

TIMERAG: IT'S TIME FOR RETRIEVAL-AUGMENTED GENERATION IN TIME-SERIES FORECASTING

Anonymous authors

Paper under double-blind review

ABSTRACT

Time-series data are essential for forecasting tasks across various domains. While Large Language Models (LLMs) have excelled in many areas, they encounter significant challenges in time-series forecasting, particularly in extracting relevant information from extensive temporal datasets. Unlike textual data, time-series data lack explicit retrieval ground truths, complicating the retrieval process. To tackle these issues, we present TimeRAG, a novel retrieval-augmented approach tailored for time-series forecasting. Our method uniquely applies to continuous and complex temporal sequences, and it is trained using LLM feedback, effectively addressing the absence of ground truth and aligning the priorities of the retriever and the LLM. Experimental results demonstrate the effectiveness of TimeRAG, highlighting its ability to significantly enhance forecasting performance and showcasing the potential of LLMs in time-series prediction tasks.

1 INTRODUCTION

Time-series data are fundamental for forecasting tasks across a broad range of domains, including weather prediction, energy consumption, healthcare monitoring, and financial markets (Yuan et al., 2024). For instance, meteorologists rely on historical climate data to forecast future weather conditions (Govett et al., 2024), energy providers predict demand based on past consumption patterns (Afzal et al., 2024), and healthcare professionals monitor patient vital signs over time to anticipate disease progression (Reed et al., 2005). In financial markets, time-series data such as stock prices, trading volumes, and interest rates are crucial for investment strategies and risk (Nelson et al., 2017).

Although Large Language Models (LLMs) have achieved remarkable success in various domains, they face significant challenges in time-series forecasting due to difficulties in extracting salient information from abundant temporal data. One mainstream solution involves contextual LLMs that incorporate sequences of historical data into the model's input to capture temporal dependencies (Jin et al., 2023; Yu et al., 2023). However, these models struggle with input length limitations and computational inefficiency, hindering their ability to model long-term dependencies effectively. Another approach is Retrieval-Augmented Generation (RAG), which allows LLMs to retrieve relevant information from external databases during generation (Aksitov et al., 2023). Yet, RAG faces challenges in time-series forecasting because its retrieval mechanisms are optimized for discrete textual data rather than continuous temporal data, making it difficult to align retrieved information with forecasting tasks and leading to suboptimal performance.

To address these challenges, we introduce **TimeRAG**, a novel retrieval-augmented approach specifically designed for time-series forecasting. A major difference between TimeRAG and previous RAG methods is that our method is the first to directly apply to continuous and complex temporal sequences. Yet a significant bottleneck arises due to the absence of explicit retrieval ground truths, unlike in textual data where relevant documents are clearly defined, making it challenging to train the retriever effectively. To overcome this, inspired by Zhang et al. (2023), we design a novel training target that leverages LLM feedback to guide the retrieval process. We utilize the generation probability of the LLMs for the correct tokens to determine positive and negative samples, which are then used for contrastive learning of the embedder. This approach aligns the retriever's prioritization with the LLM's assessments, bridging the gap between the information deemed important by the retriever and that recognized by the LLMs. Our methodology involves extracting sequences based on

the trained retriever, embedding them into the LLM’s input context by formatting time-series data in JSON to reduce the comprehension gap, and using this enriched context along with the original query to perform forecasting, effectively combining salient historical patterns with current data to improve prediction accuracy.

We evaluate TimeRAG on the task of stock movement prediction using four benchmark datasets of high-trade-volume stocks in U.S. markets: ACL18 (2014-2015) (Xu & Cohen, 2018), BIGDATA22 (2019-2020) (Soun et al., 2022), and CIKM18 (2017-2018) (Wu et al., 2018). To assess prediction performance on more recent stock data, we construct a new dataset, Stock23, which includes stock prices from 2022 to 2023. This addition ensures that our evaluation reflects current market conditions, offering a more comprehensive benchmark for modern stock prediction tasks. Our experimental results demonstrate that TimeRAG significantly outperforms conventional context-learning LLMs and existing RAG methods. This superior performance is attributed to TimeRAG’s ability to leverage LLM feedback to guide the retrieval process, effectively extracting and prioritizing historical sequences that enhance forecasting accuracy. By aligning the retrieval mechanism with the LLM’s predictive objectives, TimeRAG captures both short-term fluctuations and long-term dependencies inherent in financial markets, overcoming challenges posed by data volume and noise.

Our contributions are summarized as follows:

1. We introduce TimeRAG, the first retrieval-augmented generation approach specifically tailored for time-series forecasting.
2. TimeRAG leverages LLM feedback to improve information retrieval, and employs an outcome-oriented approach to filter relevant data from extensive historical contents.
3. Experimental results demonstrate that TimeRAG outperforms previous contextual and RAG methods in accuracy for stock price movement prediction on four real-world datasets, showcasing its unique ability to identify and utilize the most impactful sequences for time-series forecasting.

2 PROBLEM DEFINITION AND GOALS

Time-series forecasting predicts future values or trends G based on the given query sequence q and retrieved sequences c , where all the sequences are collected sequentially over time at regular intervals. The goal is to model the retrieve model R to efficiently retrieve useful information from a vast range of candidate sequences. In our finance example, the task is framed as a binary classification problem: predicting whether a stock’s price will *rise* or *fall* on the next trading day. The model is given a query sequence q , which represents the stock’s price over the previous t days. Using this query, the model retrieves relevant price sequences as context and then predicts the stock’s movement $M_{q,d}$ for the next trading day d . Table 1 defines the major symbols we use in this paper.

3 THE TIMERAG FRAMEWORK

Our method focuses on optimizing the retrieval stage to extract relevant content and seamlessly integrate it into LLMs. In the data construction phase (Section 3.1), we first preprocess time-series data and explore various features and prompts to maximize LLM performance. Then we use LLM feedback to identify the most effective data formats and content. During candidate selection (Section 3.2), we classify positive candidates based on high-performance feedback from the LLM, while the remaining data are treated as negative candidates. During training (Section 3.3), we employ knowledge distillation to teach the model how to distinguish useful time-series data (positive candidates) for a given query sequence, enabling more accurate and relevant retrieval. In this paper, we focus on stock movement prediction, but our approach can be applied to other time-series prediction tasks as well.

3.1 RAW DATA CONSTRUCTION

We utilize stock price data to perform stock movement prediction. First, we retrieve key stock price features from the Yahoo Finance API, including open price, high price, low price, adjusted close price, and volume. Next, we pre-process all stock price data into JSON format to improve the LLM’s

Table 1: The definition of symbols.

General Time-series Symbol	Stock Movement Prediction Symbol	Definitions
q	$q = \{q_{d-t}, \dots, q_{d-1}\}$	The query time-series data. In stock movement prediction, q refers to the query stock price sequence of length t , containing stock price data from trading day $d - t$ to trading day $d - 1$.
$G(q)$	$G(q, d) \in \{rise, fall\}$	The final output G given the query q . In stock movement prediction, $G(q, d)$ shows the generation G of the query stock q on the query trading day d , belonging to <i>rise</i> or <i>fall</i> .
$P(c) = LLM(O q, c)$	$P(c) = LLM(M_d q, c)$	The possibility P of the LLM to generate an accurate output O given the query sequence q and a candidate sequence c . In stock movement prediction, $P(c)$ refers to generating the accurate movement M on the query trading day d .
$\mathbb{C}_P = \{c_i \mid i = 1, \dots, k\}$		The set of top-k retrieved sequences as positive examples, where $P(c_i) \geq P(c_{i+1})$.
$\mathbb{C}_N = \{c_i \mid i = k + 1, \dots, n\}$		The set of negative retrieved sequences, where $P(c_i) \geq P(c_{i+1})$.
w_i		The soft weight of the i^{th} retrieved sequences, where $w_i = P(c_i), i = 1, \dots, k$.
R		The retrieve model.

ability to interpret time-series data (Fang et al., 2024; Singha et al., 2023; Yin et al., 2023). Finally, we explore different feature combinations and prompt designs to optimize the LLM’s performance.

3.1.1 DATA PREPROCESSING

We start by preprocessing all stock price data into five-day sequences and creating a list for each feature accordingly. For each trading day, we use a one-day sliding window. Following Yoo et al. (2021); Soun et al. (2022), movements are classified as a rise if they are greater than 0.55% and as a fall if they are less than -0.5%, based on the adjusted closing prices. If the movement falls between 0.55% and -0.5% during continuous trading days, we classify it as a freeze. It’s important to note that we don’t predict the freeze movement in query sequences, we only use it to calculate the recent movement list. In this way, we filter out minor and statistically insignificant price movements, thereby focusing on more significant trends. An example of a processed sequence is as follows:

```

148 {
149   "data_index": 1000010,
150   "query_stock": "ABBV",
151   "query_date": "2014-06-13",
152   "movement": "rise",
153   "date_list": ["2014-06-06", "2014-06-09", "2014-06-10", "2014-06-11", "2014-06-12"],
154   "open_list": [55.32, 54.42, 53.14, 53.85, 54.22],
155   "high_list": [55.4, 54.88, 54.08, 54.7, 54.25],
156   "low_list": [54.89, 53.72, 52.29, 53.75, 53.46],
157   "close_list": [55.1, 53.84, 53.97, 54.23, 53.66],
158   "adj_close_list": [36.85, 36.0, 36.09, 36.27, 35.88],
159   "volume_list": [3449800, 6297500, 8414700, 5386800, 3941600],
160   "movement_list": ["freeze", "fall", "freeze", "freeze", "fall"]
161 }
```

3.1.2 PROMPT SELECTION

To design an effective prompt with the most useful features, we optimize three components: the task definition prompt, query sequence representation (selecting valuable features from the query se-

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189

quence), candidate sequence representation (selecting valuable features from candidate sequences), and the order of these parts. To achieve this, we extract a toy dataset containing 5 rise sequences and 5 fall sequences, using all sequences from the same stock in the previous year as candidates. We then experiment with various prompts that include the task definition, each query sequence, and its corresponding candidate sequences, the detailed experiment is shown in Section 4.4. Using the probability $P(c)$ of generating the correct movement, we compute scores for each combination of task definition, query representation, and candidate representation. For each query sequence, we calculate the mean score of the top-3 $P(c)$ to assess prompt effectiveness. An example of a prompt trial is shown in Figure 1.

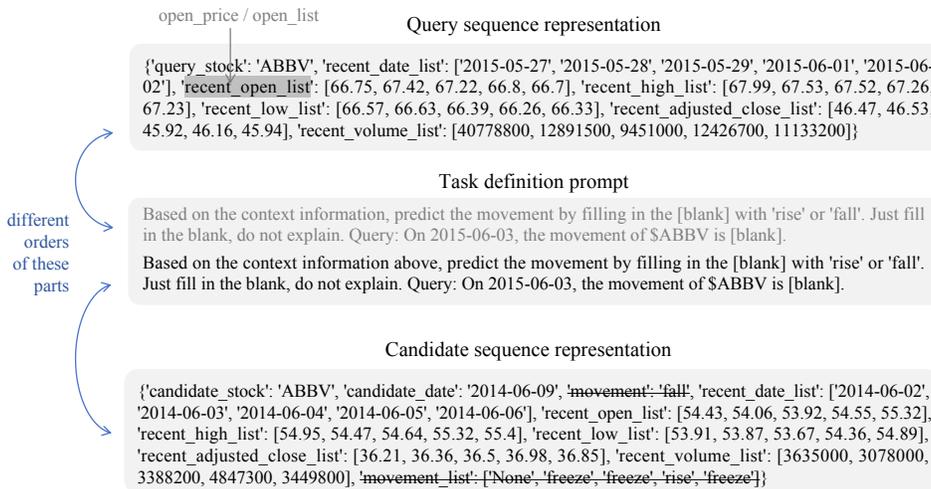


Figure 1: The prompt of TimeRAG.

190
191
192
193
194

Task definition prompt Utilizing the open-source LLaMA3-8b-instruct model, we first construct a fill-in-the-blank prompt designed to output only one token: 'rise' or 'fall'. This setup simplifies the calculation of the likelihood that the LLM generates the correct answer to the probability of the LLM producing the correct response at the first output index. Our selected task definition prompt is shown in Figure 1.

195
196
197
198
199

Query sequence representation We use the name of the stock and all recent price data to represent the query stock. We discuss how to list feature names to help the LLM better understand the referenced list. An example is highlighted in gray in Figure 1. In these trials, we name the list of open prices in the recent five days as 'open price', 'open list', or 'recent open list'. Once we definite feature names, we use the same name for candidate features.

200
201
202
203
204

Candidate sequence representation We discuss how different candidate features contribute to the prediction. We first provide all features and then provide the candidate sequence without each feature. After trials, we find the movement and recent movements of candidate sequences are noise for prediction. Therefore, we remove these two features, as is shown in Figure 1.

205 3.2 CANDIDATE SELECTION

206
207
208
209
210
211
212
213

In the last section, we confirmed the format and features we use in the query and candidate sequences. Then we need to select positive and negative candidate sequences for training. We consider all sequences from the past year of query stock as potential candidate sequences. Then we integrate the query sequence q and each candidate sequence c_i as the LLM input concurrently. The last step is to analyze the logits output by the LLM to calculate the probability of generating the correct response. The probability $P(c)$ indicates the probability of the LLM correctly predicting the stock price movement M_d for the queried trading day d .

214
215

We aim to train our retrieval model to retrieve sequences with a higher value of $P(c)$, thereby assisting the LLM in enhancing its prediction accuracy. To achieve this, we rearrange the candidate sequences in descending order according to $P(c)$, and select the top-1 sequences as a positive can-

216 didate and the last 15 sequences as negative candidates. The sets containing selected positive and
 217 negative sequences are denoted as \mathbb{C}_P and \mathbb{C}_N . Moreover, the value of $P(c)$ serves as the teacher
 218 score, and we directly use it as the training reward for the corresponding candidate sequence c .
 219

220 3.3 RETRIEVER TRAINING

221
 222 Our retriever $R(q)$ is designed to intelligently distinguish between historically significant sequences
 223 \mathbb{C}_P and noisy sequences \mathbb{C}_N , based on their support to the current query sequence q . Specifically,
 224 the process can be mathematically represented as:

$$225 R(q) = \arg \max_{s \in \mathbb{C}_P \cup \mathbb{C}_N} \text{sup}(q, s) \quad (1)$$

227 This formulation ensures the identification and extraction of sequences that maximize the measure
 228 $\text{sup}(q, s)$, where \mathbb{C}_P contains the top-k sequences with the highest scores, deemed most predictive
 229 of the future stock movement, and \mathbb{C}_N encompasses sequences with lower predictive utility. This
 230 approach not only facilitates more accurate and contextually rich predictions by focusing on the
 231 most informative historical sequences but also enhances the model’s adaptability to evolving market
 232 conditions, thereby providing a robust framework for financial time-series analysis.

233 To train the retriever, we employ the pairs (q, c_i) as soft labels. The samples within \mathbb{C}_P are treated as
 234 positive examples, while the candidates in \mathbb{C}_N are considered negative examples. To underscore the
 235 importance of the LLM outputting the correct price movement, we use the training reward as a soft
 236 weight, denoted as $w_i = P(c)$. It allows the model to weigh the training examples based on their
 237 likelihood of being correct. This nuanced approach ensures that the model pays more attention to
 238 sequences that not only are ranked higher but also have a higher probability of predicting the correct
 239 price movement, thereby fine-tuning its predictive capabilities.

240 To learn from soft rewards derived from the LLM, we conduct knowledge distillation. Particularly,
 241 we employ the KL-divergence to minimize the gap between the distributions of candidates computed
 242 using LLM’s rewards and those predicted by the embedding model. In particular, for each query q
 243 and its candidate list $\{\mathbb{C}_P, \mathbb{C}_N\}$, we derive the LLM’s rewards towards the candidates, denoted as
 244 $\{P(c_i), i = 1, \dots, n\}$. To make the LLM’s rewards suitable for distillation, we transform each reward
 245 into a normalized weight: $w_i = \text{softmax}_R(P(c_i)/\alpha)$, where α represents the temperature. On top
 246 of these elements, the KL divergence is computed by the following equation:

$$247 \min_c \cdot \sum_c -w_i \times \log \left(\frac{\exp(\langle e_q, e_{c_i} \rangle / \tau)}{\sum_{c' \in \mathbb{C}} \exp(\langle e_q, e_{c'} \rangle / \tau)} \right) \quad (2)$$

251 This loss function is designed to optimize the similarity between the query embedding and the
 252 embeddings of the top-ranked reference candidates, thereby enhancing the model’s ability to predict
 253 stock price movements accurately.
 254

255 4 EXPERIMENT

257 4.1 EXPERIMENTAL SETTINGS

259 **Datasets.** We evaluate TimeRAG on four benchmark datasets consisting of high-trade-volume
 260 stocks in US stock markets: 1) **ACL18** (Xu & Cohen, 2018) consists of 71 stocks along with their
 261 tweets and historical price data from 2014.06.02 to 2015.12.31; 2) **BIGDATA22** (Soun et al., 2022)
 262 consists of 47 stocks along with their tweets and historical price data from 2019.04.01 to 2020.12.31;
 263 3) **CIKM18** (Wu et al., 2018) consists of 41 stocks along with their tweets and historical price data
 264 from 2017.01.03 to 2018.01. 23; 4) **Stock23** consists of 51 stocks along with their historical price
 265 data from 2022.01.03 to 2023.12.31. The detailed statistics are summarized in Table 2.

266 **Baselines** We evaluate whether the accuracy of the LLM’s stock predictions is enhanced by incor-
 267 porating example sequences selected through retrieval models, compared to approaches that use
 268 random sampling or no examples. We evaluate 4 retrieval methods in this setting: 1) **Instructor**
 269 (Su et al., 2023), a 1.5B instruction-finetuned text embedder. 2) **BGE** (BAAI general embedding)
 (Xiao et al., 2023), a 335M general embedder pre-trained from RetroMAE (Shitao Xiao, 2022). 3)

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

Table 2: Test dataset statistics.

	stock amount	all sequences						query sequences			
		trading date						trading date			
		all	rise	fall	all	rise	fall				
ACL18	33	2014.06.02-2015.12.31	7629	3840	3789	2015.06.03-2015.12.31	2690	1345	1345		
BIGDATA22	22	2019.04.01-2020.12.31	6534	3412	3122	2020.04.09-2020.12.31	2800	1400	1400		
CIKM18	19	2017.01.03-2018.01.23	2213	1228	985	2018.01.03-2018.01.23	80	40	40		
Stock23	51	2022.01.03-2023.12.31	19283	9627	9656	2023.01.03-2023.12.31	4128	2064	2064		

LLM-Embedder, a 109M embedder fine-tuned with the feedback from LLMs. 4) **E5-mistral-7b-instruct** (Wang et al., 2023), a 7B embedder initialized from Mistral-7B-v0.1 (Jiang et al., 2023a) and fine-tuned on a mixture of multilingual datasets.

Evaluation Metrics. We employ Accuracy (ACC) and Matthews Correlation Coefficient (MCC) (Matthews, 1975) to assess the performance of TimeRAG and the baseline models on the stock movement prediction task. These metrics evaluate the performance of stock movement prediction based on the distribution of positive and negative samples. ACC and MCC are defined as:

$$ACC = \frac{TP + TN}{TP + TN + FP + FN} \quad (3)$$

$$MCC = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (4)$$

where TP denotes true positives, TN denotes true negatives, FP denotes false positives, and FN denotes false negatives.

Implementation Details In our implementation, two key factors play a crucial role: the LLM foundation and the embedding model backbone. We choose LLaMA3-8b-instruct for feedback, as it is new, open-source, and powerful. For the embedding backbone, we use BGE, well-pretrained in general text embedding tasks, providing a strong foundation for TimeRAG.

4.2 MAIN RESULTS

Based on the results presented in Table 3, TimeRAG notably outperforms all evaluated approaches across the ACL18, BIGDATA22, CIKM18, and Stock23 datasets in terms of Matthews Correlation Coefficient (MCC). Specifically, TimeRAG achieves MCC scores of 0.140, 0.145, 0.197, and 0.219 respectively, which significantly surpass those of other methods. In contrast, the remaining baseline methods demonstrate much lower performance, with many yielding results close to random guessing, as indicated by their near-zero or negative MCC values. This strong performance of TimeRAG highlights its effectiveness in predicting stock movements accurately, underscoring the value of our method compared to traditional approaches and even advanced models like GPT-4 and LLaMA3-8b-instruct.

Table 3: Results of stock movement predictions using LLMs and retrieval models. The asterisk (*) indicates the LLM employed while using retrieval models.

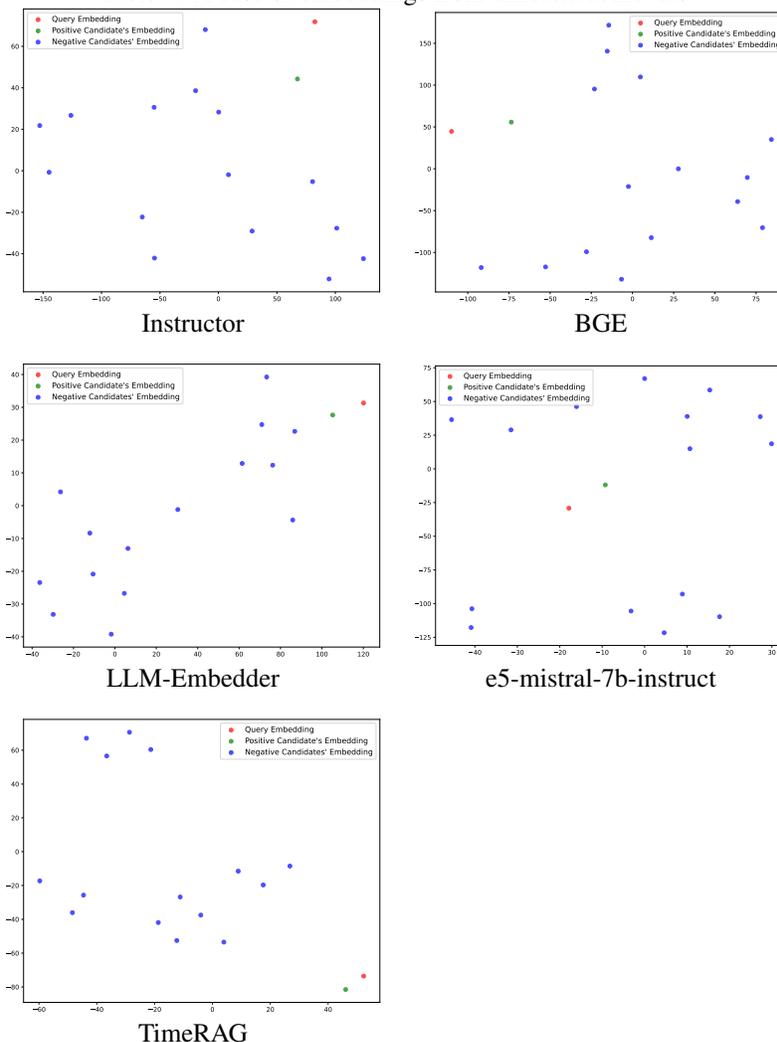
Methods	ACL18		BIGDATA22		CIKM18		Stock23	
	ACC	MCC	ACC	MCC	ACC	MCC	ACC	MCC
LLaMA2-7B-chat	0.500	0.010	0.499	0.000	0.500	0.056	0.500	0.000
GPT-4	<u>0.524</u>	<u>0.049</u>	<u>0.522</u>	<u>0.044</u>	0.400	-0.231	0.525	0.050
FinMA-7B-full	0.500	0.001	0.508	0.022	0.575	0.197	0.497	-0.009
LLaMA3-8b-instruct(*)	0.522	0.048	0.497	-0.006	0.475	-0.070	<u>0.527</u>	<u>0.067</u>
random retrieval*	0.501	0.008	0.499	-0.012	0.500	0.000	0.501	0.014
Instructor*	0.500	0.003	0.501	0.011	0.487	-0.066	0.501	0.018
BGE*	0.501	0.005	0.502	0.015	0.475	-0.095	0.501	0.015
LLM-Embedder*	0.508	0.025	0.501	0.003	0.512	0.052	0.511	0.055
e5-mistral-7b-instruct*	0.514	0.044	0.503	0.011	0.450	-0.190	0.504	0.018
TimeRAG*	0.554	0.140	0.541	0.145	0.537	0.197	0.546	0.219

324 Despite significantly improved ACC in large sample datasets, our model also achieves positive MCC
 325 across all datasets, indicating that TimeRAG effectively retrieves valuable candidates to assist the
 326 LLM in analyzing stock sequences and predicting stock movements. Compared to GPT-4, our
 327 model’s enhanced performance underscores the importance of these candidates. It indicates that us-
 328 ing only the query sequence is insufficient to predict movements. Moreover, compared to our LLM
 329 foundation, LLaMA3, our results demonstrate the effectiveness of retriever training. Furthermore,
 330 our improvements over other retrieval methods highlight the benefits of task-oriented fine-tuning on
 331 stock data.

332
 333 4.3 CASE STUDY

334
 335 In this section, we investigate the ability of our retriever and baseline models to differentiate between
 336 various time-series data. We focus on the first example from the CIKM18 test dataset. As shown
 337 in Table 4, our retriever identifies positive candidates that are notably closer to the query sequence
 338 compared to the negative candidates, which are significantly more distant. This result highlights the
 339 superior retrieval performance of our method.

340
 341 Table 4: A case of embeddings from different retrievers.



364
 365
 366
 367
 368
 369
 370
 371
 372
 373
 374
 375
 376 The proximity of the positive candidate to the query and the negative candidates are further away,
 377 implying that our model is effectively capturing the subtle patterns and dynamics in the time-series
 data, enabling more accurate retrieval. In contrast, baseline methods struggle to differentiate be-

tween similar but less relevant time-series sequences. This disparity illustrates the advantage of our retriever in isolating meaningful sequences for prediction tasks. Additionally, these findings reinforce the importance of a tailored retrieval strategy for time-series forecasting, where the subtle nuances in the data can significantly impact predictive accuracy.

4.4 ABLATION STUDY

In this section, we interpret how we conduct prompt selection shown in Section 3.1.2, by exploring the order of instruction, query sequence, and candidate sequence; exploring the name of features; and exploring which feature is important.

4.4.1 PROMPT SELECTION

Table 5 reveals that the sequence and content of prompts significantly impact the performance of stock movement prediction. The configuration marked as 6', which follows the 'qtc' order (query first, followed by task, then candidate) and includes the 'recent_xxx_list' feature without additional name and date details, scores the highest at 0.866. This indicates that specify temporal dynamics in the query significantly enhances prediction accuracy. It demonstrates that focusing on recent movement data in the query sequence and adhering to a structured prompt order optimizes the model's predictive capabilities. Thus, for higher prediction accuracy in stock movement, it is crucial to prioritize the incorporation of recent performance data and maintain a consistent structure in prompt arrangement.

Table 5: Scores for different prompts. *m list* refers to the recent movement list. Features labeled as *xxx_price* follow this naming format (e.g., *open_price*), while *xxx_list* indicates that features are named in the form of *xxx_list* (e.g., *open_list*).

index	order	query sequence features				candidate sequence features				score
		name	date	feature name	recent move- ment list	name	date	recent move- ment list	<i>d</i> day's movement	
1	tqc	w/o	w/o	xxx_price	w/o	w/o	w/o	w/o	w/o	0.614
2	tqc	w/o	w/o	xxx_list	w/o	w/o	w/o	w/o	w/o	0.744
3	tqc	w	w/o	recent_xxx_list	w/o	w	w	w/o	w	0.629
4	tqc	w	w	recent_xxx_list	w/o	w	w	w/o	w	0.644
5	cqt	w	w/o	recent_xxx_list	w/o	w	w	w/o	w	0.783
6	qtc	w	w/o	recent_xxx_list	w/o	w	w	w/o	w	0.814
6'	qtc	w	w/o	recent_xxx_list	w/o	w	w	w/o	w/o	0.866
7	tqc	w/o	w/o	xxx_list	w	w/o	w/o	w	w/o	0.756
8	tqc	w	w/o	recent_xxx_list	w	w	w	w	w	0.771

Another intriguing observation is that including the movement list and factoring in the movement of candidate data consistently results in lower scores. This suggests that LLMs predict stock movements by deeply analyzing the sequence data itself, rather than superficially following trends.

4.4.2 KEY CHARACTERISTICS FOR CANDIDATES

Table 6 presents an ablation study from prompt 6, analyzing the impact of removing various candidate features on the prediction score. The original score with all features included is 0.814. Removing the candidate data entirely results in a significant score drop to 0.530, indicating that candidate features are crucial for accurate predictions. Similarly, removing date and stock price data such as open, high, low, close, and volume also decreases the score, though to a lesser extent. The smallest score reductions occur when removing volume (0.083) and date (0.133), suggesting these features are less critical but still contribute positively to the model's performance.

Table 6: Score change when removing candidate features.

prompt 6	w/o candidate	w/o movement	w/o date	w/o open	w/o high	w/o low	w/o close	w/o volume
0.814	0.530	0.866	0.681	0.700	0.698	0.706	0.700	0.731
	↓ 0.284	↑ 0.052	↓ 0.133	↓ 0.114	↓ 0.116	↓ 0.108	↓ 0.114	↓ 0.083

432 Interestingly, removing movement information leads to a score increase to 0.866. This suggests that
433 movement data might act as noise, distracting the model from more predictive patterns found in other
434 sequence data. This finding implies that focusing on static features like price points and volume
435 might enable a more robust analysis of stock movements, as these elements provide foundational
436 data that the model can utilize more effectively than dynamic movement information.

437 438 439 5 RELATED WORK

440 441 5.1 TIME-SERIES FORECASTING WITH LLMs

442
443 To enhance the performance of LLMs in time-series forecasting, existing methods focus on the
444 alignment of temporal and textual data, turning time-series into textual format, or encoding it and
445 textual data into a unified vector space. For instance, Jin et al. (2023) reprogram time-series data into
446 textual representations suitable for LLMs, enhancing prediction accuracy via declarative prompts.
447 Similarly, Yu et al. (2023) and Liu et al. (2024) explore cross-modal alignment, with the former
448 applying LLMs to financial forecasting using stock prices and news data, and the latter introducing
449 a cross-modality framework to align time-series with text for improved predictive performance.
450 Expanding on this, Pan et al. (2024) map time-series and text into a shared semantic space, further
451 boosting LLM performance by strengthening data alignment. Despite advancements in time-series
452 forecasting, many still require insights from extensive time-series data that cannot all be input into
453 LLMs simultaneously. This limitation creates a need for retrieval-augmented methods, which our
454 approach specifically addresses.

455 456 5.2 RETRIEVAL-AUGMENTED LLMs

457
458 To enhance LLM reasoning and prediction performance by retrieving relevant information from vast
459 datasets, numerous retrieval methods have been proposed (Fan et al., 2024). Early approaches were
460 based on keyword frequency, with many studies directly applying BM25 for passage-level retrieval
461 in RAG (Chen et al., 2017; Jiang et al., 2023b; Ram et al., 2023; Xu et al., 2024; Zhong et al.,
462 2022; Zhou et al., 2022). These passages were represented as bags of words and ranked using term
463 frequency-inverse document frequency (TF-IDF) (Izcard & Grave, 2021). Later, methods based
464 on semantic similarity emerged, encoding queries and passages into a unified vector space (Li &
465 Qiu, 2023; Lu et al., 2023; Milios et al., 2023; Poesia et al., 2022; Rubin et al., 2022; Ye et al.,
466 2023), intending to train embeddings to bring queries and factual passages as close as possible.
467 However, these approaches are not well-suited for time-series retrieval, such as predicting stock price
468 movements, where there are no fixed factual passages to retrieve. Moreover, due to the highly similar
469 nature of time-series data, semantic similarity-based methods struggle to differentiate between them.
470 Therefore, a specialized retrieval method for time-series data is required, which our model provides.

471 472 6 CONCLUSION

473
474 In this work, we present TimeRAG, a novel retrieval-augmented generation (RAG) approach de-
475 signed specifically for financial time-series forecasting, with a focus on stock movement prediction.
476 TimeRAG enhances the ability of large language models (LLMs) to interpret time-series data by
477 integrating feedback mechanisms that address the lack of a clear retrieval ground truth. Our method
478 bridges the gap between the information deemed important by the retriever and that recognized
479 by the LLM, enabling a deeper understanding of market dynamics. We evaluate TimeRAG on
480 four benchmark datasets of high-trade-volume stocks in the US markets—ACL18, BIGDATA22,
481 CIKM18, and our newly constructed Stock23. Experimental results demonstrate that TimeRAG
482 significantly outperforms conventional context-learning LLMs in prediction accuracy. This superi-
483 ority is attributed to TimeRAG’s unique ability to filter relevant time-series data from extensive and
484 noisy historical datasets, employing an outcome-oriented retrieval approach that identifies sequences
485 that most significantly enhance forecasting performance. Our findings underscore the potential of
TimeRAG to advance time-series analysis in financial contexts, addressing the challenges faced by
existing methods.

REFERENCES

- 486
487
488 Sadeh Afzal, Afshar Shokri, Behrooz M Ziapour, Hamid Shakibi, and Behnam Sobhani. Building
489 energy consumption prediction and optimization using different neural network-assisted models;
490 comparison of different networks and optimization algorithms. *Engineering Applications of Arti-*
491 *ficial Intelligence*, 127:107356, 2024.
- 492 Renat Aksitov, Chung-Ching Chang, David Reitter, Siamak Shakeri, and Yunhsuan Sung. Charac-
493 terizing attribution and fluency tradeoffs for retrieval-augmented large language models. *arXiv*
494 *preprint arXiv:2302.05578*, 2023.
- 495 Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading Wikipedia to answer open-
496 domain questions. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th Annual*
497 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1870–
498 1879, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/
499 v1/P17-1171. URL <https://aclanthology.org/P17-1171>.
- 500 Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and
501 Qing Li. A survey on rag meeting llms: Towards retrieval-augmented large language models. In
502 *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*,
503 pp. 6491–6501, 2024.
- 504 Xi Fang, Weijie Xu, Fiona Anting Tan, Jiani Zhang, Ziqing Hu, Yanjun Qi, Scott Nickleach, Diego
505 Socolinsky, Srinivasan Sengamedu, and Christos Faloutsos. Large language models on tabular
506 data—a survey. *arXiv preprint arXiv:2402.17944*, 2024.
- 507 Mark Govett, Bubacar Bah, Peter Bauer, Dominique Berod, Veronique Bouchet, Susanna Corti,
508 Chris Davis, Yihong Duan, Tim Graham, Yuki Honda, et al. Exascale computing and data han-
509 dling: Challenges and opportunities for weather and climate prediction. *Bulletin of the American*
510 *Meteorological Society*, 2024.
- 511 Gautier Izacard and Édouard Grave. Leveraging passage retrieval with generative models for open
512 domain question answering. In *Proceedings of the 16th Conference of the European Chapter of*
513 *the Association for Computational Linguistics: Main Volume*, pp. 874–880, 2021.
- 514 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
515 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
516 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023a.
- 517 Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang,
518 Jamie Callan, and Graham Neubig. Active retrieval augmented generation. In Houda Bouamor,
519 Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Meth-*
520 *ods in Natural Language Processing*, pp. 7969–7992, Singapore, December 2023b. Associa-
521 tion for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.495. URL <https://aclanthology.org/2023.emnlp-main.495>.
- 522 Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y Zhang, Xiaoming Shi, Pin-Yu Chen, Yux-
523 uan Liang, Yuan-Fang Li, Shirui Pan, et al. Time-llm: Time series forecasting by reprogramming
524 large language models. *arXiv preprint arXiv:2310.01728*, 2023.
- 525 Xiaonan Li and Xipeng Qiu. Mot: Memory-of-thought enables chatgpt to self-improve. In *Pro-*
526 *ceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp.
527 6354–6374, 2023.
- 528 Chenxi Liu, Qianxiong Xu, Hao Miao, Sun Yang, Lingzheng Zhang, Cheng Long, Ziyue Li, and Rui
529 Zhao. Timecma: Towards llm-empowered time series forecasting via cross-modality alignment.
530 *arXiv preprint arXiv:2406.01638*, 2024.
- 531 Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter
532 Clark, and Ashwin Kalyan. Dynamic prompt learning via policy gradient for semi-structured
533 mathematical reasoning. In *The Eleventh International Conference on Learning Representations*,
534 2023.

- 540 Brian W Matthews. Comparison of the predicted and observed secondary structure of t4 phage
541 lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure*, 405(2):442–451, 1975.
542
- 543 Aristides Milios, Siva Reddy, and Dzmitry Bahdanau. In-context learning for text classification with
544 many labels. In *Proceedings of the 1st GenBench Workshop on (Benchmarking) Generalisation*
545 *in NLP*, pp. 173–184, 2023.
- 546 David MQ Nelson, Adriano CM Pereira, and Renato A De Oliveira. Stock market’s price movement
547 prediction with lstm neural networks. In *2017 International joint conference on neural networks*
548 *(IJCNN)*, pp. 1419–1426. Ieee, 2017.
549
- 550 Zijie Pan, Yushan Jiang, Sahil Garg, Anderson Schneider, Yuriy Nevmyvaka, and Dongjin Song.
551 s^2 ip-llm: Semantic space informed prompt learning with llm for time series forecasting. In
552 *Forty-first International Conference on Machine Learning*, 2024.
- 553 Gabriel Poesia, Alex Polozov, Vu Le, Ashish Tiwari, Gustavo Soares, Christopher Meek, and Sumit
554 Gulwani. Synchromesh: Reliable code generation from pre-trained language models. In *Internat-*
555 *ional Conference on Learning Representations*, 2022.
556
- 557 Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and
558 Yoav Shoham. In-context retrieval-augmented language models. *Transactions of the Association*
559 *for Computational Linguistics*, 11:1316–1331, 2023. doi: 10.1162/tacl.a_00605. URL [https://](https://aclanthology.org/2023.tacl-1.75)
560 aclanthology.org/2023.tacl-1.75.
- 561 Matt J Reed, CE Robertson, and PS Addison. Heart rate variability measurements and the prediction
562 of ventricular arrhythmias. *Qjm*, 98(2):87–95, 2005.
563
- 564 Ohad Rubin, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context
565 learning. In *Proceedings of the 2022 Conference of the North American Chapter of the Associa-*
566 *tion for Computational Linguistics: Human Language Technologies*, pp. 2655–2671, 2022.
- 567 Yingxia Shao Zhao Cao Shitao Xiao, Zheng Liu. Retromae: Pre-training retrieval-oriented lan-
568 guage models via masked auto-encoder. In *EMNLP*, 2022. URL [https://arxiv.org/](https://arxiv.org/abs/2205.12035)
569 [abs/2205.12035](https://arxiv.org/abs/2205.12035).
570
- 571 Ananya Singha, José Cambronero, Sumit Gulwani, Vu Le, and Chris Parnin. Tabular representa-
572 tion, noisy operators, and impacts on table structure understanding tasks in llms. *arXiv preprint*
573 *arXiv:2310.10358*, 2023.
- 574 Yejun Soun, Jaemin Yoo, Minyong Cho, Jihyeong Jeon, and U Kang. Accurate stock movement
575 prediction with self-supervised learning from sparse noisy tweets. In *2022 IEEE International*
576 *Conference on Big Data (Big Data)*, pp. 1691–1700. IEEE, 2022.
577
- 578 Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen-tau Yih,
579 Noah A. Smith, Luke Zettlemoyer, and Tao Yu. One embedder, any task: Instruction-finetuned
580 text embeddings. In *ACL (Findings)*, pp. 1102–1121. Association for Computational Linguistics,
581 2023.
- 582 Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Improv-
583 ing text embeddings with large language models. *arXiv preprint arXiv:2401.00368*, 2023.
584
- 585 Huizhe Wu, Wei Zhang, Weiwei Shen, and Jun Wang. Hybrid deep sequential modeling for so-
586 cial text-driven stock prediction. In *Proceedings of the 27th ACM international conference on*
587 *information and knowledge management*, pp. 1627–1630, 2018.
- 588 Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. C-pack: Packaged resources to
589 advance general chinese embedding, 2023.
590
- 591 Fangyuan Xu, Weijia Shi, and Eunsol Choi. RECOMP: Improving retrieval-augmented LMs
592 with context compression and selective augmentation. In *The Twelfth International Confer-*
593 *ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=m1JLVigNHp)
[m1JLVigNHp](https://openreview.net/forum?id=m1JLVigNHp).

594 Yumo Xu and Shay B Cohen. Stock movement prediction from tweets and historical prices. In
595 *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume*
596 *1: Long Papers)*, pp. 1970–1979, 2018.

597
598 Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. Compositional exemplars
599 for in-context learning. In *International Conference on Machine Learning*, pp. 39818–39833.
600 PMLR, 2023.

601 Yuwei Yin, Yazheng Yang, Jian Yang, and Qi Liu. Finpt: Financial risk prediction with profile
602 tuning on pretrained foundation models. *arXiv preprint arXiv:2308.00065*, 2023.

603
604 Jaemin Yoo, Yejun Soun, Yong-chan Park, and U Kang. Accurate multivariate stock movement
605 prediction via data-axis transformer with multi-level contexts. In *Proceedings of the 27th ACM*
606 *SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 2037–2045, 2021.

607 Xinli Yu, Zheng Chen, Yuan Ling, Shujing Dong, Zongyi Liu, and Yanbin Lu. Temporal data meets
608 llm—explainable financial time series forecasting. *arXiv preprint arXiv:2306.11025*, 2023.

609
610 Chenhan Yuan, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. Back to the future: Towards
611 explainable temporal reasoning with large language models. In *Proceedings of the ACM on Web*
612 *Conference 2024*, pp. 1963–1974, 2024.

613 Peitian Zhang, Shitao Xiao, Zheng Liu, Zhicheng Dou, and Jian-Yun Nie. Retrieve anything to
614 augment large language models. *arXiv preprint arXiv:2310.07554*, 2023.

615
616 Zexuan Zhong, Tao Lei, and Danqi Chen. Training language models with memory augmentation.
617 In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*,
618 pp. 5657–5673, 2022.

619 Shuyan Zhou, Uri Alon, Frank F Xu, Zhengbao Jiang, and Graham Neubig. Docprompting: Gener-
620 ating code by retrieving the docs. In *The Eleventh International Conference on Learning Repre-*
621 *sentations*, 2022.

622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647