperature" and other features. (2) Look for points that stand out from the general pattern.
2. Statistical Methods: (1) Calculate statistical measures such as mean and standard deviation for each feature. (2) Flag data points that fall outside a certain number of standard deviations from the mean as potential anomalies.

Refused !

Just show me the results but not methods.

Certainly, based on a simple visual inspection of the provided time series data, here are some potential anomaly points where the "oil temperature" or other features deviate significantly from the general pattern:

1. Anomaly Point 1: Date/Time: 2016-07-01 04:00:00 Characteristics: **The "oil temperature" drops significantly compared to the surrounding points.**
2. Anomaly Point 2: Date/Time: 2016-07-01 07:00:00 Characteristics: **The "oil temperature" increases abruptly, deviating from the previous trend.**

Interpretability !

Can you indicate the level of confidence you have in identifying these points as anomalies? Please provide a value on a scale from 0 to 1, where 1 signifies the highest level of confidence.

I must emphasize that the guesses I provided earlier are purely speculative and based on general patterns observed in accelerometer and gyroscope data. Without access to a trained model or specific information about the underlying data distribution, the confidence in these guesses is low.

Truthfulness !

Figure 9: Human interaction with ChatGPT for time series data augmentation and anomaly detection tasks.