# Time-MQA: Time Series Multi-Task Question Answering with Context Enhancement

Anonymous ACL submission

#### Abstract

Time series data are foundational in finance, healthcare, and energy domains. However, most existing methods and datasets remain focused on a narrow spectrum of tasks, such as forecasting or anomaly detection. To bridge this gap, we introduce Time Series Multi-Task Question Answering (Time-MQA), a unified framework that enables natural language queries across multiple time series tasks - numerical analytical tasks and open-ended ques-011 tion answering with reasoning. Central to Time-MQA is the TSQA dataset, a large-scale dataset containing ~200k question-answer pairs de-014 rived from diverse time series spanning environment, traffic, etc. This comprehensive resource covers various time series lengths and promotes robust model development. We fur-017 ther demonstrate how continually pre-training large language models (Mistral 7B, Llama-3 8B, and Qwen-2.5 7B) on the TSQA dataset enhanced time series reasoning capabilities, moving beyond mere numeric tasks and enabling more advanced and intuitive interactions with temporal data. The complete TSQA dataset, models, executable codes, user study questionnaires for evaluation, and results have all been open-sourced here<sup>1</sup>.

#### 1 Introduction

Time series analysis has long been fundamental to real-world applications in finance, healthcare, energy, and other domains (Nie et al., 2024; Xu et al., 2023; Yang et al., 2021b; Chen et al., 2024b). Before the emergence of large language models (LLMs), most research in this area focused on numerical analytical tasks such as forecasting and anomaly detection, with methods typically constrained to a single objective. Recently, the community has begun to explore ways to enhance these conventional time series approaches by incorporating LLMs (Fons et al., 2024; Zhang et al., 2024b); Healthcare Domain - Time Series Array (150 time points)



Figure 1: Example of Time-MQA with context enhancement. The LLM answers the user's questions based on the provided context and the input time series. It can conduct reasoning and handle multiple time series tasks.

however, while some efforts do leverage contextual information, they typically focus on a single task, such as forecasting, leaving a gap in broader, multitask reasoning and inference capabilities (Merrill et al., 2024; Ansari et al., 2024; Frisoni et al., 2024; Jin et al., 2024a). To bridge this gap, we propose a unified *Time Series Multi-task Question Answering (Time-MQA)* framework that integrates diverse tasks with natural language queries.

*Time-MQA* is crucial for advanced reasoning and inference, as it enables models to interpret temporal data through natural language queries and uncover deeper insights beyond mere classi-

<sup>&</sup>lt;sup>1</sup>https://anonymous.4open.science/r/Time-MQA-6085/

cal numeric tasks (Chow et al., 2024; Xu et al., 2024). However, there is a notable lack of relevant datasets - specifically, paired language with time series — which severely limits the development of models capable of dialogue and reasoning in this domain (Jin et al., 2024b; Kong et al., 2025). This significant gap underscores the urgent need for comprehensive resources that integrate multiple time series tasks under a question-answering framework, allowing users to query these tasks in natural language.

055

056

067

076

090

096

098

100

101

102

104

In this paper, we introduce *Time-MQA*, a new framework for multi-task time series question answering with context enhancement. As illustrated in Figure 1, Time-MQA expands beyond numerical analytical tasks by consolidating open-ended questions and classical time series tasks, such as forecasting, into a single framework. Users can pose queries in natural language, enabling a more intuitive and flexible interface for a wide range of time series analyses.

To support Time-MOA, we proposed TSOA, a large-scale dataset with approximately 200k question-answer pairs spanning multiple domains, including healthcare, environment, AIOps, machine sensors, finance, energy, traffic, IoT, nature, transport, human activities, and the web. This dataset covers various time series lengths and tasks, ensuring broad coverage and robustness. Notably, TSQA features open-ended reasoning questions with more elaborate text-based explanations. To the best of our knowledge, this is the first large-scale QA dataset in the time series domain that spans multiple domains and tasks, effectively bridging the gap between classical time series analysis and modern LLM-driven approaches.

We further demonstrate the utility of the TSQA dataset by employing continual pre-training techniques on representative LLMs such as Mistral 7B (Jiang et al., 2023), Llama-3 8B (Dubey et al., 2024), and Qwen-2.5 7B (Yang et al., 2024). Our experiments reveal that these models trained on the TSQA dataset can effectively acquire time series knowledge and reasoning abilities, enabling more advanced capabilities beyond basic numeric handling of time series data.

In summary, our contributions include:

• We propose *Time-MQA*, a multi-task time series question answering framework that leverages contextual enhancement to extend beyond traditional numerical analysis tasks.

• We construct TSQA, a dataset comprising  $\sim 200$ k question-answer pairs across over twelve domains (e.g., healthcare, finance, and energy) and five tasks (e.g., forecasting, anomaly detection, 108 and open-ended reasoning question answering). 109

105

106

107

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

We demonstrate that fine-tuning LLMs on the TSQA dataset equips them with time series-specific knowledge, enabling natural language queries for comprehensive time series analysis.

#### **Related Works** 2

#### 2.1 **Classical Time Series Tasks**

Time series analysis has been extensively studied in various real-world applications, such as finance, healthcare, climate, electricity, AIOps, and industrial system maintenance (Nie et al., 2024; Guo et al., 2024; Ma et al., 2024; Yang et al., 2023a, 2021a). Classical time series tasks focus on extracting insights from these time series and addressing challenges associated with temporal patterns (Fuller, 2009; Hamilton, 2020).

Forecasting is one the most fundamental tasks. It predicts future time points based on historical values and features (Lim and Zohren, 2021). Depending on the temporal horizon, forecasting can be categorized into short-term forecasting, which captures immediate fluctuations, and long-term forecasting, which models more complex temporal dependencies (Wang et al., 2024b). Common methods include statistical methods (e.g., ARIMA and exponential smoothing) and deep learning-based methods (e.g., RNN, LSTM, Transformer-based architectures) (Wen et al., 2022; Miller et al., 2024).

Anomaly detection seeks to identify abnormal patterns or deviations from expected behavior in time-series data (Zamanzadeh Darban et al., 2024). Classical approaches rely on statistical models like z-score analysis and dynamic thresholding, while contemporary methods incorporate deep learning-based frameworks, such as autoencoders, transformers, and graph neural networks, to capture complex dependencies and temporal correlations (Han et al., 2022; Yang et al., 2023b).

Imputation addresses the issue of missing or corrupted data in sequence. It is essential for ensuring data integrity in downstream analysis (Du et al., 2024). Traditional imputation techniques, such as interpolation, have been widely employed, whereas deep learning-based methods, such as variational autoencoders (VAE), generative adversarial



Figure 2: The overview of the proposed Time-MQA framework.

networks (GANs), and diffusion models, have recently demonstrated promising results in learning complex missing patterns and improving imputation performance (Wang et al., 2025).

155

156

157

159

160

161

162

163

165

166

167

168

169

170

171

173

174

175

176

177

178

181

182

185

186

187

189

190

191

193

Beyond these tasks, time series classification/regression, generation, augmentation, and decomposition are also frequently used in real-world scenarios (Mohammadi Foumani et al., 2024; Wen et al., 2020; Zhang et al., 2024a).

### 2.2 Text-Enhanced Time Series Tasks

Recent advancements in time series analysis have demonstrated the potential of incorporating textual information based on LLMs to enhance time series tasks (Jin et al., 2024b; Wang et al., 2024a). Unlike classical approaches that rely solely on numerical data, text-enhanced time series analysis leverages domain-specific textual descriptions, contextual metadata, or associated reports to improve the cognitive understanding and modeling of timedependent patterns (Liu et al., 2024c). This hybrid approach mitigates the limitations of unimodal time series models by integrating additional semantic and contextual cues that are useful for decisionmaking (Kong et al., 2025; Singh et al., 2024).

In detail, text-enhanced time series forecasting and anomaly detection tasks benefit from textual information by incorporating expert reports to refine predictions and provide anomaly causal explanations (Hollmann et al., 2025; Chen et al., 2023). Similarly, classification, imputation, and generation tasks can be enhanced by leveraging textual descriptions as auxiliary supervision or describing missing values (Bernardini et al., 2023; Moor et al., 2023). It will help models distinguish subtle variations across different categories and generate more informed reconstructions.

Recent research has explored various techniques for integrating textual and time-series data, including LLM-based alignment, cross-modal attention mechanisms, and contrastive learning strategies that jointly encode text and time-series representations (Jin et al., 2023b; Zhang et al., 2023; Liu et al., 2024a). Some approaches, such as Time-LLM (Jin et al., 2023a), directly adapt LLMs to process text and time-series data, whereas others, like Time-MMD (Liu et al., 2024b), employ weighted fusion methods to combine textual embeddings with deep time-series backbones. By enhancing classical time series tasks with textual information, text-enhanced time series models offer greater robustness and richer interpretability across diverse applications (Jin et al., 2024b; Kong et al., 2025). 194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

#### 2.3 Language Question Answering

Question Answering (QA) in Natural Language Processing (NLP) involves systems that interpret human language queries to retrieve or generate accurate answers (Biancofiore et al., 2024; Chen et al., 2024c). It evolves from rule-based systems to neural architectures driven by LLMs like GPT-4 (Achiam et al., 2023) and Llama (Touvron et al., 2023). These models leverage massive text corpora and large-scale datasets for end-toend pre-training, fine-tuned via supervised learning or reinforcement learning with human feedback (RLHF) to align responses with factual and contextual relevance (Liu et al., 2023). Innovations such as retrieval-augmented generation (RAG) combine parametric knowledge with external data sources, while benchmarks like SQuAD, HotpotQA, MuSiQue, FinTextQA, SyllabusQA, and ToolQA drive progress (Trivedi et al., 2022; Ho et al., 2020; Yang et al., 2018; Zhuang et al., 2024; Chen et al., 2024a; Fernandez et al., 2024). However, challenges remain, such as handling ambiguous queries, ensuring the accuracy of generated answers, and maintaining efficiency in processing large volumes of data. Ongoing research focuses on enhancing the reasoning capabilities of QA sys-

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

329

330

284

234

236

240

241

242

243

244

245

246

247

248

251

255

257

258

262

267

270

271

272

273

275

277

278

279

281

tems, improving their ability to handle complex and nuanced questions, and expanding their applicability across diverse domains (Singh et al., 2025).

## 3 Methodology

### 3.1 The Time-MQA Framework

The Time-MOA framework broadens traditional analysis by integrating diverse objectives - such as forecasting, imputation, classification, anomaly detection, and notably, open-ended queries - into a unified question-answering paradigm (as shown in Figure 2. Formally, let  $\mathbf{X} = \{x_1, x_2, ..., x_T\}$ denote a time series input, where  $x_t \in \mathbb{R}^d$  represents a *d*-dimensional observation at timestep *t*. Let C represent additional contextual information (e.g., textual metadata, domain-specific knowledge, or other modalities). For a question Q expressed in natural language, the goal is to generate an answer A conditioned on both X and contextual information C. Time-MQA aims to learn a function  $f: (\mathbf{X}, \mathbf{C}, Q) \to A$ , where A is the correct answer to the query Q. Depending on the nature of Q, A can take diverse forms, such as predicted values, classification labels, a set of anomalous timestamps, and textual explanations. Specifically, the model component of Time-MQA is based on continued pre-trained LLMs (i.e., Mistral 7B (Jiang et al., 2023), Llama-3 8B (Dubey et al., 2024), and Qwen-2.5 7B (Yang et al., 2024)), using the prepared TSQA dataset. To optimize parameter usage, Time-MQA employs Parameter-Efficient Fine-Tuning (PEFT) with a LoRA adapter. The supported various tasks of the Time-MQA with context enhancement are demonstrated with details in Figure 3.

**Key Distinctions.** There are three main differences between Time-MQA and traditional time series analysis tasks:

- (1) Task Scope: Traditional tasks focus on singular objectives (e.g., forecasting future values or classifying the time series). In contrast, Time-MQA unifies these under a question-driven paradigm, enabling both conventional tasks (e.g., "Forecast the next 5 values") and complex queries (e.g., "Why did the temperature drop abruptly at hour 10?") that require joint reasoning across detection, explanation, and contextual knowledge.
- (2) Context Enhancement: Traditional methods rely solely on the time series X. Time-

MQA integrates auxiliary context C to resolve ambiguities and improve robustness. For instance, even when analyzing the same time series, differences in a dataset's background information can yield contrasting predictions (Williams et al., 2024).

(3) Multi-Task Generalization: Unlike singlepurpose models, Time-MQA dynamically adapts to diverse question types through a unified architecture, eliminating the need for taskspecific pipelines. This flexibility allows it to handle various time series tasks within a single framework, promoting knowledge sharing via shared representations that enhance both performance and interpretability.

### 3.2 The TSQA Dataset

This section will introduce our self-constructed 200k-level text-enhanced time series dataset, i.e., TSQA. From perspectives such as the data collection process and types across multiple domains and tasks. We compare it with other datasets and demonstrate the advantages of TSQA.

### 3.2.1 Dataset Composition and Categorization

The raw data we use comes from various classic publicly available datasets in the time series domain, covering multiple tasks and application areas. Innovatively, we have incorporated textual descriptions, including background information, feature descriptions, etc. From a task perspective, our dataset can be categorized into the following types. Some examples of the proposed TSQA dataset can be found in the Appendix.

**Forecasting.** In the forecasting task, we utilize UTSD datasets (Liu et al.), publicly available time series forecasting datasets, such as ETTh1, ETTh2, ETTm1, ECL, Weather, etc. (Zhou et al., 2021), and the latest text-enhanced time series datasets as the raw time-series data. Additionally, for the first time, we incorporate financial datasets with earnings call transcripts (The Motley Fool, 2024). To ensure the generalization ability of algorithms and models, we set the input length and prediction length to random values between 64–256 and 8–32, respectively. Furthermore, we enhance all data with background information, feature descriptions, and task descriptions as textual information based on the original data sources. In summary, the forecasting dataset includes 42,557 data instances spanning application domains such as energy, environment,



Figure 3: The demonstration of the Time-MQA with context enhancement.



Figure 4: Distribution of data statistics in TSQA dataset.

health, IoT, nature, transport, Web, AIOps, etc.

331

334

**Imputation.** The imputation task shares the same 332 original time series data sources as the forecasting task. Additionally, we randomly set the input length between 64 and 256, randomly removed 4 to 335 12 values, and replaced them with "X". The imputation dataset consists of a total of 38,657 instances from multiple application domains. 338

Anomaly Detection. For the anomaly detection task, we utilize commonly used public datasets in the field, such as the NASA Series dataset, UCR, ECG, KPI, MGAB, NAB, SensorScope, and Yahoo, among others (Zamanzadeh Darban et al., 2024). To enhance the generalization capability of algorithms and models, we randomly set the input length between 8 and 256 while ensuring a balanced distribution between anomalous and normal data. The anomaly detection dataset comprises 37,000 instances from various application domains, including AIOps, Health, Finance, Machinery, Industrial Sensors, Environment, and Traffic.

339

340

341

342

345

350

351

353

355

356

357

359

360

361

362

363

364

366

**Classification.** For the classification task, the data sources are relatively straightforward, primarily derived from the human activity recognition application domain (Kwapisz et al., 2011; Bachlin et al., 2009). The dataset includes both binary and multi-class classification tasks. To ensure class balance, we collected a total of 37,000 data instances. Additionally, the input sequence length is randomly set between 8 and 32 to maintain diversity in data representation.

Open-Ended Reasoning QA. We used different parts of the UTSD dataset (Liu et al.) - distinct from those employed in our forecasting tasks - to avoid data leakage. To generate open-ended reasoning QA, we utilized GPT-40, instructing it

Table 1: Comparison of text-enhanced time series datasets.  $\bigcirc$  and  $\checkmark$  indicate having only real data and having both synthetic and real data.

| Dataset     | Numerical | Cognitive | Source | Domain | Size  |
|-------------|-----------|-----------|--------|--------|-------|
| TS-Insights | ~         | ×         | 0      | 4      | ~10k  |
| ChatTS      | ×         | ~         | ~      | 4      | ~2.2k |
| CiK         | ~         | ×         | ⊘      | 7      | ~2.9k |
| TimeMMD     | ~         | ×         | ⊘      | 9      | ~16k  |
| TSQA        | ~         | ~         | ~      | 12     | ~200k |

to create questions covering various topics such as trends, seasonality, cyclical patterns, summarization, volatility, anomalies, structural breaks, and other statistical properties (without limiting it strictly to these areas). We also incorporated multiple question types, including multiple-choice, true/false, and open-ended formats. The prompts used to generate these QA pairs are provided in the Appendix. In total, we generated 37,629 data instances, from which we manually reviewed and selected 1,400 QA pairs for use in the continual pretraining step.

#### **3.2.2 Data Statistics**

367

371 372

374

382

386

398

400

401

402

403

404

405

406

The TSQA dataset comprises 192,843 (~200k) question-answer pairs spanning twelve domains - healthcare, finance, energy, traffic, environment, IoT, nature, transport, human activities, machine sensors, AIOps, and the web - and five task types: forecasting, imputation, anomaly detection, classification, and open-ended reasoning (see Figure 4 for distributions). Within the open-ended reasoning QA, the dataset includes 6,919 true/false questions, 11,281 multiple-choice questions, and 12,510 openended questions, offering a broad and diverse range of question formats.

### 3.2.3 Comparison with Existing Datasets

We summarize several existing datasets compared to our proposed TSQA dataset in Table 1. Unlike prior datasets (TS-Insights (Zhang et al., 2023), ChatTS (Xie et al., 2024), CiK (Williams et al., 2024), and TimeMMD (Liu et al., 2024b)), which focus on either classical numerical analytical time series tasks (e.g., forecasting, anomaly detection, imputation, classification) or text-enhanced cognitive tasks (e.g., reasoning, QA), our TSQA dataset supports both, making it the most comprehensive benchmark. It also covers the widest range of application domains (12 vs. a maximum of 9) and is significantly larger (~200k instances vs. ~10k–16k). Furthermore, TSQA includes both real and syn-

| Hyperparameter              | Assignment                      |
|-----------------------------|---------------------------------|
| Base model                  | Mistral 7B                      |
| Computing infrastructure    | 1*A100-80GB GPU                 |
| Max steps                   | 4000                            |
| Warm-up steps               | 1000                            |
| Batch size per device       | 4                               |
| Gradient accumulation steps | 8                               |
| Learning rate               | 5e-5                            |
| Embedding learning rate     | 1e-5                            |
| Optimizer                   | AdamW (8-bit)                   |
| Learning rate scheduler     | Cosine                          |
| Weight decay                | 0.1                             |
| LoRA rank (r)               | 16                              |
| LoRA alpha                  | 16                              |
| LoRA dropout                | 0.0                             |
| LoRA target modules         | q_proj, k_proj, v_proj, o_proj, |
|                             | gate_proj, up_proj, down_proj   |
| Training time               | ~1 Day                          |

Table 2: Hyper-parameters and training time of finetuning the Mistral 7B based on the TSQA dataset.

thetic data, ensuring greater diversity and robustness for text-enhanced time series analysis. 407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

#### 4 Experiment and Result

#### 4.1 Experimental Settings

In the experiments, we considered the ratio of application domains and data sources by randomly selecting 1,400 QA pairs for each task type - forecasting, imputation, anomaly detection, and classification - and by manually reviewing and selecting 1,400 QA pairs for open-ended reasoning, leading to a total of 7,000 QA pairs. We then followed the setting from (Cheng et al., 2024). Specifically, to ensure the model adequately learns the time-series domain, we then combined our dataset with a general QA corpus sourced from OpenOrca (Lian et al., 2023) at a 70% to 30% ratio, resulting in 10k QA pairs overall. Finally, we formatted all QA pairs so that the question and answer were clearly labeled, and then we tokenized the text. An example of a formatted QA pair from the pre-tokenized text is provided in the Appendix. All training runs were conducted on a single A100 80GB GPU. For a detailed overview of the hyperparameters used and training time, please refer to Table 2, using Mistral 7B for example.

### 4.2 Results

**Main Results.** Table 3 presents the performance of three fine-tuned models, alongside the GPT-4o (Achiam et al., 2023) model and Doubao (Volcengine, 2023) results, across multiple tasks. For Table 3: Comparison of our three fine-tuned models, GPT-40, and Doubao across diverse tasks. Forecasting and imputation tasks were evaluated using average MSE, while anomaly detection, classification, and open-ended reasoning tasks (including multiple-choice questions (MCQs) and true-false judgment) were measured by accuracy. \* Doubao uses simple mean forecasting, which outputs the same value for all forecasting.

| Daalzhana   | Classical Numerical Task |                         |                              |                           | <b>Open-Ended Reasoning QA</b> |       |
|-------------|--------------------------|-------------------------|------------------------------|---------------------------|--------------------------------|-------|
| Dackbolle   | Forecasting $\downarrow$ | Imputation $\downarrow$ | Anomaly Detection $\uparrow$ | Classification $\uparrow$ | Judgment ↑                     | MCQ ↑ |
| Doubao      | *                        | 0.018                   | 0.52                         | 0.44                      | 0.78                           | 0.56  |
| GPT-40      | 1.79                     | 0.018                   | 0.64                         | 0.32                      | 0.72                           | 0.58  |
| Llama-3 8B  | 2.01                     | 0.020                   | 0.54                         | 0.24                      | 0.74                           | 0.48  |
| Qwen-2.5 7B | 1.82                     | 0.016                   | 0.68                         | 0.52                      | 0.82                           | 0.54  |
| Mistral 7B  | 1.35                     | 0.014                   | 0.58                         | 0.44                      | 0.80                           | 0.64  |

evaluation, we randomly selected 50 QA pairs for each task type (or question format), taking into account their respective application domains.

Overall, our fine-tuned models demonstrated improved performance across multiple tasks. For open-ended reasoning QA, the fine-tuned Qwen model achieved an accuracy of 82% on judgment questions, while the Mistral model reached 64% on multiple-choice questions. In forecasting tasks, the MSE values were relatively high, likely due to the long time series in our dataset, which can be challenging for models to handle. Nevertheless, our fine-tuned Mistral model outperformed GPT-40.

For imputation tasks, the MSE remained low, possibly because having more extensive time series data provides additional information that supports better imputation. Our fine-tuned Mistral model generally showed stronger performance. These results suggest that by continuously pre-training on QA time series pairs, models can learn and leverage time series patterns more effectively.

**Forecasting with Long Time Series.** Figure 5 illustrates an example of forecasting results produced by our fine-tuned Mistral 7B model and GPT-40 on a long time-series input. Notably, the Mistral 7B model's forecasts outperform those generated by GPT-40. Moreover, the fine-tuned model can provide the rationale behind its predictions. When we posed follow-up questions in natural language, the model offered explanations of its reasoning. This illustrates that the model trained using the TSQA dataset can effectively learn complex time-series patterns and seasonal trends.

TSQA-Tuned vs. Zero-Shot Model. To further
demonstrate the effectiveness of our fine-tuned
model on the proposed TSQA dataset, we compared the performance of the tuned Mistral 7B
with that of the original pretrained model on open-



Figure 5: Example of long time series forecasting.

Table 4: Performance of TSQA-tuned and Zero-shotMistral 7B on two tasks (evaluated at accuracy).

| Model      | Judgment $\uparrow$ | $MCQ\uparrow$ |
|------------|---------------------|---------------|
| Zero-shot  | 0.78                | 0.60          |
| TSQA-tuned | 0.80                | 0.64          |

ended reasoning QA tasks (including judgment and multiple-choice questions). The results, shown in Table 4, indicate that the TSQA-tuned model performs slightly better than the zero-shot model, suggesting the value of our TSQA dataset for improving question-answering performance.

475

476

477

478

479

480

481

482

483

484

**Open-Ended Reasoning QA.** Figure 6 presents an example of an open-ended reasoning question answering and compares the responses generated by our fine-tuned models with those from GPT-

467

469

| Q: Does the sequence [0.28, 0.3, 0.34, 0.39, 0.44,<br>0.49, 0.54, 0.58, 0.62, 0.65, 0.68, 0.7, 0.71, 0.72,<br>0.73, 0.73, 0.73, 0.73, 0.73, 0.72, 0.71, 0.7, 0.69, 0.68]<br>exhibits an increasing pattern?   |
|---|
| A: The sequence consistently increases until it peaks at<br>0.73, after which it slightly decreases, but remains<br>relatively high.  |
| A: The given sequence does not display a clear increasing pattern throughout. While there is an initial increase up to 0.73, followed by a period of stability, the later points decrease consistently from 0.73 to 0.68, breaking the continuity of an upward pattern.   |
| A: The sequence does show an initial increase from<br>0.28 to 0.73, but it then stabilizes at 0.73 for several<br>data points before decreasing slightly towards the end<br>(down to 0.68). Therefore, while there is a phase of<br>increase, the overall description suggests both rise and<br>stability/decline, which contrasts with an exclusively<br>increasing pattern. |
| A: The sequence initially increases from 0.28 to 0.73,<br>but after reaching 0.73, it remains constant for a few<br>terms and then starts decreasing from 0.73 to 0.68.<br>Therefore, the sequence does not exhibit a<br>continuously increasing pattern throughout.  |

Figure 6: Example of open-ended reasoning QA.

40. As shown, our fine-tuned models effectively justify the reasons behind their answers. Notably, the Qwen-2.5 7B model provides a more thorough and comprehensive reasoning process compared to the other models.

#### 4.3 User Study on Time Series QA

Given the subjectivity of time series question answering and reasoning tasks, we conducted a user study. We recruited 78 participants from diverse domains, such as data science researchers, financial practitioners, and students of liberal arts, who participated in this study anonymously. The four optional models are A) Llama-3 8B, B) Mistral 7B, C) Qwen-2.5 7B, and D) GPT-40 sequentially. Participants assessed responses based on accuracy, completeness, clarity, and user preference. The complete result and questionnaire of the user study can be found in the Appendix, and structured summaries of the findings are as follows:

504Accuracy and Completeness. Mistral achieved505the highest accuracy in tasks requiring numeri-506cal precision. For cyclical pattern identification507(Q4), it received 80.8% user preference, and 69.2%508selected its response as most comprehensive for509volatility analysis (Q6). However, in trend analysis510(Q12), 51.3% preferred Qwen over Mistral (37.2%).511In summary, Mistral's and Qwen's answers to the512numerical questions are generally considered to be

more accurate as well as more comprehensive.

**Explanation Clarity.** Explanation clarity was measured by asking participants to rate the understandability of the models' reasoning. In Question 7 ("Considering the data points [60.0, 28.0, ...], do you see any seasonal patterns?"), more participants found Mistral's explanation clear, praising its logical structure and ease of understanding. Although Qwen's responses were also well received, they trailed slightly behind in these clarity ratings.

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

**User Preference.** In terms of overall user preference, when asked which answer they preferred (as in Question 11), mistral was selected by 70.5% of participants, while Qwen received 32.1% of the votes. This reflects a general tendency among users to favor responses that combine detailed data analysis with clear, logical explanations. In contrast to the above two models, Llama and GPT-40 are slightly less well accepted.

Overall, from the user study, Qwen and Mistral emerged as the top models for accuracy-driven tasks, while Mistral excelled in generating thorough explanations and more preferences. This highlights a trade-off between precision and interpretability in open-ended time series reasoning.

### 5 Conclusion

In this paper, we introduce *Time-MQA*, a multi-task framework that unifies time series analysis through natural language question answering, overcoming the limitations of single-task approaches. By integrating forecasting, imputation, anomaly detection, classification, and open-ended reasoning, Time-MOA enables flexible interactions with time-series data. The introduced TSQA dataset comprised of  $\sim$ 200k question-answer pairs across twelve diverse domains and varied time series lengths. It is a vital resource for advancing time series question answering and reasoning in LLMs. Experimental results show that continually pre-training models such as Mistral 7B, Llama 8B, and Qwen-2.5 7B on TSQA enhance their ability to interpret temporal patterns-going beyond numeric tasks to generate contextually rich insights. This work bridges the gap between traditional time series analysis and modern LLM capabilities, making temporal information more accessible and context-aware. Future directions include expanding the TSQA dataset to dynamic real-world data streams and exploring strategies to further refine time series reasoning.

499

501

503

485

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

612

613

### Limitations

562

While Time-MQA advances multi-task time series question answering, it has limitations. Our 564 proposed TSQA dataset, though diverse, may not 565 cover all real-world scenarios, particularly highly 566 irregular or domain-specific time series. But data with those properties may be common in real sce-568 narios (Mulayim et al., 2024). Our approach relies on pre-training LLM architectures, which may 570 not be optimized for long numerical sequences, which may cause hallucination problems (Yao et al., 2023). Evaluating open-ended reasoning remains 573 challenging, as current automated metrics may not 574 fully capture response quality (Kong et al., 2025). Additionally, computational costs for fine-tuning large models are high, which may limit some real-577 time applications (Jin et al., 2024b). 578

### **Ethics Statement**

Except for the pre-training corpora in the finance domain, all datasets and language models used in 581 this work are publicly available. The TSQA dataset 582 was constructed from well-established, publicly ac-583 cessible time series benchmarks and synthetic data 584 generation followed ethical guidelines to minimize 585 biases and ensure data quality. Our approach does not involve any personally identifiable information 587 (PII) or sensitive data, and we adhere to responsible AI practices by following the ACL Code of Ethics.

### References

590

595

596

597

598

599

604

610

611

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- 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. 2024. Chronos: Learning the language of time series. arXiv preprint arXiv:2403.07815.
- Marc Bachlin, Daniel Roggen, Gerhard Troster, Meir Plotnik, Noit Inbar, Inbal Meidan, Talia Herman, Marina Brozgol, Eliya Shaviv, Nir Giladi, et al. 2009. Potentials of enhanced context awareness in wearable assistants for parkinson's disease patients with the freezing of gait syndrome. In 2009 International Symposium on Wearable Computers, pages 123–130. IEEE.
- Michele Bernardini, Anastasiia Doinychko, Luca Romeo, Emanuele Frontoni, and Massih-Reza Amini.

2023. A novel missing data imputation approach based on clinical conditional generative adversarial networks applied to ehr datasets. *Computers in Biology and Medicine*, 163:107188.

- Giovanni Maria Biancofiore, Yashar Deldjoo, Tommaso Di Noia, Eugenio Di Sciascio, and Fedelucio Narducci. 2024. Interactive question answering systems: Literature review. *ACM Computing Surveys*, 56(9):1–38.
- Jian Chen, Peilin Zhou, Yining Hua, Yingxin Loh, Kehui Chen, Ziyuan Li, Bing Zhu, and Junwei Liang. 2024a. Fintextqa: A dataset for long-form financial question answering. *Association for Computational Linguistics(ACL)*.
- Wei Chen, Yuxuan Liang, Yuanshao Zhu, Yanchuan Chang, Kang Luo, Haomin Wen, Lei Li, Yanwei Yu, Qingsong Wen, Chao Chen, et al. 2024b. Deep learning for trajectory data management and mining: A survey and beyond. *arXiv preprint arXiv:2403.14151.*
- Xiaoyang Chen, Ben He, Hongyu Lin, Xianpei Han, Tianshu Wang, Boxi Cao, Le Sun, and Yingfei Sun. 2024c. Spiral of silences: How is large language model killing information retrieval?–a case study on open domain question answering. *Association for Computational Linguistics(ACL)*.
- Zhuo Chen, Wen Zhang, Yufeng Huang, Mingyang Chen, Yuxia Geng, Hongtao Yu, Zhen Bi, Yichi Zhang, Zhen Yao, Wenting Song, et al. 2023. Teleknowledge pre-training for fault analysis. In 2023 IEEE 39th International Conference on Data Engineering (ICDE), pages 3453–3466. IEEE.
- Daixuan Cheng, Yuxian Gu, Shaohan Huang, Junyu Bi, Minlie Huang, and Furu Wei. 2024. Instruction pretraining: Language models are supervised multitask learners. *arXiv preprint arXiv:2406.14491*.
- Winnie Chow, Lauren Gardiner, Haraldur T Hallgrímsson, Maxwell A Xu, and Shirley You Ren. 2024. Towards time series reasoning with llms. *arXiv preprint arXiv:2409.11376*.
- Wenjie Du, Jun Wang, Linglong Qian, Yiyuan Yang, Zina Ibrahim, Fanxing Liu, Zepu Wang, Haoxin Liu, Zhiyuan Zhao, Yingjie Zhou, et al. 2024. Tsi-bench: Benchmarking time series imputation. *arXiv preprint arXiv:2406.12747*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Nigel Fernandez, Alexander Scarlatos, and Andrew Lan. 2024. Syllabusqa: A course logistics question answering dataset. *arXiv preprint arXiv:2403.14666*.

- 66 66
- 6 6
- 67
- 68 68 68
- 683
- 685
- 68
- 68 69

694 695 696

698 699 700

701

702

- 703 704 705
- 7
- 708 709 710

711 712

713 714 715

- Elizabeth Fons, Rachneet Kaur, Soham Palande, Zhen Zeng, Tucker Balch, Manuela Veloso, and Svitlana Vyetrenko. 2024. Evaluating large language models on time series feature understanding: A comprehensive taxonomy and benchmark. *arXiv preprint arXiv:2404.16563*.
- Giacomo Frisoni, Alessio Cocchieri, Alex Presepi, Gianluca Moro, and Zaiqiao Meng. 2024. To generate or to retrieve? on the effectiveness of artificial contexts for medical open-domain question answering. *Association for Computational Linguistics(ACL)*.
- Wayne A Fuller. 2009. *Introduction to statistical time series*. John Wiley & Sons.
- Yizhen Guo, Tian Zhou, Wanyi Jiang, Bo Wu, Liang Sun, and Rong Jin. 2024. Maximizing the impact of deep learning on subseasonal-to-seasonal climate forecasting: The essential role of optimization. *arXiv preprint arXiv:2411.16728*.
- James D Hamilton. 2020. *Time series analysis*. Princeton university press.
- Songqiao Han, Xiyang Hu, Hailiang Huang, Minqi Jiang, and Yue Zhao. 2022. Adbench: Anomaly detection benchmark. *Advances in Neural Information Processing Systems*, 35:32142–32159.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060*.
- Noah Hollmann, Samuel Müller, Lennart Purucker, Arjun Krishnakumar, Max Körfer, Shi Bin Hoo, Robin Tibor Schirrmeister, and Frank Hutter. 2025. Accurate predictions on small data with a tabular foundation model. *Nature*, 637(8045):319–326.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Chuanyang Jin, Yutong Wu, Jing Cao, Jiannan Xiang, Yen-Ling Kuo, Zhiting Hu, Tomer Ullman, Antonio Torralba, Joshua B Tenenbaum, and Tianmin Shu. 2024a. Mmtom-qa: Multimodal theory of mind question answering. arXiv preprint arXiv:2401.08743.
- 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. 2023a. Timellm: Time series forecasting by reprogramming large language models. *arXiv preprint arXiv:2310.01728*.
- Ming Jin, Qingsong Wen, Yuxuan Liang, Chaoli Zhang, Siqiao Xue, Xue Wang, James Zhang, Yi Wang, Haifeng Chen, Xiaoli Li, et al. 2023b. Large models for time series and spatio-temporal data: A survey and outlook. *arXiv preprint arXiv:2310.10196*.

Ming Jin, Yifan Zhang, Wei Chen, Kexin Zhang, Yuxuan Liang, Bin Yang, Jindong Wang, Shirui Pan, and Qingsong Wen. 2024b. Position: What can large language models tell us about time series analysis. In *Forty-first International Conference on Machine Learning*. 718

719

721

722

723

724

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

753

754

755

758

759

760

761

762

763

764

765

766

767

768

769

770

771

- Yaxuan Kong, Yiyuan Yang, Shiyu Wang, Chenghao Liu, Yuxuan Liang, Ming Jin, Stefan Zohren, Dan Pei, Yan Liu, and Qingsong Wen. 2025. Position: Empowering time series reasoning with multimodal llms. *arXiv preprint arXiv:2502.01477*.
- Jennifer R Kwapisz, Gary M Weiss, and Samuel A Moore. 2011. Activity recognition using cell phone accelerometers. *ACM SigKDD Explorations Newsletter*, 12(2):74–82.
- Wing Lian, Bleys Goodson, Eugene Pentland, Austin Cook, Chanvichet Vong, and "Teknium". 2023. Openorca: An open dataset of gpt augmented flan reasoning traces. https://https://huggingface. co/Open-Orca/OpenOrca.
- Bryan Lim and Stefan Zohren. 2021. Time-series forecasting with deep learning: a survey. *Philosophical Transactions of the Royal Society A*, 379(2194):20200209.
- Haoxin Liu, Chenghao Liu, and B Aditya Prakash. 2024a. A picture is worth a thousand numbers: Enabling llms reason about time series via visualization. *NAACL*.
- Haoxin Liu, Shangqing Xu, Zhiyuan Zhao, Lingkai Kong, Harshavardhan Kamarthi, Aditya B Sasanur, Megha Sharma, Jiaming Cui, Qingsong Wen, Chao Zhang, et al. 2024b. Time-mmd: Multi-domain multimodal dataset for time series analysis. In *The Thirtyeight Conference on Neural Information Processing Systems Datasets and Benchmarks Track.*
- Haoxin Liu, Zhiyuan Zhao, Jindong Wang, Harshavardhan Kamarthi, and B. Aditya Prakash. 2024c. LST-Prompt: Large language models as zero-shot time series forecasters by long-short-term prompting. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 7832–7840. Association for Computational Linguistics.
- Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, et al. 2023. Summary of chatgpt-related research and perspective towards the future of large language models. *Meta-Radiology*, page 100017.
- Yong Liu, Haoran Zhang, Chenyu Li, Xiangdong Huang, Jianmin Wang, and Mingsheng Long. Timer: Generative pre-trained transformers are large time series models. In *Forty-first International Conference on Machine Learning*.
- Ziqing Ma, Wenwei Wang, Tian Zhou, Chao Chen, Bingqing Peng, Liang Sun, and Rong Jin. 2024. Fusionsf: Fuse heterogeneous modalities in a vector

774

quantized framework for robust solar power forecasting. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 5532–5543.

- Mike A Merrill, Mingtian Tan, Vinayak Gupta, Tom Hartvigsen, and Tim Althoff. 2024. Language models still struggle to zero-shot reason about time series. *arXiv preprint arXiv:2404.11757*.
- John A Miller, Mohammed Aldosari, Farah Saeed, Nasid Habib Barna, Subas Rana, I Budak Arpinar, and Ninghao Liu. 2024. A survey of deep learning and foundation models for time series forecasting. *arXiv preprint arXiv:2401.13912*.
- Navid Mohammadi Foumani, Lynn Miller, Chang Wei Tan, Geoffrey I Webb, Germain Forestier, and Mahsa Salehi. 2024. Deep learning for time series classification and extrinsic regression: A current survey. *ACM Computing Surveys*, 56(9):1–45.
- Michael Moor, Oishi Banerjee, Zahra Shakeri Hossein Abad, Harlan M Krumholz, Jure Leskovec, Eric J Topol, and Pranav Rajpurkar. 2023. Foundation models for generalist medical artificial intelligence. *Nature*, 616(7956):259–265.
- Ozan Baris Mulayim, Pengrui Quan, Liying Han, Xiaomin Ouyang, Dezhi Hong, Mario Bergés, and Mani Srivastava. 2024. 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, pages 169–173.
- Yuqi Nie, Yaxuan Kong, Xiaowen Dong, John M Mulvey, H Vincent Poor, Qingsong Wen, and Stefan Zohren. 2024. A survey of large language models for financial applications: Progress, prospects and challenges. *arXiv preprint arXiv:2406.11903*.
- Aditi Singh, Abul Ehtesham, Saket Kumar, and Tala Talaei Khoei. 2025. Agentic retrieval-augmented generation: A survey on agentic rag. *arXiv preprint arXiv:2501.09136*.
- Kuldeep Singh, Simerjot Kaur, and Charese Smiley.
  2024. Finqapt: Empowering financial decisions with end-to-end llm-driven question answering pipeline.
  In *Proceedings of the 5th ACM International Conference on AI in Finance*, pages 266–273.
- The Motley Fool. 2024. Earnings call transcripts.
  - Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
    - Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Musique: Multihop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554.

Volcengine. 2023. Doubao. Accessed: [Insert Date].

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

- Chengsen Wang, Qi Qi, Jingyu Wang, Haifeng Sun, Zirui Zhuang, Jinming Wu, Lei Zhang, and Jianxin Liao. 2024a. Chattime: A unified multimodal time series foundation model bridging numerical and textual data. *Preprint*, arXiv:2412.11376.
- Jun Wang, Wenjie Du, Yiyuan Yang, Linglong Qian, Wei Cao, Keli Zhang, Wenjia Wang, Yuxuan Liang, and Qingsong Wen. 2025. Deep learning for multivariate time series imputation: A survey. *Preprint*, arXiv:2402.04059.
- Yuxuan Wang, Haixu Wu, Jiaxiang Dong, Guo Qin, Haoran Zhang, Yong Liu, Yunzhong Qiu, Jianmin Wang, and Mingsheng Long. 2024b. Timexer: Empowering transformers for time series forecasting with exogenous variables. *arXiv preprint arXiv:2402.19072.*
- Qingsong Wen, Liang Sun, Fan Yang, Xiaomin Song, Jingkun Gao, Xue Wang, and Huan Xu. 2020. Time series data augmentation for deep learning: A survey. *arXiv preprint arXiv:2002.12478*.
- Qingsong Wen, Tian Zhou, Chaoli Zhang, Weiqi Chen, Ziqing Ma, Junchi Yan, and Liang Sun. 2022. Transformers in time series: A survey. *arXiv preprint arXiv:2202.07125*.
- Andrew Robert Williams, Arjun Ashok, Étienne Marcotte, Valentina Zantedeschi, Jithendaraa Subramanian, Roland Riachi, James Requeima, Alexandre Lacoste, Irina Rish, Nicolas Chapados, et al. 2024. Context is key: A benchmark for forecasting with essential textual information. arXiv preprint arXiv:2410.18959.
- Zhe Xie, Zeyan Li, Xiao He, Longlong Xu, Xidao Wen, Tieying Zhang, Jianjun Chen, Rui Shi, and Dan Pei. 2024. Chatts: Aligning time series with llms via synthetic data for enhanced understanding and reasoning. *arXiv preprint arXiv:2412.03104*.
- Jingwen Xu, Fei Lyu, and Pong C Yuen. 2023. Densityaware temporal attentive step-wise diffusion model for medical time series imputation. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, pages 2836–2845.
- Zhijian Xu, Yuxuan Bian, Jianyuan Zhong, Xiangyu Wen, and Qiang Xu. 2024. Beyond trend and periodicity: Guiding time series forecasting with textual cues. arXiv preprint arXiv:2405.13522.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.
- Yiyuan Yang, Rongshang Li, Qiquan Shi, Xijun Li, Gang Hu, Xing Li, and Mingxuan Yuan. 2023a. Sgdp: A stream-graph neural network based data prefetcher. In 2023 International Joint Conference on Neural Networks (IJCNN), pages 1–8. IEEE.

- Yiyuan Yang, Yi Li, Taojia Zhang, Yan Zhou, and Haifeng Zhang. 2021a. Early safety warnings for long-distance pipelines: A distributed optical fiber sensor machine learning approach. In *Proceedings* of the AAAI Conference on Artificial Intelligence, volume 35, pages 14991–14999.
- Yiyuan Yang, Chaoli Zhang, Tian Zhou, Qingsong Wen, and Liang Sun. 2023b. Dcdetector: Dual attention contrastive representation learning for time series anomaly detection. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 3033–3045.

895

900

901

902

903

904

905

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

923

924

927

931

933

936

937

938

939

- Yiyuan Yang, Haifeng Zhang, and Yi Li. 2021b. Long-distance pipeline safety early warning: a distributed optical fiber sensing semi-supervised learning method. *IEEE sensors journal*, 21(17):19453– 19461.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*.
- Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan Ning, Yu-Yang Liu, and Li Yuan. 2023. Llm lies: Hallucinations are not bugs, but features as adversarial examples. *arXiv preprint arXiv:2310.01469*.
- Zahra Zamanzadeh Darban, Geoffrey I Webb, Shirui Pan, Charu Aggarwal, and Mahsa Salehi. 2024. Deep learning for time series anomaly detection: A survey. *ACM Computing Surveys*, 57(1):1–42.
- Kexin Zhang, Qingsong Wen, Chaoli Zhang, Rongyao Cai, Ming Jin, Yong Liu, James Y Zhang, Yuxuan Liang, Guansong Pang, Dongjin Song, et al. 2024a.
  Self-supervised learning for time series analysis: Taxonomy, progress, and prospects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Weiqi Zhang, Jiexia Ye, Ziyue Li, Jia Li, and Fugee Tsung. 2024b. Dualtime: A dual-adapter multimodal language model for time series representation. *arXiv preprint arXiv:2406.06620*.
- Yunkai Zhang, Yawen Zhang, Ming Zheng, Kezhen Chen, Chongyang Gao, Ruian Ge, Siyuan Teng, Amine Jelloul, Jinmeng Rao, Xiaoyuan Guo, Chiang-Wei Fang, Zeyu Zheng, and Jie Yang. 2023. Insight miner: A large-scale multimodal model for insight mining from time series. In *NeurIPS 2023 AI for Science Workshop*.
- Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. 2021. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 11106–11115.
- Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. 2024. Toolqa: A dataset for llm question answering with external tools. *Advances in Neural Information Processing Systems*, 36.

### A Examples of TSQA Dataset

Below are several examples from the TSQA dataset.

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

### A.1 Forecasting.

**Question:** The Self-regulation of Slow Cortical Potentials dataset, provided by the University of Tuebingen, involves EEG recordings from a healthy subject controlling a cursor on a screen through cortical potentials. The sampling frequency of this dataset is 256 Hz. The input Time Series are *[Time Series Data Points]*. Please predict the next nine time series points given the information above.

**Answer:** Based on the given information, the predictions are [-11.69, -10.72, -9.84, -9.88, -9.03, -8.78, -8.41, -8.28, -7.59].

**Question:** From March 31, 2014 to June 04, 2014, the weighted ILI rate was 1.74542, and the ili total count reached 14637 across 1866 providers serving a total of 797863 patients, especially 0-4 age group was 3582, 5-24 age group was 5659, 25-49 age group was 3153, 50-64 age group was 1436, 65+ age group was 807, and 25-64 age group was unavailable. The internet search results for the same period are summarized as follows: ... The prior historical average weighted ILI rate during this period was 1.286616667 with a standard deviation of 0.073690536. Please analyze the data and provide a target result on June 04, 2014 with a summary of the findings.

**Answer:** Based on the provided textual information, I predict that the influenza infection situation in the United States for the upcoming period is likely to be moderate to severe, with a potential for increased risk of complications and co-infections, particularly among vulnerable populations such as those with cystic fibrosis and critically ill patients. Therefore, the predicted weighted ILI rate for June 04, 2014 is estimated to be 1.83453.

#### A.2 Imputation.

**Question:** This dataset, sourced from the Kaggle web traffic forecasting competition, contains daily time series data representing the number of hits for 145,063 Wikipedia pages from July 1, 2015, to September 10, 2017. The input Time Series are *[Time Series Data Points with Missing Values].* Please give full-time series with missing value imputed. The missing values are represented by 'X'.

1057

1058

1059

1060

1061

1040

1041

1042

1003 1004

- 1008
- 1010

1012

1013

1014

1015

1017 1018 1019

1021

1023

1024 1025 1026

1029 1030

1031

1033

1027

1034

1035

1037

1039

991

995

996

997

999

1000

1002

**Answer:** Based on the given information, the full-time series with missing value imputed are [Time Series Data Points].

## A.3 Classification.

Question: Following data provides accelerometer data for activity recognition research. The dataset has a sampling rate of 20Hz and records accelerometer data for six activity states: walking, jogging, sitting, standing, upstairs, and downstairs. Each sample includes acceleration values for the X, Y, and Z axes, ranging from -20 to 20, where 10 represents 1g and 0 indicates no acceleration. The recorded acceleration includes gravitational acceleration, so when the phone is stationary on a flat surface, the vertical axis registers approximately  $\neg \pm 10$ . We provide 10 timestamps of accelerometer data, with each timestamp containing X, Y, and Z values, for a total of 30 values. The recorded Time Serie is [Time Series Data Points]. Please judge whether this data corresponds to 'Walking' or 'Jogging' or 'Upstairs' or 'Downstairs' or 'Sitting' or 'Standing' given the information above.

Answer: Based on the given information, the activity is Walking.

## A.4 Anomaly Detection.

Question: The following data is derived from AWS CloudWatch services and is primarily used for monitoring and managing metrics related to cloud infrastructure and application performance. These include CPU usage, memory consumption, network throughput, and disk I/O metrics. This data can be utilized to identify performance issues in cloud services, such as traffic spikes, resource exhaustion, or network anomalies, and to perform root cause analysis for rapid problem resolution. The input Time Serie is [Time Series Data Points]. Please determine whether there are anomalies in this time series given the information above.

Answer: Based on the given information, this time series includes Anomaly Points.

# A.5 Open-Ended Reasoning QA.

Question: Summarize the overall behavior observed in the data points [Time Series Data Points].

The data initially shows a slow Answer: increase followed by a decreasing trend. The middle part of the data points represents a peak with few higher fluctuations and later shows a slightly decline indicating a volatile pattern.

**Question:** Determine if there's a structural break in the data points [Time Series Data Points]. A) Yes, between 0.71 and 0.72 B) Yes, between 0.65 and 0.63 C) No structural break

Answer: C) No structural break. The data follows a gradual upward trend and subsequent flattening with slight variability, but no abrupt change in the level or direction indicating a structural break.

**Question:** The data points [Time Series Data Points] exhibit a clear seasonal pattern. True or False?

Answer: False. While there are fluctuations and some repeated values, no clear cyclical pattern emerges that would indicate seasonality. The variations appear to be more random than driven by any seasonal influences.

#### **Open-Ended Reasoning QA Generation** B

Figure 7 illustrated the prompt for generating openended reasoning QA pairs.



Figure 7: Prompt for open-ended reasoning QA generation.

#### С Training Data Format

We format our question-and-answer pairs using a specifically designed template to clearly sep-1063 arate questions from answers. The template is 1064 structured as follows: <QUE> {Question} <ANS> 1065 {Answer} </END>. Additionally, for continual 1066 pre-training of the Llama model, we prepend 1067 <|begin\_of\_text|> at the start of each sample 1068 and append <|end\_of\_text|> at the end. For the 1069 Mistral model, we use <s> at the beginning of each sample and </s> at the end. In the case of the 1071

Qwen model, only <|endoftext|> is added at the end of each sample.

### D User Study

1075A total of 78 surveys were collected, and the results1076for each question are as follows:

**Question 1.** What is your field of study?



Question 4.Which do you think is the most1088accurate answer with respect to the open-ended1089reasoning question?1090

Identify any cyclical patterns in the data1092points [0.28, 0.27, 0.26, 0.25, 0.24, 0.23, 0.24,10930.25, 0.26, 0.28, 0.3, 0.33, 0.36, 0.39, 0.43, 0.46,10940.5, 0.53, 0.57, 0.6, 0.63, 0.65, 0.67, 0.68].1095



**Question 2.** How familiar are you with time series analysis? (1 = Not at all, 5 = Very familiar)



**Question 3.** How accurate is the answer with respect to the question?

Analyse and summarise the trend and pattern in the data set [68.0, 83.0, 95.0, 103.5, 108.5, 112.0, 113.0, 114.0, 114.5, 115.5, 116.0, 116.0, 115.0, 114.0, 113.5, 111.5, 108.0, 102.5, 100.5, 93.5, 88.0, 82.0, 78.0, 73.5].

| Answer 1  | Answer 2   | Answer 3  | Answer 4  |  |
|---|--|---|---|--|
| The data are reached as power discovered in<br>several and the are rivited products of the constant<br>relation in the several and the reaction of the constant<br>end of the several and the reaction of the constant<br>end of the several several reaction of the constant<br>end of the several several reaction of the constant<br>end of the several several reaction of the several<br>power and the several several reaction of the several<br>end of the several several several reaction of the several<br>end of the several several several several<br>end of the several several several several several<br>end of the several several several several several several<br>end of the several | A set of the second set of | The data series a location by each of other<br>shares and the base of the series.<br>The data series is a series of the series<br>of the series of the series of the series of the<br>series of the series of the series of the series<br>of the series of the series of the series of the<br>series of the series of the series of the series of the<br>series of the series of the series of the series of the<br>series of the series of the series of the series of the<br>series of the series of the series of the series of the<br>series of the series of the series of the series of the<br>series of the series of the series of the series of the<br>series of the series of the series of the series of the<br>series of the series of the series of the series of the<br>series of the series of the series of the series of the series of the<br>series of the series of the | Oward Thot The data citeday increases<br>the fill or a guard of the source in the<br>source of the source of the source in the<br>source of the source of the source in the<br>source of the source of the source of the source of the<br>source of the source of the source of the source of the<br>source of the source of the source of the source of the<br>source of the source of the source of the source of the<br>source of the source of the source of the source of the<br>source of the source of the source of the source of the source of the<br>source of the source of the source of the source of the source of the<br>source of the source |  |
| 1 = Completely Incorrect  | 2 = Incorrect But Wit  | th Some Correct Elements  | 3 = Correct   |  |

| Question 5.    | How complete is the answer? Did it |
|----------------|------------------------------------|
| address all pa | rts of the question?               |

Examine the data points [19.34, 20.41, 20.38, 19.28, 19.75, 22.84, 25.09, 24.97, 24.75, 26.25, 27.0, 25.03, 22.41, 21.78, 22.44, 21.19, 19.44, 19.62, 20.62, 20.16, 18.03, 17.19, 18.31, 19.91] and summarize the overall movement trend in this data.

| lasver 1   | Answer 2  | Answer 3   | Assessor 6   |
|--|---|--|--|
| The data allows and data strategy of the data of the d | <text><text><text><text><text><text><text></text></text></text></text></text></text></text> | The def point where a point is the solution of | In periodic to a consert invester to the site of the sector is the secto |
| 0 1 = Incomplete   | 2 = Somewhat Complete 3 = Cor   | nplete   |  |
| 0  |   |  |  |
|  |   |  |  |
|  |   |  |  |

**Question 6.** Which answer is the most compre-1106 hensive in relation to the question?

Considering the data points [47.0, 51.0, 51.0, 20.0, 48.0, 44.0, 41.0, 28.0, 48.0, 38.0, 67.0, 66.0, 70.0, 54.0, 63.0, 56.0, 76.0, 58.0, 60.0, 58.0, 69.0, 56.0, 39.0, 50.0], which statement best describes the volatility of this time series? A) The volatility is low with consistent data points. B) There is moderate volatility with regular patterns. C) High volatility with sudden peaks and drops. D) The time series exhibits a perfectly linear trend.

Question 6 (Answers):



| 1117 | <b>Question 7.</b> How clear or understa | indable is the |
|------|--|----------------|
| 1118 | reasoning provided in the answer?        |                |

Considering the data points [60.0, 28.0, 46.0, 51.0, 37.0, 46.0, 41.0, 35.0, 32.0, 66.0, 10.0, 21.0, 67.0, 72.0, 83.0, 61.0, 11.0, 31.0, 62.0, 50.0, 56.0, 76.0, 53.0, 86.0], do you see any seasonal patterns? A) Yes, there is a seasonal pattern. B) No, there is no identifiable seasonal pattern.

1126Question 8. How clearly does the response1127explain its reasoning? (1 = Not Clear At All, 5 =1128Really Clear)

What is the volatility level in the dataset [0.83, 0.8, 0.76, 0.75, 0.75, 0.76, 0.78, 0.8, 0.8, 0.8, 0.81, 0.8, 0.8, 0.8, 0.8, 0.79, 0.8, 0.82, 0.83, 0.82, 0.82, 0.82, 0.82, 0.81]? A) High B) Moderate C) Low D) No volatility"

1136Answer: C) Low Explanation: The dataset1137shows moderate stability with minor fluctuations,



indicating a low volatility level. The changes between consecutive points are generally small and consistent rather than erratic or extreme. Thinking process: The dataset has some movement, but it is not highly variable or erratic. There are no significant spikes or drops across the data points. The changes are steady and minor, suggesting a low level of volatility. Comparing to other options:

A) High: The dataset does not show high levels of volatility with major fluctuations.

B) Moderate: While there is some movement, it is not as severe as to classify it as moderate volatility.

D) No volatility: There is some movement, but it is not so minimal that it would be considered no volatility.

Therefore, the appropriate response and explanation for this question is C) Low volatility.





1163Summarise the overall movement trend in1164the data points [2.16, 1.86, 1.84, 2.49, 1.98, 1.96,11651.73, 1.79, 2.13, 1.43, 1.49, 1.31, 0.55, -0.16, -0.12,11660.63, 0.66, 0.51, 0.28, -1.14, -2.81, -3.21, -2.26,1167-1.65].



**Question 10.** How would you rate the quality of the explanation and reasoning with respect to the question? (1 = Really Poor, 5 = Excellent Reasoning)

In the time series data [0.29, 0.29, 0.29, 0.29, 0.29, 0.29, 0.29, 0.29, 0.29, 0.29, 0.28, 0.28, 0.28, 0.28, 0.28, 0.28, 0.28, 0.29, 0.3, 0.32, 0.34, 0.37, 0.4, 0.44, 0.47, 0.51], is there evidence of any structural breaks?

Answer: There is a subtle structural break between the stability around 0.28 and the upward shift starting around 0.3.

Explanation: A structural break in time series data is a sudden change in the data pattern or mean level. In the given data, after initially fluctuating around 0.29, there is a persistent drop to 0.28, and then a change occurs towards an upward trend from 0.3 onwards. Although subtle, this transition from 0.28 to 0.3 indicates a structural break where the mean level shifts. It signals the end of one data pattern and the beginning of another distinct phase.

1193Reasoning: The data stability around 0.291194suggests a constant level, followed by a sudden1195shift to 0.28, and then an upward trend from 0.31196onwards. This progressive change in behavior1197indicates there is indeed evidence of a structural1198break in the sequence.



**Question 11.** Overall, which answer do you prefer for this question? (You can choose one or more)

Given the data points [6.64, 1.71, 27.49, 1.7, 0.2, 6.05, 33.49, 21.13, 0.0, 0.99, 5.34, 10.74, 2.6, 11.04, 14.36, 9.08, 3.69, 0.26, 8.4, 6.37, 2.35, 6.02, 114.07, 49.47], describe the level of volatility found in this data.



**Question 12.** Which answer do you consider is the best?

The sequence [0.28, 0.3, 0.34, 0.39, 0.44, 0.49, 0.54, 0.58, 0.62, 0.65, 0.68, 0.7, 0.71, 0.72, 0.73, 0.73, 0.73, 0.73, 0.73, 0.72, 0.71, 0.7, 0.69, 0.68] exhibits an increasing pattern?

Answer 1: The sequence consistently increases until it peaks at 0.73, after which it slightly decreases, but remains relatively high.

Answer 2: The given sequence does not display a clear increasing pattern throughout. While there is an initial increase up to 0.73, followed by a period of stability, the later points decrease consistently from 0.73 to 0.68, breaking the continuity of an upward pattern.

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234 1235

1236

1237

1238

1239

1240

1241

Answer 3: The sequence does show an initial increase from 0.28 to 0.73, but it then stabilizes at 0.73 for several data points before decreasing slightly towards the end (down to 0.68). Therefore, while there is a phase of increase, the overall description suggests both rise and stability/decline, which contrasts with an exclusively increasing pattern. Hence, the correct response is False regarding continuous increase throughout.

Answer 4: The sequence initially increases from 0.28 to 0.73, but after reaching 0.73, it remains constant for a few terms and then starts decreasing from 0.73 to 0.68. Therefore, the sequence does not exhibit a continuously increasing pattern throughout.



 Answer 1: The sequence consistently increases until it peaks at 0.73, after which it slightly decreases, but remain...
 Answer 2: The given sequence does not display a clear increasing pattern throughout. While there is an initial inc...
 Answer 3: The sequence does show an initial increase from 0.28 to 0.73, but it...
 Answer 7: The sequence initially increases from 0.28 to 0.73, but after r...