ZEROTS: ZERO-SHOT TIME SERIES FORECASTING VIA MULTI-PARTY DATA-MODEL INTERACTION

Anonymous authors

Paper under double-blind review

Abstract

Time series forecasting (TSF) is a fundamental task in artificial intelligence, with applications ranging from weather prediction, stock market analysis to electricity demand forecasting. While existing models, particularly large language models (LLMs) tailored for TSF, primarily focus on improving accuracy and generalization through pre-training and fine-tuning, zero-shot prediction without taskspecific fine-tuning, still remains underexplored. This limitation arises from the restricted scalability and flexibility of current LLMs for TSF, which struggle to fully capture the interactions between data and model. In this work, we introduce ZeroTS, a novel approach that bridges open-world knowledge with inherent data regularities by constructing multi-party interactions between data and models. On the data side, we propose a TS-RAG (Retrieval-Augmented Generation for Time Series), which efficiently retrieves both meta and series information, enabling diverse domain-specific time series to be used as prompts. On the model side, we develop a reinforcement learning framework that treats ground-truth as environments, providing error feedback to optimize a smaller model and harnessing the capabilities of LLMs. This allows ZeroTS to incrementally approach inherent data regularities while iteratively refining its outputs. We validate ZeroTS via extensive experiments on zero-shot and long-horizon forecasting. ZeroTS achieves best or second best results with comparative parameters, 1/4 memory and 1/7 inference speed, demonstrating its efficiency and effectiveness. Our results highlight the potential of Data-LLM interactions for zero-shot learning with acceptable parameters, opening new avenues on research of this underexplored area.

031 032 033

034

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

1 INTRODUCTION

Time series forecasting (TSF) is one of the fundamental capacities for data-driven artificial intelli-035 gence, which facilitates the progress of various urban and scientific applications such as weather forecasting (Jin et al., 2023), stock prediction (Singh & Srivastava, 2017) and electricity plan-037 ning (Reddy et al., 2023; Contreras et al., 2003). Zero-shot prediction, an emerging technique for object detection, enables learning scheme free of repeated fine-tuning and usually output the data point as a new class or derive a description based on learning with historical data (Bansal et al., 040 2018; Li et al., 2019). In the field of time-series forecasting, zero-shot time series prediction is 041 also seriously required in data-scarce scenarios such as new cities, extreme weather conditions and 042 cross-domain adaptation. However, due to the large-scale patterns and trends in real-world series, it 043 is a challenging and under-explored task.

044 Current efforts on time series forecasting can be categorized into traditional model like moving average-based ARIMA (Contreras et al., 2003), recurrent-based LSTM (Yu et al., 2019), 046 convolution-based Temporal Convolution Network (TCN) (Hewage et al., 2020), and Transformer-047 based methods (Zerveas et al., 2021), where they can mostly well address the specific single task 048 in single domain. With the rising of Large Language Models (LLM), researchers start to exploit the open-world knowledge from LLM to help time series prediction. For the first time, Time-LLM transforms numerical series into characters and realizes prediction with LLM backbones (Zhou et al., 051 2023b). Even so, general LLM solutions solely rely on parameterized textualized knowledge and fail to filtrate specific and conditional patterns for achieving more accurate results, which is so-052 called *hallucination* in LLM. To this end, we can summarize two issues that limit existing LLMs to implement zero-shot prediction. 1) Insufficient scalability and flexibility. With increasing production of time series in open world, parameterized model usually fails to dynamically incorporate
the increasingly real-time data into predictions. And most LLM solutions without any learnable
parameters cannot adapt themselves to new target data (Zhou et al., 2023b; Gruver et al., 2024). 2) *Lacking data-model interactions*. Existing methods usually fail to provide feedback of how LLM
match the task, i.e., how to evaluate the quality of results from LLM and adjust LLM prompts for
achieving better results. Thus, interactions of coupling inherent knowledge within structural data
and textualized knowledge in LLM are significantly neglected in an unified training pipeline.

Actually, constructing multi-party in-

teractions between LLM and origi-063 nal structural data to iteratively im-064 prove the coherence between large model, small model and targeted 065 data, can activate the power of LLM 066 for achieving zero-shot time-series 067 prediction. To this end, we pro-068 pose our ZeroTS from both data as-069 pect and model aspect as follows. As illustrated in Figure. 1, from data 071 aspect, data from different domains 072 can potentially share similar evolu-073 tion patterns. The more diverse and 074 richer time-series, the more patterns



Figure 1: Motivation. (a) Series from different domains share similar patterns, where they may experience similar context (e.g., peak followed by fluctuations). (b) Compared to none feedback for optimization, involving feedback can connect the model output with arrived new time-series for model optimization, improving flexibility and adaptability.

075 are available. Fortunately, RAG integrates capacity of information retrieval with generative ability of LLM, enabling the prediction system to automatically pick up the key samples of interest from the 076 database and provide potential cues to suppress the hallucinations and deviations of the model (Fan 077 et al., 2024; Peng et al., 2024). On model aspect, most LLM solutions either directly take the output of LLM as final results (Zhou et al., 2023b; Gruver et al., 2024) thus lacking the specified optimiza-079 tion, or explore an Low-Ranking Adaptation (LoRA) to fine-tune partial layers of LLM (Hu et al., 2021; Zhou et al., 2023a) thus leading to excessive computational loads. In contrast, designing a 081 lightweight learnable module to couple the feedback of LLM and inherent data patterns can enable 082 the learnable model to fully assimilate both textualized knowledge from large model and real-world 083 data regularities, further improving the quality of prediction. 084

Even so, coupling the RAG with time series is an emerging topic that has never been reported in pre-085 vious literature. Given properties of time series, i.e., informative meta values, multi-grained pattern 086 and complex trends, the coupling is still faced with following challenges. First, How to construct 087 time series-oriented database (dataset) for supporting effective and efficient retrieval augmented 088 generation? Second, how to optimize the dual pipeline to improve quality of LLM output, i.e., 089 optimizing the representation of targeted time-series and constructing the informative and help-090 ful prompts from retrieved time series to help better implement zero-shot time series prediction? 091 Thirdly, how to construct the lightweight module, to not only assimilate the intelligence from data 092 side, but also simultaneously exploit knowledge from LLM to cooperatively realize feedback-based optimization for achieving high-quality prediction results? 093

094 To address above challenges in zero-shot time-series prediction, we propose a novel time-series learning framework ZeroTS, which leverages the data and model to cooperatively realize interac-096 tions and enables a lightweight parameter-efficient model to activate the power of LLM. The tech-097 nical contributions are three-fold. 1) We first introduce RAG into zero-shot time-series prediction, 098 devise a meta-value structure to accommodate informativeness of various series, and a modified quick retrieval HSNSW with hybrid affinity measurement for efficient retrieval. 2) We propose a 099 reinforcement scheme to assimilate knowledge from LLM, and interact LLM with ground-truth as 100 feedback to optimize the coherence between RAG and LLM, so that it can constantly absorb the 101 law of real data, and guide the output of LLM closer to the facts. 3) We evaluate our ZeroTS over 102 both zero-shot and long-term prediction settings. It achieves satisfactory performances, improving 103 at most 6% and obtaining almost best or second best results with limited memory and quickest test 104 time, i.e., 1/4 memory and 1/7 speed against other parallel models. 105

106

107

108 2 RELATED WORK

Time-series forecasting. Time Series Forecasting has evolved significantly from traditional sta-110 tistical approaches to modern machine learning techniques. Classical methods like Autoregressive 111 Integrated Moving Average Model (ARIMA) and Exponential Smoothing (ETS) provided a foun-112 dational framework by modeling linear relationships within time series (Box et al., 2015). How-113 ever, these methods often struggled with capturing non-linear dynamics in complex systems. Ma-114 chine learning algorithms such as decision trees and support vector machines later improved perfor-115 mance by addressing non-linearities (Hyndman, 2018). Deep learning techniques like Long-Short 116 Term Memory networks (LSTM) (Graves & Graves, 2012) and Temporal Convolutional Networks (TCNs) (Bai et al., 2018) has emerged, excelling in identifying complex temporal dependencies and 117 long-range patterns within sequential data. Besides, to advance the multi-variate series forecasting, 118 Graph neural network (GNN)-based series forecasting are proposed, such as MTGNN (Wu et al., 119 2020; Zhou et al., 2023b) adaptively captures variable-wise relations, and CrossGNN dynamically 120 captures scale-level and variable-level correlations (Huang et al., 2023). These models have been 121 widely adopted for various tasks such as finance, weather prediction, and traffic demand. 122

LLM4TS and RAG. LLM4TS (Large Language Models for Time Series) represents a novel ap-123 proach that leverages the capabilities of large language models for time series forecasting. Unlike 124 traditional methods, LLM4TS utilizes self-attention mechanisms to capture global patterns and de-125 pendencies in sequential data, providing enhanced performance over longer time horizons (Chang 126 et al., 2023). Recent advancements have shown that large models, such as Transformers, are highly 127 effective in understanding the temporal structure of time series, enabling them to generate more ac-128 curate forecasts (Vaswani, 2017). For example, FPT (Zhou et al., 2023a) fine-tunes a BERT encoder 129 to perform time series forecasting. Similarly, Zhang et al, (Zhang et al., 2023) introduce Meta-130 Transformer, a framework for finetuning a language model for non-text modalities, including time 131 series. Huang, et, al (Huang et al., 2024) devises a LeRet to integrate language knowledge with 132 structural data to realize accurate LLM-empowered forecasting.

Retrieval-Augmented Generation (RAG) enhances large language models by allowing them to retrieve external information dynamically. This technique combines the strengths of generative models, such as GPT, with retrieval mechanisms that automatically retrieves relevant knowledge from
external databases (Lewis et al., 2020). By incorporating external data into the forecasting process,
RAG improves the model's ability to make informed predictions, especially when facing scenarios
with limited historical data or requiring domain-specific knowledge (Karpukhin et al., 2020).

The combination of LLM4TS and RAG can offer significant advantages. LLM is equipped with great open-world knowledge in parameters for forecasters. RAG enhances LLM4TS's forecasting by incorporating relevant external knowledge, like economic indicators or weather series, leading to better contextual understanding and more precise predictions in specific domains. Designing such hybrid model can potentially tackle the data sparsity and uncertainty by filling missing information with external knowledge. But unfortunately, there is very limited works investigating such new learning scheme, i.e., LLM4TS with RAG.

146 Zero-Shot Forecasting. Zero-Shot Forecasting has emerged as an exciting area in time series fore-147 casting, offering the ability to make predictions where task-specific training data is either unavailable 148 or extremely limited (Bansal et al., 2018; Li et al., 2019). Zero-shot forecasting utilizes pre-trained 149 models that can generalize across domains, allowing them to make accurate predictions in novel 150 contexts without retraining. This capability is particularly useful in cases where collecting labeled time series data is costly or impractical. LLMTime (Gruver et al., 2024) explored the potential of 151 zero-shot forecasting by demonstrating how pre-trained language models can be applied to time se-152 ries data, generating reliable forecasts without the need for extensive fine-tuning. The combination 153 of LLM4TS and zero-shot forecasting allows for the flexible application of models to new domains, 154 greatly enhancing the practical utility of these systems in real-world forecasting tasks. However, 155 given the exciting recent research, how to incorporate the RAG bonus into LLM to enable efficient 156 and effective zero-shot forecasting still remains challenging and unclear. 157

158 159

3 PRELIMINARIES AND PROBLEM DEFINITION

161 **Time-series forecasting (TSF).** Given time series $\{x_1, x_2, ..., x_t, ...\}$, and a well-trained foundation time-series forecasting model \mathcal{M} , the goal of time series forecasting is to predict the following

162 consecutive Q steps by optimizing all learnable parameters Θ in \mathcal{M} , i.e., $\hat{Y} = \mathcal{M}(X; \Theta)$, where 163 $\hat{Y} = {\hat{y}_{T+1}, \hat{y}_{T+2}, ..., \hat{y}_{T+Q}}, X = {x_1, x_2, ..., x_T}, \Theta$ represents all learnable parameters ${W_i}$. 164 Retrieval Augmented Zero-shot TSF. Given a set of series from different domains 165 $\mathcal{G} = \{\mathcal{G}_1, \mathcal{G}_2, ..., \mathcal{G}_K, ...\},$ where each domain consists of corresponding time-series $\mathcal{G}_k \leftarrow$ 166 $\{X^k | x_1^k, x_2^k, ..., x_T^k\}$. We construct them into an auxiliary database (dataset) \mathcal{D} where the do-167 main and series can be dynamically updated and increasing. Given one targeted series X^{T} , we 168 retrieve K most helpful auxiliary series $\mathbb{X}^R = \{X_0^R, X_1^R, ..., X_K^R\} \in \mathcal{D}$, and implement forecasting $\widehat{Y}^T = \mathcal{M}^*(X^T | \mathbb{X}^R)$, where \mathcal{M} is the modified learning pipeline with LLM, $\widehat{Y}^T = \{\widehat{y}_{T+1}, \widehat{y}_{T+2}, ..., \widehat{y}_{T+Q}\}, \mathbb{X}^R = \{X_1^R, X_2^R, ..., X_K^R\}$. The condition can be considered as the re-170 171 trieved augmented series for implementing RAG prompts. 172



Figure 2: Framework overview of ZeroTS.

4 METHODOLOGY

Our ZeroTS addresses the zero-shot time series forecasting, which leverages the power of LLM 194 and coheres the gap between the parameterized knowledge in LLM and structural real-world data 195 via devising a parameter-efficient adapter. Our ZeroTS consists of two major components, a Time-196 series Retrieval Augmented Generation (TS-RAG) for efficiently retrieving high-quality exogenous 197 series from an auxiliary database (corresponding to an external dataset) \mathcal{D} , and a ReinLLM, an reinforcement scheme that simultaneously calibrates parameter-efficient adapter and forces LLM 199 to approach the data property. The framework overview is illustrated in Figure. 2. More technical 200 details, theoretical guarantee as well as the efficiency are comprehensively discussed in Sec. A.1 and 201 A.3 of Appendix.

202 203

204

173

174

175

176

177

178

179

181

183

186

187

188

189 190 191

192 193

EFFICIENT AUGMENTED TIME SERIES RETRIEVAL 4.1

205 As discussed, auxiliary series can help provide informative hints for guiding prediction of targeted 206 series. To ensure augmented prediction, the first task is to construct a retrieval dataset and efficiently 207 retrieve the augmented auxiliary time-series.

208 Structural key-value generation for RAG dataset. Time series is a kind of structural data with 209 various meta information such as domain category, timestamps, and identity of corresponding series. 210 The construction of our RAG database for auxiliary series retrieval can be two-fold, the meta infor-211 mation and deterministic value of series. To standardize the data formats and facilitate the retrieval 212 process, we first generate a structural key-value sequence in RAG database for efficient retrieval. 213 We argue that both textualized meta information and numerical statistics are important in retrieval, thus our keys consist of the meta information and especially the numerical statistics. We select the 214 domain category, timestamps and spatial location or user identification as the meta information, and 215 calculated mean and variance are also considered meta information for similarity-based retrieval.

216 For the values of numerical series, we normalize them to eliminate the influences of fluctuation 217 ranges of each series in different domain categories. Given the *i*-th time series in training set X_i , we 218 update each element into the normalized version by extracting the expectation $avg(X_i)$ and standard deviation $\operatorname{std}(X_i)$ of corresponding sequence, $x_i' = \frac{x_i - \operatorname{avg}(X_i)}{\operatorname{std}(X_i)}$. For easy notation, we let $x_i \leftarrow x_i'$ 219 220 update sequences into normalized version. 221

Data-driven series retrieval. As time-series with similarities tend to share common evolution pat-222 terns, we explore the RAG series database for augmented learning. To facilitate the retrieval, we 223 devise a compressed representation to maximumly reduce the computation loads and propose a hy-224 brid similarity to extract the optimal sequence for augmentation. Given series $X_0^R \in \mathbb{R}^{1 \times Q}$ in 225 the retrieved database (dataset), the compressed representation can be obtained by minimizing the 226 following reconstruction objective, 227

$$\min || \mathbf{X}_0^R - (\mathbf{X}_0^R * \mathbf{W}_c) * \mathbf{W}_c' ||_2$$
(1)

where the $W_c \in \mathbb{R}^{q \times p}$, $W_c' \in \mathbb{R}^{p \times q}$, $p = \rho q (0 < \rho < 0.5)$. When the model is well-trained, we adopt the compressed matrix $(X_0^R)_c = X_0^R * W_c$ as the compressed representation and transfer the 229 230 231 W_c to any other series for reducing the dimension of series.

232 Regarding the computation of similarity, we adopt a hybrid similarity metric by incorporating both 233 cosine similarity for trend modeling and Euclidean distances for value discrepancy. Given a targeted 234 sequence X^T and series from auxiliary database (dataset) X_i^R , we can obtain the similarity by,

235 236

249 250

251

253

228

$$\operatorname{Sim}(\boldsymbol{X}^{T}, \boldsymbol{X}_{i}) = \cos(\boldsymbol{X}^{T}, \boldsymbol{X}_{i}) + 1/dist(\boldsymbol{X}^{T}, \boldsymbol{X}_{i})$$
(2)

where $\cos(\mathbf{X}^T, \mathbf{X}_i)$ and $dist(\mathbf{X}^T, \mathbf{X}_i)$ denote cosine similarity and Euclidean distance. Eq. 2 237 emphasizes more on cosine similarity as we focus more on trends and evolution patterns. 238

239 For retrieval, we exploit a Hierarchical Series-level Navigable Small World (HSNSW), which 240 slightly modifies the HNSW from Faiss (Johnson et al., 2019) into series level with a series-level 241 similarity measurement Dist(a, b). The core idea of HSNSW is to organize different series into a 242 multi-layered graph, where each layer is a network of small world, with upper layer providing a 243 quick search and the lower layer providing a more fine granular search (Malkov & Yashunin, 2018; 244 Zhong, 2020). We consider the series-level affinity as the hybrid similarity in Eq. 2, and implement 245 the series clustering to construct the layer-wise network. Then the advantage of HSNSW lies in its efficient hierarchical search from coarse to fine granularity and its scalability for dynamically 246 inserting or deleting the records in dataset \mathcal{D} without re-constructing the searching indexes. Given 247 the targeted series X^T , we select the Top-K affinity series by Eq. 2, 248

$$\mathbb{X}^{R} = \{\boldsymbol{X}_{1}^{R}, \boldsymbol{X}_{2}^{R}, \dots, \boldsymbol{X}_{K}^{R}\} = \underset{\text{Min}-\text{K}}{\text{Dist}}(\boldsymbol{X}^{T}, \boldsymbol{X}_{i}^{R})$$
(3)

where the retrieved K series are further processed as representation for constructing the prompts in ZeroTS prediction. 252

4.2 **REINLLM ADAPTER** 254

255 Existing LLM-based solutions with RAG usually consider the retrieved series as prompts and di-256 rectly take them into LLMs parallelly without any correspondence and learnable mechanism, lead-257 ing to lacking of data-model interactions for model adaptation (Jin et al., 2023; Zhou et al., 2023a). 258 In this work, with the help of time series patterns from constructed series dataset, ZeroTS leverages 259 external knowledge to prompt the LLM for better prediction. We propose an RAG-LLM adapter, 260 designated as ReinLLM, a lightweight reinforcement model, which not only introduces learnable 261 parameters to cohere the LLM with retrieved and targeted series representations, but also allows data-model interactions by introspecting the output errors of LLM. For zero-shot prediction, the 262 core task can be divided into two-fold. First, how to sufficiently capture the previous evolution 263 patterns of targeted series, while second one is how to identify the underlying evolution trends and 264 patterns of targeted unseen series by maximumly exploiting the auxiliary series. With above analy-265 sis, as illustrated in Figure. 3, we model our Adapter as a reinforcement architecture equipped with 266 a small scale learnable module, instantiated with a policy network for action selection and a value 267 network for deriving the rewards. 268

Policy network for action selection. We determine our policy network by considering the impor-269 tant elements during learning as actions and update the representation with selected action as the

277

281 282

283

284

297

317

320 321

322 323

270 new states. Our policy network is composed of two branches, i.e., temporal kernel selection for 271 target series representation and series fusion coefficient selection for retrieved series representation. 272 Regarding target series modeling, we aim to output a series of best temporal convolution sizes as 273 actions. We take the available targeted time series X^T as inputs, then they are fed into a tempo-274 ral convolution layer with an initialized convolution kernel size τ , and an MLP is followed with 275 Softmax. The action derivation process can be formulated as, i.e.,

$$\boldsymbol{H} = \operatorname{Conv}_{\tau}(\boldsymbol{X}^T; \boldsymbol{W}_{cov}) \tag{4}$$

where H is the hidden representation of the target series. Assuming there are potentially l convolution kernels for selection, then an ArgSort {Softmax(·)} is imposed to derive the best action via descending order,

$$[\tau_1, \tau_2, ..., \tau_l] = \operatorname{ArgSort}\{\operatorname{Softmax}(\operatorname{MLP}(\boldsymbol{H}; \boldsymbol{W}_{ker})\}$$
(5)

The kernel associated with largest probability τ_1 can be selected to implement the convolution for obtaining series-level representation. Our end-to-end learning process can help our policy network select the optimized kernel size for series representation.



Figure 3: Illustration of ReinLLM adapter.

298 Regarding the policy network of series retrieval branch, we aim at extracting sufficient cues 299 for target series and finding out the optimal fusion coefficients for K retrieved series \mathbb{X}^R = 300 $\{X_1^R, X_2^R, ..., X_K^R\}$ with reinforcement learning. First, for pattern extraction of retrieved series, we 301 construct a triple unit including three subsections, i.e., domain category, trend pattern, text descrip-302 tion, as well as series representation, to maximumly extract sufficient cues from available series. 303 In detail, the meta information is tokenized and concatenated by textual attributes into X_m . The 304 trend pattern is extracted by element-wise differential vector $(\widetilde{X}_t)_i = \{X_i^R(t+1) - X_i^R(t)\}^{-1}$, 305 where t + 1 and t are two adjacent temporal steps. Since detailed data description of time-series 306 can provide informative knowledge for their regularities, such as 'traffics on rainy days' indicating 307 more congestions and accidents within road networks. Finally, we impose a series-level repre-308 sentation parameters $(\mathbf{X}_r)_i = \mathbf{X}_i^R * \mathbf{W}_a$, where \mathbf{W}_a is the learnable parameters for representation 309 transformation. We further concatenate the transformed $(\tilde{X}_r)_i$ with previous two subsections into 310 $\boldsymbol{X}_i = [(\boldsymbol{X}_m)_i, (\boldsymbol{X}_t)_i, (\boldsymbol{X}_r)_i].$ 311

To fuse K auxiliary series, we take advantage of the reinforcement architecture via our policy network where the aggregation weights $\{\alpha_i\}, (i \in \{1, 2, ..., K\})$ are selected actions. Specifically, we feed the concatenated representation of each series \widetilde{X}_i into our retrieval policy network with parameters W_{agg} , and output the predicted action $\{\alpha_i\}$,

$$(\alpha_1, \alpha_2, ..., \alpha_K) = \mathrm{MLP}(\mathrm{Concat}[\mathbf{X}_1, ..., \mathbf{X}_K]; \mathbf{W}_{agg})$$
(6)

The α_i will be updated in different actions in a step-by-step manner, i.e., $\alpha_i = \alpha_i - \eta$, where η is the step size. With well-learned actions $\{\alpha_i\}$, we can update the representation of retrieved ones X_R ,

$$\widetilde{\mathbf{X}}^{R} = \sum_{i=1}^{K} \alpha_{i} \widetilde{\mathbf{X}}_{i}$$
(7)

¹We omit the superscript R for the modified sequence of X.

Then \widetilde{X}^R can well indicate the patterns of series close to targets, for constructing prompts. Our end-to-end learning with error-based minimization allows the automatic optimum selection of α_i .

LLM-feedback Value Network. Finally, we construct the value network to derive the potential 327 reward for each action in the ReinLLM Adapter. It has been demonstrated that introspecting the 328 errors and differences from outside environment can make more sense than learning from the model 329 itself (Li et al., 2024). To this end, the ground-truth of zero-shot series, can be viewed as the open-330 world environment, and enable our lightweight learnable module to actively interact with the LLM 331 and outside environment (Ground-truth), thus fully activating the power of LLM. We realize it by 332 seamlessly integrating a parameterized learnable module before our Time Series-oriented LLM (TS-333 LLM). Assuming W_T and W_R are two set of parameters for respectively aligning target series and 334 retrieved ones with TS-LLM, we can respectively update their representations, {Updated targeted 335 series: \widetilde{X}^T , Retrieved augmented series for prompts: P}, with selected actions in value network,

$$\widetilde{\boldsymbol{X}}^T = \mathrm{MLP}(\mathrm{Conv}_{\tau_1}(\boldsymbol{X}^T; \boldsymbol{W}_c); \boldsymbol{W}_T)$$
(8)

$$\boldsymbol{P} = \mathrm{MLP}(\widetilde{\boldsymbol{X}}^{R}; \boldsymbol{W}_{R}) \tag{9}$$

Then we take the errors of LLM predictor, i.e., the negative of absolute differences between LLM output and ground-truth, and such negative values can be considered as the reward of our learning process. The higher the values are, the performances are better. Specifically, we exploit the action learned from Policy Network, and update the state by re-computing the series-level representation of targets and retrieved auxiliary ones. Therefore, all learnable parameters tend to be optimized to achieve best actions and competitive prediction performances.

Regarding the large language model (LLM), we adopt GPT-2 as the backbone to avoid the label legacy and heavy training burdens, where GPT-2 is developed by OpenAI in 2019, and all testing datasets in our work is published after 2019. To enable the cooperation between series and potential prompts, we take targeted representation and prompts respectively as inputting steps \widetilde{X}^T and informative prompts P. We can then obtain the following conditional output from TS-LLM,

$$\widehat{Y}^T = \text{TS} - \text{LLM}(\widetilde{X}^T | P)$$
(10)

The $\widehat{\mathbf{Y}}^T = {\widehat{y}_1^T, ..., \widehat{y}_Q^T}$ is the final prediction of our ZeroTS learning system. To calibrate the prediction results, we compare the output of LLM against ground-truth of time-series, and take such differences based on MAE back to our value network and optimize the representation thus the learnable parameters can be led to approaching the LLM.

$$Value(\widetilde{\boldsymbol{X}}^T, \boldsymbol{P}) = -MAE(\widehat{\boldsymbol{Y}}, \boldsymbol{Y})$$
(11)

The higher MAE leads to lower value and vice versa. Then our value network updates the parameters to estimate the value with learnable parameters, and impose the introspection whether the learnable small module activate the maximum power of LLM for predicting reasonable and true values approaching ground-truth. With prompts *P*, we can easily suppress the hallucination of LLMs.

5 EXPERIMENT

336 337 338

339

351 352

353 354 355

356

357 358 359

360

361

362

363 364

365 366

367

368

369 370

371

We conduct extensive experiments on various datasets under the settings of both zero-shot and longterm prediction to verify the effectiveness of ZeroTS, where long-term prediction can be considered as a scenario with potential distribution shifts.

5.1 DATASETS

Datasets for retrieval (RAG). Since time-series for real-world usually share commonality among
each other, we thus construct the database (dataset) from three large-scale datasets for retrieving
series-level regularity, supporting our zero-shot forecasting. a) UCR: It is a dataset for time-series
classification, including 128 subsets and covering the areas from medical science, and electricity,
to geography (Tan et al., 2020). b) Monash: A comprehensive dataset, consisting of 30 subsets of
time series, with totally more than 100,000 time series. It covers fields of health, retail, ride-sharing,
and demographics (Tan et al., 2020). c) TSB-UAD: TSB-UAD contains 18 subsets from real-world

575	2	-	0
	0	1	J
	_	_	

382

384 385 386

387

388

389

390

391

392

393

394 395

396

 Table 1: Performance comparison on zero-shot forecasting

	ETTh1-	→ETTh2	ETTh1	→ETTm2	ETTh2	→ETTh1	ETTh2	→ETTm2	ETTm1	→ETTh2	ETTm1	→ETTm2	ETTm2	→ETTh2	ETTm2	2→ETTm1
Metrics	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE
FimesNet	0.431	0.421	0.361	0.327	0.621	0.865	0.376	0.342	0.454	0.457	0.354	0.332	0.443	0.435	0.567	0.769
Time-LLM	0.387	0.353	<u>0.340</u>	0.273	0.474	0.479	0.341	0.272	0.412	0.381	0.320	0.268	0.400	0.354	0.438	0.414
GPT4TS	0.422	0.406	0.363	0.325	0.578	0.757	0.370	0.335	0.439	0.433	0.348	0.313	0.443	0.435	0.567	0.769
LLMTime	0.708	0.992	0.869	1.867	0.981	1.961	0.869	1.867	0.768	0.992	0.869	1.867	0.708	0.992	0.984	1.933
DLinear	0.488	0.493	0.452	0.415	0.574	0.703	0.386	0.328	0.475	0.464	0.389	0.335	0.471	0.455	0.537	0.649
PatchTST	0.405	0.380	0.360	0.314	0.513	0.565	0.365	0.325	0.438	0.439	0.296	0.334	0.409	0.425	0.568	0.492
ZeroTS	0.379	0.350	0.327	0.272	0.522	0.575	0.335	0.267	0.483	0.395	0.313	0.312	0.374	0.351	0.442	0.416

data science applications, which totally accounts for 12,686 time series with labeled anomalies spanning different domains with high variability of anomaly types, ratios, and sizes (Paparrizos et al., 2022b;a).

Datasets as targeted time-series for forecasting. We exploit ETTh1, ETTh2, ETTm1, ETTm2 as datasets for zero-shot predictions (Jin et al., 2023). Regarding long-term prediction, we take totally 8 datasets for evaluation, where we further employ additional four datasets those are benchmarking long-hrozion prediction models, i.e., Weather, ECL, Traffic, ILI (Wu et al., 2023). Details of datasets can be found in Appendix. A.4.

5.2 BASELINES

397 **Zero-shot forecasting.** We exploit 6 strong baselines for evaluating the zero-shot prediction setting. 398 1) TimesNet. It transforms the one-dimension time-series into a two-dimension series and employ 399 2D-convolution for temporal pattern extraction (Wu et al., 2022). 2) Time-LLM. A state-of-the-art 400 series learning model that reprogrammes large language model (LLM) to adapt for time series fore-401 casting (Jin et al., 2023). 3) GPT4TS. It is a one-fits-all model, that exploits the frozen parameters 402 from pretraining in computer vision and NLP to adapt time series learning (Zhou et al., 2023a). 4) LLMTime. An LLM-based time series forecaster without inputting retrieved series as auxiliary 403 information (Gruver et al., 2024). 5) DLinear. A linear model for long-term series forecasting 404 with a rolling prediction strategy. 6) PatchTST. This baseline simultaneously explores the patch 405 division and channel independence for multi-variate series learning (Nie et al., 2022). Among them, 406 Time-LLM (Jin et al., 2023) achieves the best competitive performances under zero-shot predictions. 407

Long-horizon forecasting. Given that LLMTime is not utilized to perform long-term prediction, we
 additionally take AutoFormer (Wu et al., 2021) and Informer (Zhou et al., 2021) as baselines for
 evaluating long-horizon prediction. These two baselines are inherited from Transformer (Vaswani,
 with respectively considering auto-correlation for long-term series forecasting and probabilis tic sparse self-attention.

413 414

415

5.3 IMPLEMENTATION DETAILS

421 Setup for long-term forecasting. The long-term prediction set as training and testing on the same 422 dataset with augmented series. We adopt 8 datasets for an extensive evaluation. The input length 423 of time series is 512, and also output the horizons of respective {96, 192, 336, 720} steps for all 424 datasets except for ILI. Since ILI is a dataset with a smaller scale, we set the forecasting horizons 425 $Q \in \{24, 36, 48, 60\}$ for it.

426 427

5.4 PERFORMANCE COMPARISON

The prediction results for zero-shot and long-term settings are respectively shown in Table. 1 and 2. The best results are **bold** and the second best are <u>underlined</u> for comparison.

Zero-shot prediction. As shown, our ZeroTS achieves either best or second best performances on all 8 settings, and the improvement is ranging from 2.07% to 6.5%. Compared with well-known

400	Δ	3	2
400			-
	л	0	0

T 1 1 A	DC	•	1 /	c
Table 7.	Performance	comparison	on long-fei	m forecasting
1 a 0 1 c 2.	1 chroninance	companson	on long ter	In rorecasting

							-			-			U			
	Time	sNet	Time-	LLM	GPT	TATS	DLi	near	Patch	nTST	AutoF	ormer	Info	mer	ZeroT	S(Ours)
Metrics	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE
ETTh1	0.450	0.458	0.423	0.408	0.455	0.465	0.437	0.422	0.430	0.413	0.487	0.496	0.795	1.040	0.421	0.403
ETTh2	0.427	0.414	0.383	0.334	0.412	0.381	0.446	0.431	0.379	0.330	0.459	0.450	1.729	4.431	0.375	0.315
ETTm1	0.406	0.400	0.372	0.329	0.403	0.388	0.378	0.357	0.380	0.351	0.517	0.588	0.734	0.961	0.377	0.347
ETTm2	0.333	0.291	0.313	0.251	0.339	0.284	0.333	0.267	0.315	<u>0.255</u>	0.371	0.327	0.810	1.410	0.331	0.276
Weather	0.287	0.259	<u>0.257</u>	0.225	0.270	0.237	0.300	0.248	0.264	0.225	0.382	0.338	0.548	0.634	0.250	0.225
ECL	0.295	0.192	0.252	0.158	0.263	0.167	0.263	0.166	0.252	0.161	0.338	0.227	0.397	0.311	0.260	0.160
Traffic	0.336	0.620	<u>0.264</u>	0.388	0.294	0.414	0.295	0.433	0.263	0.390	0.379	0.628	0.416	0.764	0.261	0.384
ILI	0.931	2.139	0.801	1.435	0.903	1.925	1.041	2.169	0.797	1.443	1.161	3.006	1.544	5.137	0.795	1.436

439 440

441 442

443

444

445

LLM-based Time-LLM, GPT4TS, PatchTST and LLMTime, ZeroTS achieves comparative and even slightly better performances than them without any fine-tune on LLM. Along with complexity and efficiency comparison in Table. 5, we believe ZeroTS is a superior lightweight zero-shot prediction model, benefiting from retrieved augmented representation, and our actively feedback scheme.

Long-term prediction. In Table. 2, our ZeroTS surpasses most other solutions with Retrieved
 Augmented Series, achieving 5 out 8 best and 2 second best. It is worth noting that ZeroTS is with
 fairly fewer learnable parameters to achieve such satisfactory results. And it is also interesting to
 find the improvement over long-term setting is less significant than zero-shot ones. It is reasonable
 that ZeroTS is more effective on data scarcity scenario, where retrieved ones provide addition hints
 to constrain and guide LLM generation.

452 453 454

5.5 ABLATION STUDY

To investigate the superiority of each well-designed component, we devise various ablative variants via replacing the component with vanilla one or directing removing corresponding one.

1) ZeroTS-RAG. We remove the auxiliary sequence retrieval and exploit the LLM for zero-shot forecasting to confirm the motivation of RAG exploitation. This scheme potentially degenerates to the LLMTime (Gruver et al., 2024). 2) ZeroTS-Para. We remove the learnable parameters in our adapter and only utilize the concatenation of retrieved series as prompts for TS-LLM. This variant also degenerate it into a model without feedback from ground-truth. 3) ZeroTS-Rein. To verify the effectiveness of coupling reinforcement learning and LLM, we degenerate the policy network into an end-to-end learning scheme with immediate error back-propagation.

464 Main results. We perform our ablation study on two 465 zero-shot settings in Table. 3 to verify the rationales of each design in ZeroTS. Most significantly, by re-466 moving the retrieval augmented series, our ZeroTS-467 RAG degenerates to an LLM-based model without 468 any augmentation except for inputting targeted se-469 ries. The performance experiences the largest drop, 470 ranging from 9.63% to 14.76% on MAE and MSE, 471 it suggests that retrieved series explicitly provide ad-472 ditional signals and regularity to help prediction, es-

Table 3:	Ablation	studies	on	zero-shot	set-
tings					

	ETTh1 -	\rightarrow ETTh2	ETTh1 -	→ ETTm2
Metrics	MAE	MSE	MAE	MSE
ZeroTS-RAG	0.435	0.398	0.364	0.301
ZeroTS-Para.	0.407	0.375	0.349	0.282
ZeroTS-Rein	0.412	0.380	0.351	0.286
ZeroTS	0.379	0.353	0.332	0.274

473 pecially on zero-shot predictions. With learnable parameters removed, our ZeroTS cannot receive 474 the feedback from LLM and cannot be trained in an end-to-end manner, leading to lacking flexibility 475 of correspondence between learnable module, LLM, as well as learnable module and retrieved data. 476 Then our ablative Zero-Para has 2%-6% performance degeneration, which verifies the importance of 477 imposing feedback from ground-truth for optimizing learnable module. In addition, Similarly, when modifying the reinforcement architecture into an ordinary end-to-end architecture, our learning mod-478 ule cannot find a more optimized action thus becomes inferior to integrated ZeroTS, demonstrating 479 the rational promotion of reinforcement scheme. 480

481 482

483

5.6 Hyperparameter study and detailed analysis

We provide three important hyperparameters for observing how can we obtain their best performances. 1) The number of retrieved time series from auxiliary series dataset $K = \{3, 6, 9, 12\}$. 2) Representation dimension of targeted time series ranging from $\{8, 16, 24, 32, 64\}$. 3) Representation dimension of retrieved time series (meta information and time series observations) ranging from {8, 16, 24, 32, 64}. Due to the space limitation, we only illustrate the adjustment process on zeroshot setting over ETTh1 \rightarrow ETTm1 on Figure. 4. As observed, K = 6 achieves best results, as more retrieved series will tend to involve more noise while fewer series will become less informative. For dimension of representations, we can see larger dimensions can result in better performances but we have to make the trade-off between performances and efficiency. We thus choose K = 6, dimensions for target and retrieved ones are set as 64.

494 5.7 CASE STUDY

493

495

514 515

516

517

518

519

521

522

523 524

529

To illustrate how retrieved se-496 ries help better prediction, we 497 provide intuitive analysis of in-498 termediate results during predic-499 tion process. Given a targeted 500 series on ETTh2, as shown in 501 Figure. 5, our TS-RAG in Ze-502 roTS implements a retrieval and to select Top-6 series with most 504 proximity where three series are 505 illustrated. The retrieved series



Figure 4: Hyperparameter analysis on ETTh1→ETTm1.

exactly share similar regularity and evolution patterns through the hybrid metric of proximity mea-506 surement and the similarities with targets are marked. To enhance forecasting, the auxiliary series 507 are aggregated into a representation and we tokenize the representation as prompts into LLM, along 508 with main input of target series. We visualize the LLM output with and without (w/o) the series for 509 an intuitive comparison. Specifically, prediction with RAG reveals the averaged pattern of retrieved 510 series thus overcoming the overfitting on targeted series and become smoother, while prediction w/o 511 RAG is with more fluctuations and many details against facts. This also suggests RAG can suppress 512 the hallucination of LLM on time series learning. Our analysis provides further interpretability and 513 understanding of ZeroTS and highlight the contribution of RAG enhanced zero-shot prediction.



Figure 5: Case studies on ETTh2 prediction.

6 CONCLUSION AND FUTURE WORK

530 In this work, we propose a novel time-series learning framework ZeroTS, which enables multi-party 531 model-data interactions, including retrieved series, targeted series, LLM as well as learnable adapter 532 in an end-to-end learning manner. Specifically, we provide a series dataset construction protocol, 533 a small-scale learnable adapter for cohering RAG and LLM, and the first reinforcement architecture 534 for zero-shot time series prediction considering environment feedback. Extensive evaluation results suggest the effectiveness and superiority of our ZeroTS under zero-shot setting as well as complex 536 long-term forecasting. We believe ZeroTS is an interesting work that provides new insights into series learning under extreme scenarios such as data scarcity, long-horizion and zero-shot predictions with small learnable parameters for modeling interactions. Future works can respectively lie in 538 generalizing our solution to more complex scenarios as well as efficient retrieved series towards better prediction. Detailed discussion can be found in Appendix. A.5.

540 REFERENCES

554

560

561

562

- Ali N Akansu and Richard A Haddad. *Multiresolution signal decomposition: transforms, subbands, and wavelets.* Academic press, 2001.
- Ali Naci Akansu and Yipeng Liu. On-signal decomposition techniques. *Optical Engineering*, 30 (7):912–920, 1991.
- Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen,
 Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, et al.
 Chronos: Learning the language of time series. *Machine learning*, 1(2):3.
- ⁵⁵⁰ Robert B Ash. *Information theory*. Courier Corporation, 2012.
- Anonymous ICLR25 Authors. Unitst: Effectively modeling inter-series and intra-series dependencies for multivariate time series forecasting.
- Shaojie Bai, J Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- Ankan Bansal, Karan Sikka, Gaurav Sharma, Rama Chellappa, and Ajay Divakaran. Zero-shot
 object detection. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 384–400, 2018.
 - George EP Box, Gwilym M Jenkins, Gregory C Reinsel, and Greta M Ljung. *Time series analysis: forecasting and control.* John Wiley & Sons, 2015.
- 563 Ching Chang, Wen-Chih Peng, and Tien-Fu Chen. Llm4ts: Two-stage fine-tuning for time-series
 564 forecasting with pre-trained llms. *arXiv preprint arXiv:2308.08469*, 2023.
- Javier Contreras, Rosario Espinola, Francisco J Nogales, and Antonio J Conejo. Arima models to predict next-day electricity prices. *IEEE transactions on power systems*, 18(3):1014–1020, 2003.
- Timothy DelSole and Michael K Tippett. Predictability: Recent insights from information theory. *Reviews of Geophysics*, 45(4), 2007.
- Xin Ding, Lu Chen, Yunjun Gao, Christian S Jensen, and Hujun Bao. Ultraman: A unified platform
 for big trajectory data management and analytics. *Proceedings of the VLDB Endowment*, 11(7):
 787–799, 2018.
- Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 6491–6501, 2024.
- Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski.
 Moment: A family of open time-series foundation models. In *Forty-first International Conference* on Machine Learning.
- Alex Graves and Alex Graves. Long short-term memory. Supervised sequence labelling with recurrent neural networks, pp. 37–45, 2012.
- Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew G Wilson. Large language models are zero-shot time series forecasters. *Advances in Neural Information Processing Systems*, 36, 2024.
- Pradeep Hewage, Ardhendu Behera, Marcello Trovati, Ella Pereira, Morteza Ghahremani,
 Francesco Palmieri, and Yonghuai Liu. Temporal convolutional neural (tcn) network for an effective weather forecasting using time-series data from the local weather station. *Soft Computing*, 24:16453–16482, 2020.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint* arXiv:2106.09685, 2021.

- 594 Qihe Huang, Lei Shen, Ruixin Zhang, Shouhong Ding, Binwu Wang, Zhengyang Zhou, and Yang 595 Wang. Crossgnn: Confronting noisy multivariate time series via cross interaction refinement. 596 Advances in Neural Information Processing Systems, 36:46885–46902, 2023. 597
- Oihe Huang, Zhengyang Zhou, Kuo Yang, Gengyu Lin, Zhongchao Yi, and Yang Wang. Leret: 598 Language-empowered retentive network for time series forecasting. In Proceedings of the International joint conference on artificial intelligence, 2024. 600
- 601 RJ Hyndman. Forecasting: principles and practice. OTexts, 2018.

619

623

630

635

- Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y Zhang, Xiaoming Shi, Pin-Yu Chen, Yux-603 uan Liang, Yuan-Fang Li, Shirui Pan, et al. Time-Ilm: Time series forecasting by reprogramming 604 large language models. arXiv preprint arXiv:2310.01728, 2023. 605
- 606 Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. IEEE 607 *Transactions on Big Data*, 7(3):535–547, 2019.
- 608 Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi 609 Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. arXiv 610 preprint arXiv:2004.04906, 2020. 611
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, 612 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-613 tion for knowledge-intensive nlp tasks. Advances in Neural Information Processing Systems, 33: 614 9459-9474, 2020. 615
- 616 Xiaoyuan Li, Wenjie Wang, Moxin Li, Junrong Guo, Yang Zhang, and Fuli Feng. Evaluating mathe-617 matical reasoning of large language models: A focus on error identification and correction. arXiv 618 preprint arXiv:2406.00755, 2024.
- Zhihui Li, Lina Yao, Xiaoqin Zhang, Xianzhi Wang, Salil Kanhere, and Huaxiang Zhang. Zero-shot 620 object detection with textual descriptions. In Proceedings of the AAAI Conference on Artificial 621 Intelligence, volume 33, pp. 8690-8697, 2019. 622
- Yong Liu, Haoran Zhang, Chenyu Li, Xiangdong Huang, Jianmin Wang, and Mingsheng Long. 624 Timer: Transformers for time series analysis at scale. arXiv preprint arXiv:2402.02368, 2024.
- 625 Yu A Malkov and Dmitry A Yashunin. Efficient and robust approximate nearest neighbor search 626 using hierarchical navigable small world graphs. IEEE transactions on pattern analysis and 627 machine intelligence, 42(4):824-836, 2018. 628
- Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 629 words: Long-term forecasting with transformers. arXiv preprint arXiv:2211.14730, 2022.
- 631 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong 632 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-633 low instructions with human feedback. Advances in neural information processing systems, 35: 634 27730-27744, 2022.
- John Paparrizos, Paul Boniol, Themis Palpanas, Ruey S Tsay, Aaron Elmore, and Michael J 636 Franklin. Volume Under the Surface: A New Accuracy Evaluation Measure for Time-Series 637 Anomaly Detection. Proceedings of the VLDB Endowment, 15(11):2774–2787, 2022a. 638
- 639 John Paparrizos, Yuhao Kang, Paul Boniol, Ruey S Tsay, Themis Palpanas, and Michael J Franklin. 640 Tsb-uad: an end-to-end benchmark suite for univariate time-series anomaly detection. Proceedings of the VLDB Endowment, 15(8):1697–1711, 2022b. 641
- 642 Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and 643 Siliang Tang. Graph retrieval-augmented generation: A survey. arXiv preprint arXiv:2408.08921, 644 2024.645
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea 646 Finn. Direct preference optimization: Your language model is secretly a reward model. Advances 647 in Neural Information Processing Systems, 36, 2024.

648 649 650	Kashif Rasul, Arjun Ashok, Andrew Robert Williams, Arian Khorasani, George Adamopoulos, Rishika Bhagwatkar, Marin Biloš, Hena Ghonia, Nadhir Hassen, Anderson Schneider, et al. Lag- llama: Towards foundation models for time series forecasting. In <i>R0-FoMo: Robustness of Few</i> -
651	shot and Zero-shot Learning in Large Foundation Models.
652	G Vijendar Reddy, Lakshmi Jaswitha Aitha, Ch Poojitha, A Naga Shreya, D Krithika Reddy, and
654	G Sai Meghana. Electricity consumption prediction using machine learning. In E3S Web of Conferences, volume 391, pp. 01048, EDP Sciences, 2023
655	Conjerences, volume 591, pp. 01048. EDF Sciences, 2025.
656 657	Ritika Singh and Shashi Srivastava. Stock prediction using deep learning. <i>Multimedia Tools and Applications</i> , 76:18569–18584, 2017.
658 659 660	Chang Wei Tan, Christoph Bergmeir, Francois Petitjean, and Geoffrey I Webb. Monash university, uea, ucr time series extrinsic regression archive. <i>arXiv preprint arXiv:2006.10996</i> , 2020.
661 662	Taosdata. Td-engine, 2020. URL https://github.com/taosdata/TDengine.
663 664	A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.
665 666 667	Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition trans- formers with auto-correlation for long-term series forecasting. <i>Advances in neural information</i> <i>processing systems</i> , 34:22419–22430, 2021.
668 669 670	Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. <i>arXiv preprint arXiv:2210.02186</i> , 2022.
671 672 673 674	Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In <i>The Eleventh International</i> <i>Conference on Learning Representations (ICLR)</i> , 2023.
675 676 677 678 679	Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, Xiaojun Chang, and Chengqi Zhang. Con- necting the dots: Multivariate time series forecasting with graph neural networks. In <i>Proceedings</i> <i>of the 26th ACM SIGKDD international conference on knowledge discovery & data mining</i> , pp. 753–763, 2020.
680 681 682	Zhongchao Yi, Zhengyang Zhou, Qihe Huang, Yanjiang Chen, Liheng Yu, Xu Wang, and Yang Wang. Get rid of isolation: A continuous multi-task spatio-temporal learning framework. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> , 2024.
683 684 685	Yong Yu, Xiaosheng Si, Changhua Hu, and Jianxun Zhang. A review of recurrent neural networks: Lstm cells and network architectures. <i>Neural computation</i> , 31(7):1235–1270, 2019.
686 687 688	Yuan Yuan, Jingtao Ding, Jie Feng, Depeng Jin, and Yong Li. Unist: a prompt-empowered universal model for urban spatio-temporal prediction. In <i>Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , pp. 4095–4106, 2024.
689 690 691 692	George Zerveas, Srideepika Jayaraman, Dhaval Patel, Anuradha Bhamidipaty, and Carsten Eickhoff. A transformer-based framework for multivariate time series representation learning. In <i>Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining</i> , pp. 2114–2124, 2021.
693 694 695 696	Yiyuan Zhang, Kaixiong Gong, Kaipeng Zhang, Hongsheng Li, Yu Qiao, Wanli Ouyang, and Xi- angyu Yue. Meta-transformer: A unified framework for multimodal learning. <i>arXiv preprint</i> <i>arXiv:2307.10802</i> , 2023.
697 698 699	Yu Zhong. Efficient Implementation of Hierarchical Navigable Small World Similarity Matching Algorithm. University of California, San Diego, 2020.
700	Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.

Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.
 Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings* of the AAAI conference on artificial intelligence, volume 35, pp. 11106–11115, 2021.

Tian Zhou, Peisong Niu, Liang Sun, Rong Jin, et al. One fits all: Power general time series analysis by pretrained lm. Advances in neural information processing systems, 36:43322–43355, 2023a. Zheng yang Zhou, Hao Liu, Kun Wang, Peng kun Wang, Xu Wang, and Yang Wang. A teacher-student spatiotemporal semi-supervised method for urban event forecasting. Acta Electronica Sinica, 51(12):3557-3571, 2023b.

756 A APPENDIX

In this section, we provide more analysis regarding specific techniques, theoretical details, experimental results and efficiency comparison, to further support the superiority of our ZeroTS and help authors and reviewers better understand the operation mechanism of our idea. Finally, a brief discussion on our future work is elaborated. The codes of our ZeroTS are available at https://anonymous.4open.science/r/ZeroTS-7F7E/.

763 764

765

A.1 TECHNICAL DETAILS

In this subsection, we will provide substantial technical details for better understanding the utilized
 techniques of ZeroTS.

768 769

770

A.1.1 THE ARCHITECTURE OF ZEROTS

In this subsection, we provide a more detailed discussion on the architecture of ZeroTS.

771 Actually, since our ZeroTS aims to realize the data-model interactions, enabling our model to extract 772 external series and how to fuse the external series with a learning objective is a natural intuition in 773 our task. To this end, we first propose TS-RAG to realize an efficient retrieval from external series 774 dataset. To fuse retrieved series and maximumly activate the power of LLM with RAG in time 775 series, we devise ReinLLM adapter, it receives retrieved series as input and interconnects with LLM 776 to derive the results, then the errors backpropagate to learnable representations on both targets and 777 prompts, thus guiding the fusion of retrieved time series and activating the maximal power of LLM. 778 The overview and data flow illustration are in Fig. 2 and Alg. 1, where TS-RAG and ReinLLM can 779 jointly implement the data-model interactions from data and model sides, activating the output of LLM towards the ground-truth. 780

⁷⁸¹ Specifically, given a targeted time-series, we activate the retrieval process with an efficient heuristic ⁷⁸² (ρ - ρ -cut-off) and hierarchical (Hierarchical Series-level Navigable Small World, HSNSW) strategy ⁷⁸³ to obtain a series of auxiliary series for prompt construction.

784 **Regarding ReinLLM**, our **policy network** consists of two major components, i.e., temporal kernel 785 selection for target series representation and auxiliary series fusion coefficient selection for time-786 series fusion. It takes the role of selecting the critical hyperparameters for targeted time-series 787 representation and auxiliary time series fusion. Update process. When the action (important hy-788 perparameter) selected, the representation will be updated with new hyperparameters, and then feed 789 the updated representation into LLM, where targeted representation as main input to LLM and fused 790 auxiliary series as prompt to LLM for supplementary cues. After that, TS-LLM, instantiated by GPT-2, outputs the prediction results and deliver the outcome and errors to the value network. Our 791 value network transfers the errors into reward and backpropagate to the policy network for strat-792 egy optimization. Actually, our lightweight learnable module is optimized from both data side and 793 model side. On data side, we exploit retrieval process from auxiliary real-world data, and let up-794 datable real data to continuously supplement the cues for model prediction. And on model side, 795 the model (LLM) is considered as the environment and we let our lightweight learnable module 796 ReinLLM interact with LLM and derive the error feedback as reward to optimize the whole process, 797 contributing to the real-world corrections and enabling seamless data-model-world interactions.

798 799

800

A.1.2 STRUCTURE OF RAG DATASET AND HOW IT UTILIZE FOR FORECASTING

Here, we provide a table to illustrate the detailed description of this key-value based structural retrieval dataset for large-scale time series, as shown in Table. 4.

Utilization of meta information for retrieval. The key-value in Table.4 can be considered as the meta information. In this subsection, we provide details on how we utilize Regarding the retrieval process, we utilize the meta information including domain category (e.g., traffics/electricity/weather...), timestamps (2024/11/11/10:00) and spatial location/user identification to help retrieve related auxiliary series. Of course, calculated mean and variance are also considered meta information for similarity-based retrieval. Specifically, given a targeted series with its domain category, timestamps, and identification (user/location id.) where it is available when data is collected, we can retrieve time series from auxiliary series base by indexing corresponding meta information (keywords), if the meta information is matched, e.g., two series share the same domain category, share similar timestamps (day of week, hour of day), this series will be selected. And we can compute the similarity between selected series and targeted series, where the top-K similar series will be utilized for following fusion and representation as prompts.

814 Meta information utilization in ReinLLM. In ReinLLM, meta information is considered as the 815 features of retrieved series, which consists of domain category, timestamps and data description, and 816 we take the concatenation of meta information, trend vector as well as deterministic series values as 817 the integrated representations. The abovementioned meta information (domain category, timestamps 818 and data description) is available when the auxiliary base is constructed, and it corresponds to the 819 dataset and series-level property. Then the words are transferred into vectors with the tokenizer in 820 LLM (i.e., GPT-2 in our ZeroTS), where the tokenizer is trained by the strategy named Byte-Pair Encoding (BPE), i.e., the tokenizer trained with occurrence frequency of words and their semantics. 821

822 823

A.1.3 COMPLEXITY ANALYSIS OF RETRIEVAL

Assume there are *M* times of computation for retrieving the matched Top-*K* series. To reduce the complexity to a reasonable range, we propose a $\rho - \rho$ -cut-off over sampling and dimensions, thus the full computational loads are reduced to $\rho^2 M$. Regarding retrieved series, we every time sample ρ % of overall all samples for similarity comparison with a stratified sampling strategy, thus the pair-wise similarity computations are reduced. The stratified strategy ensures such sampling with fairness and diversity. For representation dimension, we further reduce the dimension of representation to ρ % of original vectors, then the computing loads are further significantly decreased.

831 832

833

834 835

836

837

838

839

840

841

842

843

844

845

846

850

851

A.1.4 CLARIFICATION ON POTENTIAL INFORMATION LEAKAGE

We have taken the issue of information leakage into consideration with the following two designs.

- First, to avoid the information leakage of LLM, and fairly investigate the effectiveness of our idea, we exploit GPT-2 (released in 2019) instead utilizing a most recent model as backbone, and the testing/inference datasets are all released after 2019. That's to say, the LLM has not seen the training/testing set when experiencing the pre-training.
- Second, to avoid the information leakage of RAG, we utilize the disjoint datasets in retrieval ones and inference ones and ensure the ordered timestamps in the retrieval. More specifically, in our testing stage, to ensure fair testing and imitate the real-status, we filter the available series with meta information ('timestamps') and directly mask the 'future observations' corresponding to this targeted series as unseen ones (the remaining $1 - \rho$). To this end, we can ensure our solution undergo a fair and trustworthy test. And in deployed real-world applications, if we are expected to predict future steps, actually only previous events and observations are available as no one can reach the future beyond the present.

With above well-designed strategies, we can declare that our solution will not encounter the information leakage in both real applications and testing periods.

A.1.5 **DISCUSSION ON SELECTION OF AUXILIARY DATASET**

As discussed in our main text, our model ZeroTS does not require intensive domain-related data but 852 exploit the intuition in our Fig.1 that series share similar trends in the world can reinforce each other. 853 To this end, we further provide some discussions on how can we address challenges when similar 854 series with targets are limited. As we know, from a classical signal decomposition theory Akansu & 855 Liu (1991); Akansu & Haddad (2001), time-series tends to be decomposed of different components, 856 or a complex series can be the combinations of several simple components. Thus, we can take this theory to interpret our auxiliary series based solution. Our retrieval and fusion take a Top-K selection 858 mechanism and then dynamically learn the fusion coefficients with reinforcement scheme. Even 859 though series have limited direct similarity, their combinations may make sense. In our solution, we can always find the Top-K series and dynamically combine them in a learnable manner and further provide the cues of potential trends for output of the backbone language model. Regarding data 861 selection, there must be unavoidable performance variations on different collection of dataset. Our 862 intuition is that when the scale and diversity of external series set is increasing, the more patterns are 863 included, but we cannot numerate all time-series in the world and finding additional series will also

Table 4: Structural key-value instance

		Meta i	information			Deterministic observation
Key	Domain	Timestamps	Spatial/User ID	Mean	Variance	Series value
Value	Electricity	20140101	(119.52E, 31.05N)	0.35	0.42	[0.32,0.34,0.45,]

lead to further unexpected complexity and burden on collection. The perspective of our solution is to alleviate the influences of variations on external dataset selection and mine the potential relationships between retrieval series and target ones with trade-off. Moreover, for reducing complexity and selecting more related series, clustering external series and using prototype series for pre-selection can be investigated as future works.

875 876 877

864

866 867 868

870

871

872

873

874

Algorithm 1 The training process of ZeroTS

Input: Target series X^T , world-wide series external dataset \mathcal{D} . 878 **Output:** Retrieved Top-K series $\mathbb{X}^R = \{X_0^R, X_1^R, ..., X_K^R\}$, Predicted target series $\widehat{Y}^T =$ 879 $\mathcal{M}^*(\boldsymbol{X}^T | \mathbb{X}^R).$ 880 Initial: the TS-LLM model, learnable ReinLLM adapter, the number of total epochs M. 881 for i = 1 to M do 882 Retrieve time-series from databse $\mathbb{X}^R = \mathrm{TS} - \mathrm{RAG}(\boldsymbol{x}_i), \mathbb{X}^R \in \mathcal{D}$ 883 Run Policy Network in ReinLLM $[\tau_1, ..., \tau_l] = \operatorname{ArgSort} \{\operatorname{Softmax}(\operatorname{MLP}(\boldsymbol{H}; \boldsymbol{W}_{ker}))\}$ 885 $(\alpha_1, ..., \alpha_K) = \mathrm{MLP}(\mathrm{Concat}[\widetilde{X}_1, ..., \widetilde{X}_K]; W_{aqq})$ Run Value Network in ReinLLM $\boldsymbol{X}^{T} = \mathrm{MLP}(\mathrm{Conv}_{\tau_{1}}(\boldsymbol{X}^{T}; \boldsymbol{W}_{c}); \boldsymbol{W}_{T}); \boldsymbol{P} = \mathrm{MLP}(\widetilde{\boldsymbol{X}}^{R}; \boldsymbol{W}_{R})$ 887 888 Feed into TS-LLM, $\widehat{Y}^T = \text{TS} - \text{LLM}(\widetilde{X}^T | P)$ 889 Compute the reward Value (\tilde{X}^T, P) 890 **Optimization** 891 end for 892 **Return** \mathbb{X}^R , Predicted \mathbf{Y}^T

893 894 895

912 913

A.1.6 ALGORITHM OF HSNSW

896 Efficient series retrieval is vital for RAG as well as our zero-shot learning. Here, we exploit the 897 Hierarchical Series-level Navigable Small World (HSNSW) algorithm to help implement our retrieval, where we slightly modify HNSW into a series-level retrieval with a hybrid similarity metric 899 of Eq. 2. Here we briefly introduce the architecture of HNSW. For HNSW, its core idea is to 900 organize different samples into a multi-layered graph, where each layer is a network of small 901 world, with upper layer providing a quick search and the lower layer providing a finer granu-902 lar search (Malkov & Yashunin, 2018; Zhong, 2020). We then let it modify into a series-level retrieval process, name this mechanism as Hierarchical Series-level Navigable Small World 903 (HSNSW). Actually, we implement it by the Faiss which is an open-sourced package devel-904 oped by Facebook (Johnson et al., 2019). 905

Specifically, it involves several key parameters and probability calculations in its mathematical for mulation and construction process. This section outlines some of the core mathematical concepts
 and formulas that underpin the HNSW algorithm. We illustrate it in Figure. 8.

Layer Probability Calculation. During the construction process, the insertion probability of each vector into different layers is given by the following probability function,

$$\operatorname{Prob}(\operatorname{level} l) = \exp\left(-\frac{l}{m_L}\right) \times \left(1 - \exp\left(-\frac{1}{m_L}\right)\right) \tag{12}$$

where m_L is the level multiplier used to balance the distribution of vectors across different layers. $m_L = 0$ indicates that all vectors are inserted only at layer 0.

917 Optimal Level Multiplier. To minimize the overlap of shared neighbors across layers and balance the average number of traversals during the search process, the rule of thumb for the optimal level

918 multiplier is:

$$m_L = \frac{1}{\ln(M)},$$

921 where M is the number of neighbors each vertex has.

923 **Neighbor Selection.** At each layer, the algorithm greedily traverses edges to find the ef nearest 924 neighbors of the inserted vector q, with ef initially set to 1. In the second phase of construction, ef925 is increased to $ef_{\text{Construction}}$, thus returning more nearest neighbors as candidate points for links.

Link Addition. When adding links, two parameters M_{max} and M_{max0} are considered, which define the maximum number of links a vertex can have at non-zero layers and at layer 0, respectively.

Stopping Condition. The stopping condition for the insertion operation is to find a local minimum at layer 0, i.e., no closer neighbors are found at this layer.

930 931 932

933

934

928

929

920

A.2 DISTINGUISH ZEROTS AGAINST OTHER RELATED WORKS

As ZeroTS is at the intersection of reinforcement learning, utilization of LLM and time-series learning, we then provide a systematical discussion regarding RLHF, LLM choice and why not fine-tune the LLM to enhance and verify our design motivations.

935 936 937

938

A.2.1 DISCUSSIONS BETWEEN ZEROTS AND RLHF

Actually, both ZeroTS and RLHF (Reinforcement Learning from Human Feedback) utilize rein-939 forcement scheme to enhance the model performance. However, they are different on which side 940 of information to receive and how to implement further optimization. Specifically, ZeroTS receives 941 the feedback from LLM and optimizes the representation of fused retrieved time series, where the 942 feedback can be automatically computed via the differences between ground-truth and LLM output, 943 and it enjoys the nice property of delivering the model side feedback to previous data side. In con-944 trast, RLHF in InstructGPT explores a fine-tune process (PPO) to align the awareness between GPT 945 and human annotated labels (Ouyang et al., 2022). DPO exploits a scheme to remove the unstable 946 reinforcement process and direct optimizes the reward from human label rewards (Rafailov et al., 947 2024). To this end, RLHF requires intensive data and extensive human labors for LLM fine-tune, 948 while our ZeroTS optimizes the retrieved series (both representation and fusion mechanisms) fusion with gaining awareness from LLM output side, which is more efficient, automatic, and lightweight 949 with fewer human workloads, yielding more practicality. 950

951 952

A.2.2 DISCUSSIONS ON MOTIVATION OF LLM AND THE VERSION CHOICE OF LLM

953 Actually, the architectures (backbones) for time-series forecasting in are all different variants of 954 transformer, i.e., Timer (Liu et al., 2024), UniTST (Authors), MOMENT (Goswami et al.), 955 Chronons (Ansari et al.) are transformer-based and Lag-Llama (Rasul et al.) is a recent LlaMA-956 based model. Noted that these above-mentioned transformer architectures require training from 957 scratch and the parameter scales of these models range from Million level to Billion level accord-958 ing the scale of datasets. In contrast, our ZeroTS is a rather lightweight model compared to them, which inherits the architecture of GPT-2 without training or fine-tuning such GPT. Our model is 959 with 50M Trainable Param with 700M of GPT-2 for inference, which is lighter than above methods. 960 Furthermore, in fact, there is not an explicit boundary between the conventional foundation models 961 and LLM, the full name of GPT-2 is generative pretraining Transformer at Version 2, where the 962 inner structure is also transformer. We further conduct empirical tests to demonstrate the perfor-963 mance difference among LLM, conventional foundation model and a more recent LLM LlaMA-3 in 964 Section A.4.4, where the results suggest the superior effectiveness of LLM.

965 966 967

A.2.3 **COMPARISON WITH FINE-TUNING MODELS**

It is acknowledged that LLM without fine-tuning leads to less flexibility of language model. But
 actually, fine-tuning the LLM requires substantial sequence observations and large-scale computa tional costs, this cannot well adapt to the zero-shot time series learning scenarios. To this end, we
 consider detour the fine-tuning of LLM and propose a lightweight learnable module to remedy the
 non-fine-tuning of LLM, which allows interactions among real-world external data, LLM model,

as well as model error feedback, thus guiding the fusion of retrieved time series and activating the maximal power of LLM. We elaborate the motivation and reasons for our designs.

- Data availability and computational costs are challenges cannot be neglected. LLM is a well pre-trained model with large number of parameters equipped with sufficient semantic patterns within sequences and documents. Therefore, fine-tuning and impacting them will require feeding into large-scale new data and substantial computational costs. However, in in our zero-shot research, the data availability and computational power are big challenges and cannot satisfy enabling the changes of LLM. That's to say, limited available data can nearly fail to enable changes of parameters in LLM. Let alone for zero-shot learning, it is practically impossible to fine tune LLM (parameters ranging from 700M to 8B) with new arriving data in a real-time manner. Thus, we devise a more lightweight manner which directly utilize the LLM model and learn a small model for prompt construction. Regarding efficiency, we have provided analysis in Table.5 of manuscript and we can observe that our ZeroTS only consumed approximately 1/4 memory overloads and 1/7 inference speed against peer models, indicating its superior efficiency.
- 987 It is necessary to incorporate additional inputs as prompts to supplement sequential 988 pattern cues. To remedy the model adaptation capacity, we especially introduce and pre-989 serve an external time-series base (data set), which can be dynamically updated to contin-990 uously supplement the new knowledge from real world. To implement this, we couple the 991 emerging Retrieval Augmented Generation (RAG) technique with time-series learning and devise TS-RAG with diverse efficient retrieval strategies. With dynamically updated and 992 retrieved series that share similarity with targeted ones, the series can provide supplemen-993 tary patterns and leading evolution patterns incorporated into retrieved ones for prompting 994 LLM outputs. 995
 - It is necessary to devise active feedback for optimizing targeted and retrieved series representation. To adapt model with real-world data and ground-truth and equip the model with adaptation capacity, we design a ReinLLM to cohere the feedback from LLM with prepositive series representation learning.

Overall, our ZeroTS is designed tailored for zero-shot time series learning from a new perspective of data-model interactions with a great trade-off, i.e., designing a lightweight learnable ReinLLM adapter. We believe such trade-off has well leveraged the power of LLM and allowed interactions between data and model sides, thus providing sufficient and flexible complementary to existing LLM for series forecasting.

1005 1006 1007

1008

975

976

977

978

979

980

981

982

983

984

985

986

996

997

998

999

A.2.4 **DISTINGUISH ZEROTS WITH OTHER CROSS-DOMAIN LEARNING AND** MULTI-TASK LEARNING WORK

It is interesting to compare our ZeroTS with other cross-domain learning and multi-task learning 1009 work, which is also prevailing in time-series learning and spatiotemporal forecasting. It is a good 1010 idea to train cross-domain model using data from different domains, and it has been verified for its 1011 effectiveness in pioneering research of UniST (Yuan et al., 2024) and CMuST (Yi et al., 2024). We 1012 can name it as unified learning or continuous multi-task learning. For UniST (Yuan et al., 2024), it 1013 falls into the idea of training a cross-domain model using multi-sourced/domain data and it exactly 1014 demonstrates the improvement and zero-shot capacity. For training and fine-tune scheme, CMuST 1015 devises a task-level continuous learning to iteratively fine-tune the new task with partial neuron 1016 stable (Yi et al., 2024), which also provides additional capacity in zero-shot and cold-start scenarios 1017 for dynamic learning. Even so, we argue that the unified model cannot adapt all scenarios and lack flexibility to continuous update (Yuan et al., 2024) while continuous learning lacks the flexibility 1018 to retrieve any external knowledge and requires additional fine-tune parameters, limiting the zero-1019 shot learning capacity. We think the key point of this research line is to decouple the invariant 1020 and variable parts, and transfer the invariances to new domain (task) while assimilate changes from 1021 environments. 1022

A toy example test on cross-domain learning. A toy experiment on NYC (with crowd in/out and taxi pick/drop four domain series) can be obtained from Ref.ppyiget, when training on other 3 tasks, and the 4th task with only 25% samples, the cross-domain learning can improve performance of MAPE by 5% from best baseline 0.473 (PromptST) to multi-task learning setting 0.450 (cross-

1028 1029 Table 5: Empirical complexity comparison on long-term prediction of ETTh1

Length	E	TTh1-96		E	TTh1-336	
Metric	Trainable Param.(M)	Mem.(MiB)	Speed(s/iter)	Trainable Param.()	Mem.(MiB)	Speed(s/iter)
Llama (32)-QLoRA	50.29	45,226	0.697	50.37	49,374	0.732
Llama (32)-Reprogram	6.39	32,136	0.517	6.48	37,988	0.632
ZeroTS (GPT2)	53.44	8,244	0.064	55.51	8,416	0.096

1030 1031 1032

domain). Even so, cross-domain learning still suffers two critical issues 1) only correlated tasks and
domains can be utilized to reinforce each other and blindly involving various tasks can lead to model
suffering noise and distractions for main predictions, 2) the unified model may not accommodate
patterns from all tasks that limits the extendibility.

Summary of discussions against other related works. To conclude, our ZeroTS is a totally different perspective to address zero-shot and generalization task, which detours the large model fine-tune, and is with more flexibility and efficiency without considering task-level or domain-level similarity. Our solution allows any time-series to be retrieved, and incorporated into the model dynamic update, utilizes the additional temporal evolution cues (mostly similar to target one) to prompt the LLM for output prediction. ZeroTS can simultaneously exploit the power of external knowledge in a plug-and-play manner and informative semantics and regularity in LLM in a lightweight manner.

1044 1045

A.3 THEORETICAL ANALYSIS AND COMPLEXITY EFFICIENCY ANALYSIS

In this subsection, we provide some theoretical analysis from information theory and complexity efficiency analysis to demonstrate the effectiveness and efficiency in a formal manner.

1048 1049 1050

1074

A.3.1 EFFECTIVENESS OF RAG.

We provide a brief but powerful analysis from the perspective of information theory. Assume all learnable parameters $\Theta = \{W_i\}$, the targeted series $X^T = \{X_1^T, ..., X_T^T, Y_{T+1}^T, ..., Y_{T+Q}^T\}$ where X^T denotes the input features while Y^T is the output target. Since entropy can measure the degree of chaos in data, which is equivalent to the learning difficulty of corresponding dataset and thus their predictability. The smaller entropy reflects less discrepancy and larger predictability (DelSole & Tippett, 2007). Then the conventional prediction model \mathcal{M} maximizes the mutual information between the input previous series of $H(X^T, Y^T)$.

When RAG is introduced for zero-shot prediction, the retrieved series $\mathbb{X}^R = \{X_1^R, ..., X_K^R\}$ can be considered as the informative knowledge from an external database (dataset). If the similarity between targeted series and retrieved ones $\text{Dist}(X_i^R, X^T) > th_{sim}$ where th_{sim} is a significant threshold indicating the overlapping patterns between retrieved and target series, we can derive the objective of \mathcal{M}^* is equivalent to minimize the entropy $H(X^T, Y^T | \mathbb{X}^R)$. According to the condition reduction principle in entropy (Ash, 2012), we first have $H(X) > H(X | \mathbb{X}^R)$ and $H(Y) > H(Y | \mathbb{X}^R)$. We then demonstrate the joint entropy also satisfies the condition reduction principle.

1066

$$H(X, \mathbf{Y} | \mathbb{X}^R) = H(\mathbf{Y} | \mathbf{X}, \mathbb{X}^R) + H(\mathbf{Y} | \mathbb{X}^R)$$
1067
1068
1069
1070

$$H(\mathbf{Y} | \mathbf{X}) + H(\mathbf{X} | \mathbb{X}^R)$$

$$(13)$$

1071 We can then systematically derive that introducing the retrieved series with explicit overlapping 1072 patterns can improve the informativeness of prediction model and squash the information from data 1073 aspect into learnable parameters Θ , thus enabling quickly adapting to LLM and new unseen data.

1075 A.3.2 COMPLEXITY ANALYSIS.

The complexity analysis is two-fold with two stages in our ZeroTS, i.e., TS-RAG for series retrieval and training process of ReinLLM adapter.

Theoretical analysis. First, as discussed in above technical details, there are M times of computation for retrieving the matched Top-K series and the proposed ρ - ρ -cut-off over sampling and



Figure 6: Analysis of cut-off coefficient ρ and step size η

1105

1108 dimensions, thus the full computational loads have been reduced to $\rho^2 M$, where ρ is expected to sat-1109 isfy $0 < \rho < 0.5$. Second, considering the training process of learnable parameters in our Adapter, 1110 the trainable parameters are in a small scale, consisting of temporal convolutions and corresponding 1111 action selection, retrieved series representation, aggregation learner of retrieved series where it is 1112 noted that none fine-tune for LLM.

1113 **Empirical analysis.** The comprehensive empirical comparisons on complexity are reported in Ta-1114 ble. 5. We report the total number of trainable parameters (in million), GPU memory (in mebibyte, 1115 MiB), and running timeseconds per iteration (s/iter) by taking ETTh1-96 and ETTh1-336 as in-1116 stances. The compared models are retrieved from TimeLLM (Jin et al., 2023). The number of 1117 parameters are slightly more than others, where the reason lies in introducing learnable scheme 1118 into retrieval of RAG, which makes the trade-off between the high-dimension retrieval process and 1119 low-dimension but learnable retrieval. Even though, we can observe that our ZeroTS has an overwhelmingly superiority on memory and speed metrics, nearly 1/4 memory overloads and 1/7 1120 inference speed against peer models, indicating our work has significantly improved the mem-1121 ory space and inference speed. Such superiority can enable ZeroTS extremely friendly to edge 1122 computing and resource-limited services. We also believe the number of parameters can be reduced 1123 by adjusting the hyperparameters of retrieval process, and learnable weights for efficiency tradeoff. 1124

Then we can conclude our our retrieval is efficient. Even though, we believe that efficiency is 1125 a permanent issue in retrieval and deep learning training process. In the future, we are going to 1126 investigate more efficient strategy to retrieve more accurate series as prediction prompts and reduce 1127 the learnable module for quick data adaptation. 1128

- 1129 1130
- A.4 ADDITIONAL EXPERIMENTAL RESULTS 1131
- 1132
- In this subsection, we first provide the detailed datasets and metrics for evaluations, and then provide 1133 more analysis regarding hyperparameter settings.

			TIOTIZOII-WI	se results	S OII Zelo-S	not preui	CHOI	
	MAE	TimesNet	Time-LLM	GPT4TS	LLMTime	DLinear	PatchTST	ZeroTS(Ours)
	96	0.387	0.337	0.374	0.576	0.4	0.35	0.336
	192	0.429	0.374	0.417	0.586	0.460	0.400	0.368
$ETTh1 \rightarrow ETTh2$	336	0.451	0.415	0.444	0.637	0.505	0.428	0.396
	720	0.458	0.420	0.452	1.034	0.589	0.443	0.417
	Avg	0.431	0.387	0.422	0.708	0.488	0.405	0.379
	96 0	0.313	0.293	0.315	0.563	0.357	0.304	0.271
	192	0.342	0.312	0.342	0.654	0.413	0.339	0.31
ETTh1→ETTm2	336	0.371	0.365	0.374	0.728	0.465	0.373	0.341
	720	0.419	0.39	0.422	1.531	0.573	0.424	0.387
	Avg	0.361	0.34	0.363	0.869	0.452	0.36	0.327
	96	0.601	0.452	0.577	0.777	0.555	0.465	0.485
	192	0.61	0.461	0.559	0.82	0.568	0.509	0.501
ETTh2→ETTh1	336	0.626	0.482	0.578	0.864	0.577	0.515	0.544
	720	0.648	0.502	0.597	1.461	0.596	0.561	0.558
	Avg	0.621	0.474	0.578	0.981	0.574	0.513	0.522
	96	0.324	0.276	0.329	0.563	0.336	0.309	0.261
	192	0.352	0.315	0.346	0.654	0.369	0.345	0.324
ETTh2→ETTm2	336	0.383	0.337	0.376	0.728	0.397	0.379	0.365
	720	0.446	0.417	0.429	1.531	0.442	0.421	0.391
	Avg	0.376	0.341	0.37	0.869	0.386	0.365	0.335

Table 6: Horizon wise results on zero shot prediction

1150 1151

1134

4405

1152

1153 A.4.1 DATASET.

1154 The datasets we exploit for zero-shot predictions are ETTh1, ETTh2, ETTm1, ETTm2, following 1155 literature (Jin et al., 2023). These datasets provide electricity consumption indexes of urban users, 1156 where ETTh1 and ETTh2 are collected on hour-level, while ETTm1 and ETTm2 are collected on 1157 15-minute levels. They are adopted to be trained on one set and tested on another one, and this 1158 is a cross-transfer process within any other two datasets to imitate (implement) zero-shot settings. Regarding long-term prediction, we take totally 8 datasets for evaluation, where we further employ 1159 additional four datasets those are benchmarking long-hrozion prediction models, i.e., Weather, ECL, 1160 Traffic, ILI (Wu et al., 2023). Specifically, the dataset Weather is comprised of one-year records from 1161 21 meterological stations in Germany with sampling rate of 10 minutes. The Traffic dataset includes 1162 862 traffic sensors across the State of California, with a sampling rate of 1 hour. The influenza-like 1163 illness (ILI) contains records of patients experiencing severe influenza with complications. 1164

1165 A.4.2 EVALUATION METRIC.

We exploit mean absolute error (MAE) and mean square error (MSE) as metrics for both zero-shot and long-horizon forecasting.

1169 1170 1171

1172

$$MAE = \frac{1}{Q} \sum_{t=T+1}^{T+Q} |Y_t - \hat{Y}_t|; \quad MSE = \frac{1}{Q} \sum_{t=T+1}^{T+Q} (Y_t - \hat{Y}_t)^2$$
(14)

1173 A.4.3 **PERFORMANCE COMPARISON ON HORIZONS.**

1174 Drawing the results for each prediction horizon brings more informativeness and helps better understanding of our model. Here, we report the results for each horizon on long-term forecasting and zero-shot forecasting, where MAE results of four setting are reported in respective tasks, i.e., ETTh1 \rightarrow ETTh2, ETTh1 \rightarrow ETTm2, ETTh2 \rightarrow ETTh1, ETTh2 \rightarrow ETTh1 for zero-shot forecasting, and ETTh1, ETTh2, ETTm1, ETTm2 for long-term prediction. The results can be found in Table. 6 and Table. 7.

As observed in above tables, we can conclude our solution outperforms baselines on most horizons, and there is no skewing on one horizon, which enhance the effectiveness of our ZeroTS.

1183

A.4.4 ABLATION ON RECENT LLAMA-3 AND PATCHTST

To further verify the motivation of exploiting LLM as a tool for prediction and the choice of GPT-2, we conduct a series of further experiments by combining RAG with respective recent Llama-3 (with 8 Billion para.), and PatchTST. We also demonstrate the performances of pure PatchTST as comparison to enhance the rationality of utilizing GPT models. The results are shown in Table. 8.

	MAE	TimesNet	Time-LLM	GPT4TS	DLinear	PatchTST	AutoFormer	Informer	ZeroTS
	96	0.402	0.392	0.397	0.399	0.399	0.459	0.713	0.3
	192	0.429	0.418	0.418	0.416	0.421	0.482	0.792	0.3
ETTh1	336	0.469	0.427	0.433	0.443	0.436	0.496	0.809	0.4
	720	0.500	0.457	0.456	0.49	0.466	0.512	0.865	0.4
	Avg	0.450	0.423	0.455	0.437	0.430	0.487	0.795	0.4
	96	0.374	0.328	0.342	0.353	0.336	0.388	1.525	0.3
	192	0.414	0.375	0.389	0.418	0.379	0.452	1.931	0.3
ETTh2	336	0.452	0.409	0.407	0.465	0.380	0.486	1.835	0.3
	720	0.468	0.42	0.441	0.551	0.422	0.511	1.625	0.4
	Avg	0.427	0.383	0.412	0.446	0.379	0.459	1.729	0.3
	96	0.375	0.334	0.346	0.343	0.342	0.475	0.571	0.3
	192	0.387	0.358	0.372	0.365	0.369	0.496	0.669	0.3
ETTm1	336	0.411	0.384	0.394	0.386	0.392	0.537	0.871	0.3
	720	0.45	0.411	0.421	0.421	0.42	0.561	0.823	0.4
	Avg	0.406	0.372	0.403	0.378	0.380	0.517	0.734	0.3
	96	0.267	0.253	0.262	0.269	0.255	0.339	0.453	0.2
	192	0.309	0.293	0.301	0.303	0.292	0.34	0.563	0.
ETTm2	336	0.351	0.329	0.341	0.342	0.329	0.372	0.887	0.3
	720	0.403	0.379	0.401	0.421	0.385	0.432	1.338	0.3
	Arra	0 333	0 313	0 3 3 9	0.333	0.315	0.371	0.810	0.3

MSE RAG+Llama-3		RAG+PatchTST PatchTS		RAG+GPT2(ZeroTS)		
ETTh1→ETTh2	0.345	0.361	0.380	0.350		
ETTh1→ETTm2	0.274	0.286	0.314	0.272		
ETTh1	0.402	0.410	0.413	0.403		
ETTh2	0.312	0.319	0.330	0.315		

¹²¹¹ 1212 1213

It is noted that a more powerful LLM can improve the performances maximumly by 1.42% (obtained on ETTh1 \rightarrow ETTh2 scenario from 0.350 to 0.345, others are ranging from 0.2% to 0.9%), but the inference efficiency of larger LM (e.g., Llama-3: 8 Billion) has dropped significantly. When replacing LLM with PatchTST, the performance will drop by 1% to 5% (maximumly arrive at ETTh1 \rightarrow ETTm2 from 0.272 to 0.286). We can witness the improvement of LLM is relatively deserve for its efficiency and enhance the choice of GPT-2 in ZeroTS. Thus, we can conclude that our GPT-2 is sufficient for our task and more powerful LLM can further sacrifice efficiency which may not deserving for our predictions.

1221 1222

1235

A.4.5 OUT-OF-THE-BOX ZERO-SHOT TESTS

To further verify the effectiveness, we have expanded experiments to a more general setting, i.e., with no particular task-specific exemplars for tuning. In our ZeroTS, we utilize the non-training GPT-2 with RAG (retrieval strategy TS-RAG) for implementation. ZeroTS has degenerated into a pure RAG model without learnable parameters and strategies. We also take Llama-3, PatchTST as comparisons with fair settings. Three typical datasets, i.e., ETTh1, ETTm1, and Traffic are selected to illustrate the out-of-the-box prediction results of our non-trainable version of ZeroTS. The results are illustrated in Table. 9.

Actually, to this end, our design of RAG+GPT-2 can outperform conventional model PatchTST, and achieve comparable performances with RAG+Llama-3 but with sacrificing large-scale computational loads. Therefore, we can conclude that the choice of GPT-2 is fine and superior to others regarding efficiency and effectiveness trade-off.

1236 1237							
	Table 9: Out-of-the-box zero-shot tests						
1238		RAG+Llama-3	RAG+PatchTST	RAG+GPT2 (ZeroTS)			
1239	Full zero-shot on ETTh1	0.442	0.480	0.457			
1240	Full zero-shot on ETTm1	0.440	0.465	0.443			
1241	Full zero-shot on Traffic	0.456	0.504	0.472			

2								
3	Table 10: Significance test on zero-shot forecasting							
		Time	-LLM	Patch	nTST			
		p-Value: MAE	p-Value: MSE	p-Value: MAE	p-Value: MSE			
	ETTh1→ETTh2	0.050*	0.065	0.045*	0.058			
	$ETTh1 \rightarrow ETTm2$	0.048*	0.067	0.038*	0.045			
	ETTh2→TTh1	0.074	0.062	0.075	0.065			
	$ETTh2 \rightarrow ETTh1$	0.050*	0.074	0.043*	0.040*			
	T	able 11: Significa	nce test on long-	term forecasting				
		Time-LLN	1	PatchTST				

	Time-	LLM	PatchTST		
	p-Value: MAE	p-Value: MSE	p-Value:MAE	p-Value: MSE	
ETTh1	0.0003*	0.00025*	0.05*	0.05*	
ETTh2	0.0001*	0.0002*	0.065	0.03*	
Traffic	0.075	0.063	0.05*	0.03*	
ILI	0.055	0.072	0.06	0.0355*	

1255 1256 1257

10/0

1260 A.4.6 SIGNIFICANCE TEST

To show the significance, we further conduct the significance tests by comparing ZeroTS against two best baselines Time-LLM and PatchTST on four selected tasks of both zero-shot learning and long-term forecasting. We implement it via t-test with scipy package of Python day and night, and demonstrate the empirical significance of ZeroTS where the p-value is denoted p* when p < 0.05. The results are demonstrated as below Table. 10 and 11.

As observed in above two tables, on zero-shot forecasting, the results on 3 tasks out of 4 are of significance, while almost all results are of significance on selected long-term forecasting tasks (9 p-valuer results out of 16 are significant for MSE and MAE) where ZeroTS can exactly outperform Time-LLM and PatchTST on most scenarios by passing t-test. The minor scenarios our ZeroTS is inferior to these two baselines may be the dataset property and the fine-tune process on model adaptation where Transformer layers in our ZeroTS are not fine-tuned.

Thus, we still believe our work is novel insight that enables trade-off among data intelligence, model
intelligence and efficiency issue. And our contributions lie in two aspects, i.e., not only the new perspective with RAG for data side adaptation and the ReinLLM coupling reinforcement learning with
LLM to facilitate the optimization of a small representation model, but also the ultra-lightweight
learnable parameters in our ZeroTS, providing possibility in edge computing and resource-limited
services.

1279

1280 A.4.7 MORE ANALYSIS ON HYPERPARAMETERS

The vital hyperparameters of our model can be five-fold illustrated in Table.6, number of retrieved auxiliary series K, representation dimensions of both target series and auxiliary series, cut-off coefficient ρ , step size for action selection η , as well as dimension of retrieved series.

The most important one is the number of retrieved auxiliary time series, which is concerned with both effectiveness and efficiency. The more series are retrieved, the more informativeness of prompts in LLM and increases the overall performance and generalization. But with more information retrieval, the computational costs on retrieval and fusion become higher. We can see that when *K* arrives 6, our ZeroTS becomes pretty good performances and satisfies the trade-off.

Regarding representation dimensions for both targeted and retrieved series, we let them range from {8, 16, 24, 32, 64} where it also requires the trade-off. It is known that larger dimensions usually account for more learning and fitting capacity, which corresponds to factual performance illustrated in Fig. 4. We thus take 64 as the final dimension as larger dimension will significantly increase the computational burden, e.g. 128.

Considering the coefficient for cut-off retrieval ρ **.** Such parameter determines which ranges of retrieved series can be obtained by our model and input into LLM as prompts. The larger ρ , sug-

	ETTh1 \rightarrow	$ETTh1 \rightarrow ETTm2$	$ETTh2 \rightarrow ETTh1$	$ETTh2 \rightarrow ETTm2$	$ETTm1 \rightarrow ETTh2$	$ETTm1 \rightarrow ETTm2$	$ETTm2 \rightarrow ETTm2$	ETTm2
	ETINZ	ETTM2	EIIni	ETTM2	ETINZ	ETTm2	ETTN2	EIIm
Top-K retrieved series	6	6	6	6	6	6	6	6
Dimension of targeted	64	64	64	64	64	64	64	64
Dimension of retrieved TS representation	64	64	64	64	64	64	64	64
Cut-off coefficient ρ	0.3	0.3	0.3	0.3	0.4	0.4	0.4	0.
Step size for action selection η	0.05	0.05	0.05	0.05	0.05	0.05	0.05	0.0

Table 12: Hyperparameter setting on zero-shot prediction

1296

100-

1307 gests more available series, and thus indicate more cues and evolution patterns obtained for LLM, 1308 yielding increasing model performances. To achieve the trade-off, we search the optimal one rang-1309 ing from $\{0.2, 0.3, 0.4, 0.5\}$ and find 0.3 or 0.4 as optimal point for respective datasets. Actually, 1310 finer granularity can be explored but it requires more costs.

Step size for action selection η . It is the parameter for discretizing the coefficient α selection process in Eq.(7) thus for its adjustment. Since the coefficient α ranges from 0 to 1, the step stride is selected from {0.05, 0.1, 0.15, 0.2} where larger step can lead to leap over optimal point while smaller step leads to inefficient dilemma. Finally, according experimental analysis and tests, we choose $\eta = 0.05$.

Empirical results. The hyperparameter adjustment on two selected scenarios are in Figure. 6 and Figure. 7. The final settings for zero-shot predictions are provided in the following Table. 12.



Figure 7: Hyperparameter analysis on ETTh1→ETTm2.

1330 1331

1333

1332 A.5 LIMITATIONS AND FUTURE WORK.

Our work focuses on retrieving high-quality time series to enhance long-term and zero-shot series 1334 prediction with help of LLM, and enabling the data-model interactions with a small-scale learnable 1335 adapter. However, even effective and efficient, every coin has two sides. In this subsection, we dis-1336 cuss the potential limitations from three aspects. Limitations. First, due to limited computational 1337 resources and complexity constraints, we take GPT-2 as the backbone for prediction, which can be 1338 extendable to an update-to-date LLM such as Llama-3 when computational resources are sufficient. 1339 Second, reducing the complexity is a permanent topic for large-scale retrieval process and how to 1340 provide a selective process for filtering out irrelevant series. Third, as ZeroTS reveals less com-1341 petitive on generalization from longer horizon to shorter horizon in Table. 1, it is suggested that the further improvement on how to transfer and adapt the coarse-grained data to fine-grained ones 1342 with series learning techniques. Therefore, our future work are also decomposed as three aspects on 1343 model architectures, high-efficient series retrieval and cross-granularity sequence learning. We will 1344 further point out a practical deployment solution to our ZeroTS. 1345

Future works. 1) Higher version of LLM. given that our work provides a scheme integrating
reinforcement with time series, modeling more optimization options in series forecasting as discrete
actions in reinforcement architecture, and it is possible to utilize higher version of GPT to achieve
more accurate performances if time and resources are available. Then our system can support for
more diverse and challenging scenarios such as out-of-distribution setting, cross-domain settings.

1374

1403

1350 2) Potential solution to more efficient retrieval. A potential alternative towards the retrieval 1351 process can be pre-clustering the series prototypes, namely prototype-based retrieval. When the 1352 auxiliary series base constructed, we can first pre-cluster all the series into clustering as where one prototype represents a corresponding clustering. When a new series is added into the base, it will be 1353 1354 assigned with one clustering and corresponding prototype can be updated. Given a targeted series, we can retrieve auxiliary series by matching target one to prototypes and reduce the complexity 1355 from $\rho^2 M$ to linear to M*log(N), where N is the number of total series in auxiliary base, M is 1356 the number of clustering of total series. 3) Overcoming long-short term transfer. We speculate 1357 the reason behind this phenomenon is the changes of forecasting granularity. On overcoming such 1358 challenge, we can design a multi-resolution prediction scheme in training period, and enable the 1359 model to aware of the multi-grained evolution patterns and interpolation in the sequential elements 1360 with new learning objective. Then when the transfer is implemented, the coarse-to-fine pattern can 1361 be transferred from known granularity to unknown ones. We can also devise a temporal interpolation 1362 strategy and sequence augmentation to imitate series in fine-grained resolution, thus contributing to 1363 the adaptation on following transferred domain.

4) Practical deployment of ZeroTS. Furthermore, implementing the AI solutions in an efficient 1365 database framework is exciting (Ding et al., 2018). Since ZeroTS is an effective zero-shot forecast-1366 ing foundation model, the practicality of its deployment is of importance. Considering the Time-1367 Series oriented ZeroTS, especially with external time-series dataset by TS-RAG, we introduce an 1368 open-sourced temporal database named TD-engine (Taosdata, 2020), which is an industrial-grade 1369 database tailored for time-series and dynamic streaming data, to support high-performance, dis-1370 tributed IoT, industrial big data platform. Due to its open-sourced property and effciency on tem-1371 poral data processing and indexing, it is expected to deploy our solution and algorithm on the such grant platform and we also believe the TD-engine can help our solution become more efficient and 1372 effective. 1373

1375	Algorithm 2 HNSW Algorithm Pseudocode
1276	1: procedure BUILDHNSW(data, M, L, efConstruction)
1370	2: for all vectors in data do
1377	3: $levels \leftarrow CalculateLevels(L, vector)$
1378	4: for all level in levels do
1370	5: If level is not in graph then
1070	7: end if
1380	8: AddVectorToLevel(graph[level], vector, M)
1381	9: ConnectVectorWithNeighbors(graph[level], vector, efConstruction, M)
1382	10: end for
1383	11: end for
1303	2: end procedure
1384	(1): procedure SEARCHHNSW (graph, query, M, ef Search)
1385	14. current_node ← ChooseStartNode(graph[current_level])
1386	16: results \leftarrow empty_list()
1387	17: visited $\leftarrow empty_set()$
1307	18: while current_level ≥ 0 do
1388	19: candidates ← FindNeighbors(graph[current_level], current_node, query, efSearch)
1389	20: candidates \leftarrow SortByDistance(candidates, query)
1390	21: Ior all candidate in candidates do
1201	22. In candidate is not in visited then 73: visited add(candidate)
1391	24: if length(results) < M then
1392	25: results.append(candidate)
1393	26: else
1394	27: break
1205	28: end if
1355	29: end II
1396	r_{0} , r_{0} $r_{$
1397	Ω break
1398	33: end if
1200	34: $current_node \leftarrow BestNeighbor(results, query)$
1222	35: $current_level \leftarrow current_level - 1$
1400	36: end while
1401	5/: return results
1402	oo. enu procedure

Figure 8: Pseudocode of HNSW