

A novel attention-enhanced LLM approach for accurate power demand and generation forecasting

Zehuan Hu^a, Yuan Gao^{b,*}, Luning Sun^a, Masayuki Mae^a

^a*Department of Architecture, Graduate School of Engineering, The University of Tokyo, Tokyo, Japan*

^b*International Institute for Carbon-Neutral Energy Research (WPI-I2CNER); Kyushu University, 744 Motooka, Nishi-ku, Fukuoka-shi, Fukuoka, 819-0395, Japan*

Abstract

Accurate forecasting of electricity demand and generation is crucial for efficient grid management and sustainable energy planning. While large language models (LLM) have shown promise in various fields, their application to time series forecasting presents challenges, including limited cross-channel information capture and the complexity of prompt design. In this study, we propose a novel framework that combines multiple attention mechanisms with LLM, enabling effective feature extraction from both target and non-target variables without the need for prompt engineering. We conducted extensive experiments using real-world electricity demand and generation data from multiple regions in Japan to evaluate the proposed model. The results demonstrate that our model outperforms state-of-the-art LLM-based and other time-series forecasting models in terms of electricity demand and generation forecast task, achieving better performance than the latest LLM-based models without using prompts or increasing model size. Compared with the Long short-term memory network (LSTM), the mean absolute error (MAE) is reduced by 20.8%. Compared with the previous time-series LLM, the proposed model reduces memory usage by 49.3% and shortens training time by 35.7%. Additionally, the proposed model exhibits superior generalization ability, maintaining high performance even in

*Corresponding author
Email address: yuangao1120@gmail.com (Yuan Gao)

zero-shot learning scenarios. Compared with the LSTM, MAE on the four test datasets is reduced by 16.6%.

Keywords: Large language model; Multi-Attention mechanism; Electricity demand and generation forecasting; Zero-shot learning

1. Introduction

2. 1.1. Background

3 In recent years, global electricity consumption has surged due to rapid ur-
4 banization, industrialization, and the proliferation of digital technologies [1].
5 This increasing demand for power places immense pressure on energy providers
6 to meet consumption needs while adhering to sustainability goals [2]. With the
7 looming threat of climate change, governments and organizations worldwide are
8 focusing on reducing carbon emissions and transitioning toward cleaner energy
9 sources [3, 4]. Achieving these ambitious targets requires not only the develop-
10 ment of renewable energy infrastructures but also the implementation of smart
11 energy management systems that can efficiently balance consumption with pro-
12 duction [5, 6].

13 Accurate forecasting of electricity demand and generation is crucial to en-
14 suring a stable and efficient power grid [7]. Mismatches between power supply
15 and demand can lead to significant operational challenges, including grid insta-
16 bility, power outages, and economic losses [8]. Traditional forecasting methods,
17 while effective to some extent, struggle to capture the complexity of modern
18 power systems, which are influenced by numerous factors such as fluctuating
19 renewable energy outputs, weather conditions, and dynamic market demands
20 [9]. As a result, there is an increasing need for advanced predictive models that
21 can enhance the accuracy and reliability of electricity forecasting, enabling more
22 proactive energy management and decision-making [10].

23 1.2. Electricity demand and generation forecasting

24 Currently, a variety of methods are employed for electricity demand and
25 generation forecasting, ranging from traditional statistical approaches to more

²⁶ advanced machine learning techniques [11]. Classical methods, such as autoregressive integrated moving average (ARIMA) and exponential smoothing, have
²⁷ been widely used due to their simplicity and effectiveness in modeling time
²⁸ series data methods often assume linear relationships and may struggle to cap-
²⁹ ture the non-linear complexities of modern energy systems [12]. On the other
³⁰ hand, machine learning methods, such as support vector machines (SVM) and
³¹ random forests, have been introduced to address these limitations by modeling
³² non-linear relationships more effectively [13]. Despite their statistical methods,
³³ these approaches still require extensive feature engineering and may not always
³⁴ generalize well to unseen data [14].

³⁶ In recent years, deep learning models, particularly artificial neural networks
³⁷ (ANN) and their variants, have gained significant attention in electricity fore-
³⁸ casting due to their ability to learn complex patterns directly from raw data [15].
³⁹ Recurrent neural networks (RNN), especially long short-term memory (LSTM)
⁴⁰ networks, have shown great promise in time series forecasting as they can cap-
⁴¹ ture temporal dependencies in the data [16, 17]. Additionally, hybrid models
⁴² that combine LSTM or other techniques, such as attention mechanisms are be-
⁴³ coming increasingly popular due to their enhanced prediction performance and
⁴⁴ adaptability to complex energy systems [18, 19]. A summary and review of
⁴⁵ relevant research on different machine learning methods of electricity demand
⁴⁶ forecasting is shown in Table 1.

Table 1: Literature review for electricity demand & generation forecasting.

Ref.	Methods	Datasets
[20]	ANN	Electricity demand data of Turkey
[21]	Autoregressive	Electricity demand data from the Nordic elec- tricity market
[22]	LSTM; SVM	Electrical load data; weather factor data of school building

Ref.	Methods	Datasets
[23]	Classical statistical; autoregressive models	Electrical demand data of Ukraine
[24]	K-nearest neighbor (KNN), linear regression, random forest (RF); SVM	Demand and renewable power generation data from South African electricity public utility company
[25]	Transformer-based model	Electricity demand data of one city in China
[26]	Elman neural network	Power consumption statistics of Serbia
[27]	Hybrid deep learning	Four buildings data from Mendeley and Kaggle
[28]	ANN	Meteorological, operational, and economic data from Mexico
[29]	Deep-Autofomer	Microgrid system data of 24 families
[30]	CNN-LSTM	Load data from individual household
[31]	KNN; SVM; RF; ANN	Electricity demand data of low-energy house

47

48 1.3. Large Language Models for time-series forecasting

49 Large language models (LLM), such as GPT (Generative Pretrained Trans-
50 former) and BERT (Bidirectional Encoder Representations from Transform-
51 ers), have revolutionized numerous fields, including natural language processing
52 (NLP), machine translation, and text generation [32]. These models, character-
53 ized by their vast number of parameters and ability to process massive datasets,
54 have shown remarkable capabilities in understanding and generating human-like

55 text [33]. Beyond traditional NLP tasks, LLMs have been increasingly adopted
56 in domains like healthcare, law, and finance, where complex, context-dependent
57 decision-making is required [34]. The rapid advancements in transformer archi-
58 tectures, which LLMs are based on, have made it possible to capture intricate
59 patterns in sequential data, offering potential applications in areas that extend
60 beyond text [35].

61 Recent studies have started exploring the application of LLMs for time se-
62 ries forecasting, leveraging their ability to model sequential dependencies and
63 long-range correlations effectively [36]. Traditional time series models, such as
64 ARIMA or LSTM, are often limited in their ability to capture the full context
65 of long sequences[37]. LLMs address these limitations by allowing parallel pro-
66 cessing of sequential data and better handling of long-term dependencies [38]
67 Research has demonstrated that LLM-based approaches, such as Time-series
68 Transformers, can outperform classical methods in various forecasting tasks.

69 At present, approaches to employing LLMs in time series forecasting can
70 be broadly classified into two categories. The first relies on prompt-based re-
71 programming of pretrained LLMs, whereby carefully crafted instructions are
72 prepended to raw time series to steer the model’s predictions [39], [40]. In-
73 spired by their remarkable conversational capabilities, these methods leverage
74 the LLM’s inherent pattern-recognition power to interpret temporal dynamics
75 without any weight updates [41]. The second paradigm involves fine-tuning the
76 LLM itself on time series data [42][43]. Chang et al. [44] subsequently proposed
77 a two-stage fine-tuning pipeline—first aligning the backbone network to time se-
78 ries representations, then training a lightweight prediction head. Sun et al. [45]
79 introduced contrastive learning objectives to sharpen the model’s temporal em-
80 beddings. A summary and review of relevant research on LLM for time-series
81 forecasting is shown in Table 2.

Table 2: Literature review for time-series forecasting by LLM.

Ref.	Target	Approaches to employing LLMs
[46]	Energy load	Directly use with predefined template prompt
[47]	Wind speed	Directly use with spatial prompts and temporal prompts
[48]	CNY-USD change rate	Directly use with prompt adjusted by reinforcement learning with human feedback
[49]	Wind power	Fine-tuning
[50]	NASDAQ-100 stock price	Directly use with predefined template prompt \$ fine-tuning
[51]	Energy load	Fine-tuning
[52]	Climate	Directly use with predefined template prompt
[53]	Wind power	Directly use with proposed soft and hard prompts
[54]	Methanol price	Fine-tuning
[55]	Stock return	Fine-tuning

82

83 *1.4. Research contribution*

84 Despite some initial applications of LLMs in the field of time series forecast-
 85 ing, significant challenges and limitations remain:

- 86 1) Pretrained LLMs have already undergone extensive training and valida-
 87 tion, so naive fine-tuning often yields only marginal gains and can even
 88 degrade performance due to overfitting; moreover, the multiple modules
 89 within modern LLM architectures make identifying which components to
 90 fine-tune a nontrivial challenge.
- 91 2) While prompt design has been widely used in other LLMs applications
 92 to enhance model performance, designing effective prompts for time se-

ries forecasting is particularly challenging. Due to the diverse nature of datasets and the complexity involved in describing historical feature relationships and interdependencies, creating a simple and generalized prompt for time series prediction tasks remains elusive.

3) LLM applications in electricity demand and generation forecasting are still largely unexplored. Given the multivariate nature of electricity data, which typically involves numerous features, developing an LLM-based model that can efficiently handle and predict these dynamics is a significant challenge.

4) Japan spans a vast distance from north to south, and there are significant regional differences in power generation, electricity demand, and climate. These variations place high demands on the generalization ability of models. However, there is limited research focused on Japan's electricity data.

To address these challenges and fill the research gaps, this study proposes a novel LLM-based model for forecasting electricity demand and generation. The key contributions of this research are as follows:

1) To enhance the model's ability to understand the relationships between different features, we design a specialized framework that allows the LLM to extract meaningful interdependencies between the target variables and other features in a simple yet effective manner.

2) We introduce a preprocessing method based on cross-attention that eliminates the need for prompt engineering, enabling multiple time series of target features to be transformed and input into the LLM simultaneously. This approach simplifies the application of LLMs to multivariate time series forecasting.

3) To validate the effectiveness and generalizability of the proposed framework, we conduct extensive experiments using historical electricity demand, generation, and weather data from four different regions in Japan.

121 **2. Methodology**

122 *2.1. Related works and Challenges*

123 In the field of time series forecasting, there are currently two mainstream
124 training strategies: Channel Independent (CI) (Fig. 1 (a)) and Channel De-
125 pendent (CD) (Fig. 1 (b)) approaches [56]. Each of these strategies has its
126 strengths and limitations. The CD strategy theoretically offers higher modeling
127 capacity because it allows the model to learn complex relationships between
128 different variables. In contrast, the CI strategy is more robust, as it focuses on
129 the self-correlation characteristics within each channel, making it less sensitive
130 to noise and more generalizable [57].

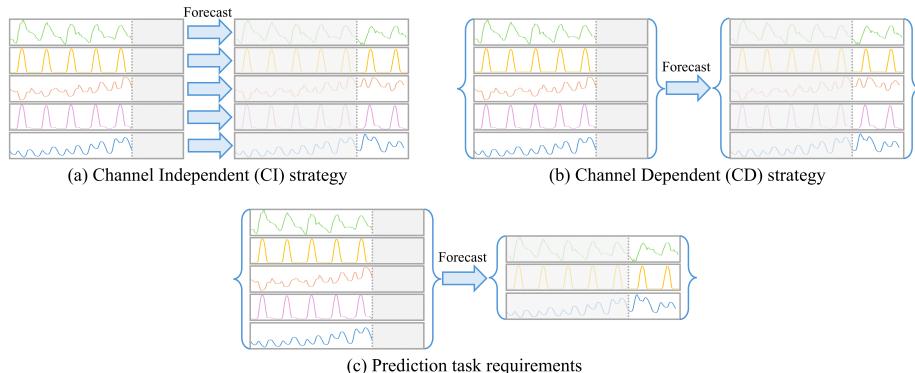


Figure 1: Training strategies and prediction requirement.

131 Meanwhile, in practical engineering applications, forecasting tasks often re-
132 quire using multiple input variables to predict a subset of target variables, such
133 as predicting future power generation and electricity demand using historical
134 power and weather data (Fig. 1 (c)). Unlike tasks that require predictions for
135 all variables or single-variable forecasting, these tasks benefit from the model’s
136 ability to consider the relationships between variables without needing to pre-
137 dict every variable. However, considering the large scale of modern LLMs and
138 their substantial computational costs, most existing applications of LLMs to
139 time-series forecasting adopt a channel-independent training strategy [36, 43].
140 This approach computes a loss for every variable—even when only a subset

141 of variables requires prediction—resulting in wasted training resources. Con-
142 versely, restricting the input to only the target variables prevents the model
143 from leveraging information in other channels, which can lead to reduced fore-
144 casting accuracy.

145 Furthermore, prompt engineering plays a crucial role in optimizing the per-
146 formance of LLMs, expanding their application range, and improving interaction
147 efficiency. For example, Cao et al. [41] introduced a Semi-Soft Prompt strategy,
148 where prompts are divided into explicit text-based prompts (hard prompts) and
149 vector-based prompts (soft prompts). The authors propose a semi-soft prompt-
150 ing strategy that generates distinct prompts corresponding to key time series
151 components: trend, seasonality, and residuals. Jin et al. [36] attached prompts
152 as prefixes to the input time series, providing background information, task
153 instructions, and data statistics. However, in practical applications, different
154 forecasting tasks and time series often have unique characteristics, making it
155 challenging to design effective prompts tailored for each scenario. This greatly
156 increases the complexity of utilizing LLMs in time series forecasting.

157 2.2. Multi-attention large language model (MultiAttLLM)

158 To address these challenges, we propose a novel model that combines the ad-
159 vantages of CI and CD strategies while eliminating the need for prompt design,
160 which called Multi-attention large language model (MultiAttLLM), as shown in
161 Fig. 2. In this model, the LLM focuses on the self-correlations of target vari-
162 ables during the initial modeling phase, while cross-channel relationships are
163 considered in a subsequent stage after the LLM output. This approach not only
164 reduces the computational resources required for training but also enhances the
165 model’s ability to capture complex interactions between variables, resulting in
166 improved forecasting performance. The architecture is composed of six main
167 components:

168 ① Word Projection: The first component of the model is the word projec-
169 tion layer, which addresses the issue of vocabulary redundancy in large

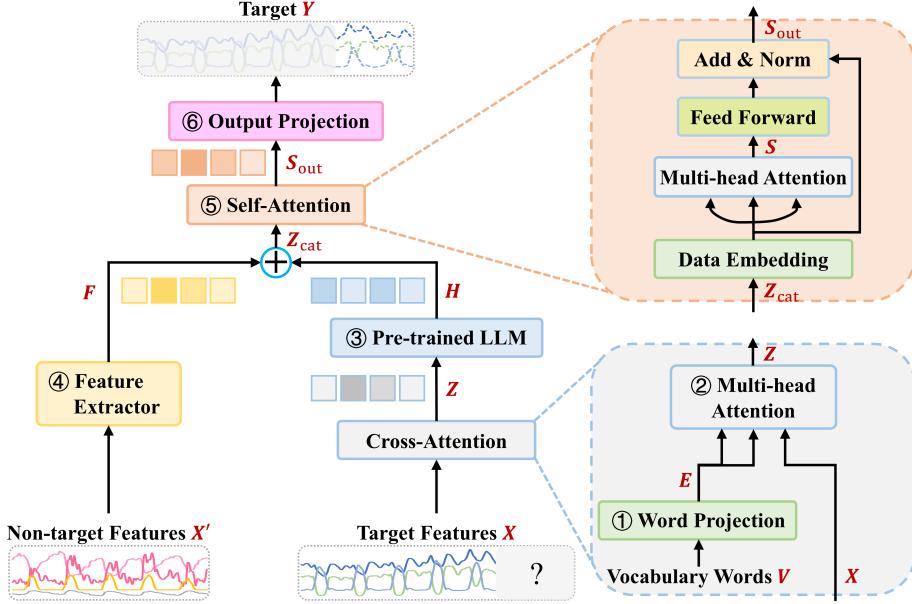


Figure 2: Topological general structure of the proposed LLM framework.

language models when applied to time series forecasting. Since the original vocabulary of LLM contains a significant amount of words unrelated to the context of time series data, the word projection layer maps the original vocabulary to a smaller, domain-specific word vector set. This allows the model to focus on the words and representations most relevant to describing time series, improving its efficiency and accuracy in forecasting tasks.

In this layer, we map the original LLM vocabulary embeddings $\mathbf{V} \in \mathbb{R}^{|V| \times d_{\text{orig}}}$ into a smaller, domain-specific space of dimension d_{wproj} , and the output \mathbf{E} of word projection layer can be calculated as follow:

$$\mathbf{E} = \mathbf{V}\mathbf{W}_{\text{wproj}} + \mathbf{b}_{\text{wproj}}, \quad \mathbf{W}_{\text{wproj}} \in \mathbb{R}^{d_{\text{orig}} \times d_{\text{wproj}}}, \quad \mathbf{b}_{\text{wproj}} \in \mathbb{R}^{d_{\text{wproj}}} \quad (1)$$

where d_{orig} is the dimension of original LLM vocabulary; d_{wproj} is the dimension of vocabulary after word projection, which is 2000 in this study.

② Cross-Attention: To enable the large language model to capture the rela-

183 tionships between different features and convert the time series data into a
 184 format that can be understood by the LLM, we employed a cross-attention
 185 module (Fig. 3). This module integrates the features by using a query-
 186 key-value (QKV) structure, where each feature can attend to all other
 187 features. This allows the model to capture complex interdependencies be-
 188 tween the target features and other input variables, effectively bridging
 189 the gap between time series data and natural language processing.

190 In this study, the target feature is used as Q , and the word vectors are
 191 used as K and V . Given an input feature matrix $X \in \mathbb{R}^{T \times f}$ (where T is
 192 sequence length, f is number of features), we compute:

$$Q = XW_q, K = XW_k, V = XW_v \quad (2)$$

193 where Q , K , and V represent the query, key, and value matrices re-
 194 spectively; $W_q, W_k, W_v \in \mathbb{R}^{f \times d}$, represent the learnable weights. The
 195 attention map and output are

$$A(Q, K, V) = \text{softmax}(QK^T / \sqrt{d_k})V \quad (3)$$

196

$$Z = A(Q, K, V)W_o, \quad W_o \in \mathbb{R}^{d \times d}, \quad Z \in \mathbb{R}^{T \times d} \quad (4)$$

197 where the attention mechanism is applied to each time step in the input
 198 sequence; d_k is the dimension of queries, keys and values, which used to
 199 reduce the impact of the input data dimension on the results.

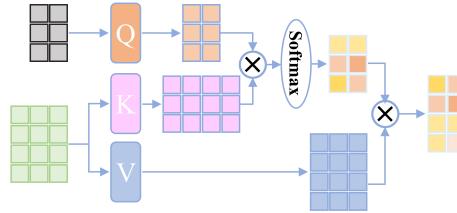


Figure 3: Topological general structure of attention mechanism.

200 ③ Frozen pre-trained LLM: After cross-attention block is a pre-trained large
 201 language model. By keeping the LLM intact, we maintain the rich lan-

202 guage understanding capabilities of the pre-trained model while focusing
203 on effective reprogramming of the input to align with its strengths.

204 In this layer, we feed the cross-attention outputs \mathbf{Z} (after adding positional
205 encodings) into a pre-trained LLM of L layers:

$$\mathbf{H} = \text{LLM}_{\text{frozen}}(\mathbf{Z} + \text{PE}), \quad \mathbf{H} \in \mathbb{R}^{T \times d} \quad (5)$$

206 where no weights inside the LLM are updated; we merely reprogram its
207 input to our time-series format.

208 ④ Feature Extractor: The feature extractor module is responsible for cap-
209 turing the relationships between non-target features, providing auxiliary
210 information that helps improve the overall forecasting performance. This
211 module can use various techniques, such as RNN or self-attention lay-
212 ers, to extract these relationships. In this study, we opt for a simple linear
213 layer to validate the effectiveness of the proposed algorithm without intro-
214 ducing unnecessary complexity. This feature extraction helps the model
215 understand the dynamics of auxiliary variables, which can influence the
216 target features.

217 To model non-target feature interactions, we apply a simple linear extrac-
218 tor:

$$\mathbf{F} = \mathbf{X}' \mathbf{W}_f + \mathbf{b}_f, \quad \mathbf{W}_f \in \mathbb{R}^{f \times d}, \quad \mathbf{b}_f \in \mathbb{R}^d, \quad \mathbf{F} \in \mathbb{R}^{T \times d} \quad (6)$$

219 This provides auxiliary embeddings capturing global feature dynamics.

220 ⑤ Self-Attention: The self-attention module is designed to enhance the un-
221 derstanding of the relationships between the target feature and the other
222 input features. By applying self-attention, the model can learn the impor-
223 tance of each feature in relation to the target, thus improving its predictive
224 accuracy. The self-attention mechanism assigns weights to each feature,
225 allowing the model to focus more on the features that have a significant
226 impact on the target. The calculation for the self-attention layer is the
227 same as the cross-attention layer. The difference here is that we merge the

228 outputs of the Feature Extractor layer and pre-trained LLM layer along
 229 the feature dimension, and then use them as the query, key, and value
 230 inputs to compute self-attention. In addition, this layer also includes a
 231 feedforward component. In the feedforward layer, a feedforward neural
 232 network processes the output embeddings \mathbf{Z} from the self-attention mech-
 233 anism to generate the final output of the transformer model.

234 In this layer, we concatenate \mathbf{H} and \mathbf{F} along the feature dimension to
 235 form $\mathbf{Z}_{\text{cat}} \in \mathbb{R}^{T \times 2d}$, then project back to d :

$$\mathbf{Z}' = \mathbf{Z}_{\text{cat}} \mathbf{W}_{\text{cat}}, \quad \mathbf{W}_{\text{cat}} \in \mathbb{R}^{2d \times d} \quad (7)$$

236 and compute standard self-attention:

$$\mathbf{Q}' = \mathbf{Z}' \mathbf{W}_q, \quad \mathbf{K}' = \mathbf{Z}' \mathbf{W}_k, \quad \mathbf{V}' = \mathbf{Z}' \mathbf{W}_v, \quad (8)$$

$$\mathbf{A}' = \text{softmax} \left(\frac{\mathbf{Q}' \mathbf{K}'^\top}{\sqrt{d}} \right), \quad \mathbf{S} = \mathbf{A}' \mathbf{V}' \quad (9)$$

238 followed by a two-layer feedforward network:

$$FFN(\mathbf{S}) = \max(0, \mathbf{S} \mathbf{W}_1 + \mathbf{b}_1) \mathbf{W}_2 + \mathbf{b}_2, \quad \mathbf{W}_1 \in \mathbb{R}^{d \times d_{ff}}, \quad \mathbf{W}_2 \in \mathbb{R}^{d_{ff} \times d} \quad (10)$$

239 where \mathbf{W}_1 , \mathbf{b}_1 , \mathbf{W}_2 , and \mathbf{b}_2 represent learnable weights and biases.

240 ⑥ Output Projection: The output projection layer consists of two linear
 241 layers and serves to convert the final output of the self-attention module
 242 into the desired feature dimension and sequence length.

243 In last layer, we map the self-attention output \mathbf{S}_{out} back to the original
 244 target dimension f_{out} :

$$\mathbf{Y} = \mathbf{S}_{\text{out}} \mathbf{W}_{\text{out}} + \mathbf{b}_{\text{out}}, \quad \mathbf{W}_{\text{out}} \in \mathbb{R}^{d \times f_{\text{out}}}, \quad \mathbf{b}_{\text{out}} \in \mathbb{R}^{f_{\text{out}}}, \quad \mathbf{Y} \in \mathbb{R}^{T \times f_{\text{out}}} \quad (11)$$

245 2.3. Benchmark

246 To validate the effectiveness and generalization ability of our proposed method,
 247 we selected five baseline models, including one traditional models, LSTM, a

248 simple linear model, DLinear, two latest novel transformer-based model, iTrans-
249 former and TimesNet, and a novel LLM-based model, TimeLLM.

250 *2.3.1. Long-short term memory network (LSTM)*

251 LSTM is a type of RNN designed to overcome the vanishing gradient prob-
252 lem that often occurs in traditional RNN when learning long-term dependencies
253 in sequential data [58]. LSTM incorporates memory cells and gating mecha-
254 nisms—input, forget, and output gates—that regulate the flow of information
255 within the network. These gates allow LSTM models to selectively retain or
256 forget information over time, making them highly effective for time series fore-
257 casting tasks that require capturing both short-term and long-term patterns in
258 the data. The detail content and calculation equations of LSTM are provided
259 in Appendix A.1 because of space constraints.

260 *2.3.2. DLinear*

261 Decomposition Linear (DLinear) is a simple linear model designed specifi-
262 cally for time series forecasting [59]. It improves forecasting accuracy by decom-
263 posing time series. Unlike the currently popular complex Transformer-based
264 models, DLinear achieves excellent performance by processing the trend and
265 cyclical components in time series through simple linear layers.

266 *2.3.3. Informer*

267 Informer is an efficient Transformer variant specifically designed for long-sequence
268 time-series forecasting [18]. It introduces a ProbSparse self-attention mechanism
269 that selects only the most informative query–key pairs, reducing the quadratic
270 complexity of standard attention. To further accelerate processing, Informer
271 employs a self-attention distilling operation that progressively shortens the se-
272 quence at intermediate layers, allowing it to scale to very long input horizons
273 while maintaining high accuracy.

274 *2.3.4. Autoformer*

275 Autoformer advances Transformer architectures by embedding series decom-
276 position directly into each model block [19]. It splits the input into trend and
277 seasonal components using a learnable decomposition layer, then applies an
278 auto-correlation mechanism to capture period-aware dependencies. This design
279 both denoises the input and enables the model to learn long-term temporal pat-
280 terns more effectively, yielding superior performance on long-horizon forecasting
281 tasks.

282 *2.3.5. iTransformer*

283 iTransformer is a novel adaptation of the traditional Transformer architec-
284 ture specifically designed for time series forecasting [60]. Unlike conventional
285 Transformer-based models that process temporal tokens (where each token cor-
286 responds to a time step with multiple feature variables), the iTransformer adopts
287 an inverted approach. It treats each time series as a separate variate token and
288 uses self-attention mechanisms to capture the interdependencies between these
289 variate tokens. This design helps the model more effectively learn multivariate
290 correlations, making it particularly well-suited for time series data with complex
291 feature relationships.

292 *2.3.6. TimesNet*

293 TimesNet is a novel model designed for general time series analysis by trans-
294 forming one-dimensional (1D) time series data into two-dimensional (2D) ten-
295 sors [17]. Traditional time series models often struggle with capturing intricate
296 temporal variations because of the limitations of processing 1D data directly.
297 TimesNet tackles this by leveraging the concept of multi-periodicity in time
298 series, where complex variations occur both within and between periods. To
299 better represent these temporal variations, TimesNet reshapes 1D time series
300 into 2D tensors, where intraperiod-variations are represented by columns and
301 interperiod-variations by rows.

302 2.3.7. TimeLLM

303 TimeLLM is a novel framework designed to adapt LLM for time series fore-
304 casting by reprogramming the input time series data [36]. Instead of fine-tuning
305 the pre-trained LLM or altering their internal architectures, TimeLLM lever-
306 ages a reprogramming technique that transforms the time series data into a
307 format compatible with LLM. This is achieved by converting time series into
308 text-like prototypes that the LLM can understand. To further enhance the
309 model’s reasoning capabilities, TimeLLM incorporates a ”Prompt-as-Prefix”
310 (PaP) approach, which enriches the input data with additional context, such
311 as task-specific instructions and domain knowledge, allowing the LLM to better
312 interpret and predict time series trends.

313 2.4. Model setup

314 For all baseline models except the LSTM, hyperparameters were set to the
315 optimal values reported in their original publications. For the LSTM and our
316 proposed model, to assess model robustness and determine optimal hyperpa-
317 rameter settings, we performed sensitivity analyses on both the LSTM baseline
318 and our MultiAttLLM framework.

- 319 • LSTM: We evaluated combinations of layer depth $L \in \{1, 2, 3, 4\}$ and
320 hidden-state dimension $h \in \{64, 128, 256, 512\}$. For each configuration,
321 we recorded MSE on the validation set. The resulting performance grid is
322 plotted in Fig. 4.

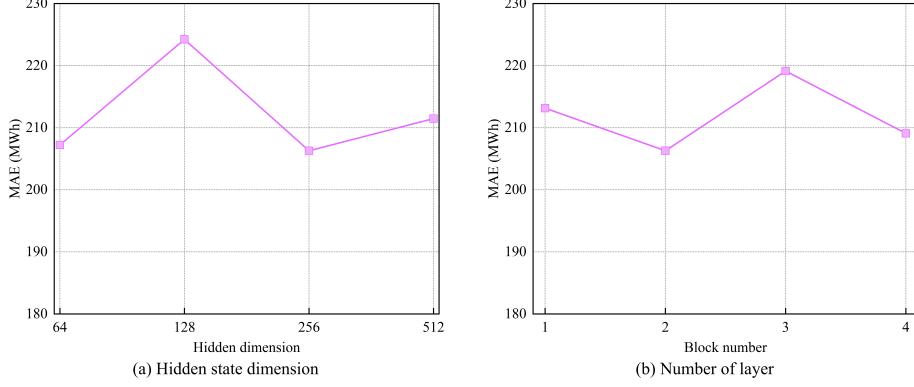


Figure 4: Hyperparameter sensitivity with respect to the hidden dimension size and number of LSTM layers (lookback window length: 72; forecast horizons: 168).

323 • MultiAttLLM: We varied the transformer model dimension $d_{model} \in \{8, 16, 32, 64, 128\}$
 324 and the number of LLM layers $N \in \{2, 4, 6, 12, 16\}$. For each configuration,
 325 we recorded MSE on the validation set. The outcomes are illustrated
 326 in Fig. 5. These experiments show that the proposed MultiAttLLM model
 327 exhibits strong robustness to hyperparameter selection, achieving optimal
 328 or near-optimal accuracy across nearly all tested configurations.

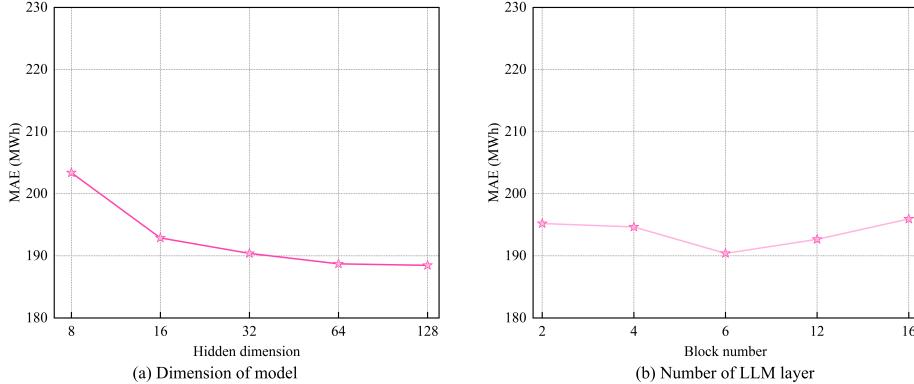


Figure 5: Hyperparameter sensitivity with respect to the dimension size of model and number of LLM layers (lookback window length: 72; forecast horizons: 168).

329 Drawing from these settings and further experimentation, we identified the
 330 parameter sets that delivered the best performance in terms of predictive accu-

331 racy and computational efficiency. Table 3 summarizes the chosen parameters
332 for each model, including the number of layers, hidden dimensions, attention
333 heads, and other key settings.

334 Each parameter set was chosen to balance model complexity with compu-
335 tational feasibility, ensuring that models could be effectively trained within the
336 limits of available resources while achieving optimal forecasting performance.
337 The parameters were adjusted based on the characteristics of the time series
338 data, such as the number of features and the frequency of observations, to en-
339 sure that the models were well-suited for the forecasting tasks.

340 To ensure fair comparisons across all models, we used consistent training
341 settings, employing the Adam optimizer with an initial learning rate of 0.001.
342 Each model was trained for a total of 10 epochs, with the learning rate reduced
343 to 95 percent of its value from the previous epoch after each epoch to facilitate
344 convergence. The mean squared error (MSE) loss function was employed as the
345 optimization objective for all models. The MSE is defined as:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (12)$$

346 where y_i and \hat{y}_i represent the measured value and the predicted value; n is the
347 number of batch size which is 24 in this study. This loss function penalizes large
348 prediction errors more heavily, making it particularly effective for minimizing
349 overall prediction deviations. All experiments were conducted on a consistent
350 computational environment to ensure reproducibility. All the code of models
351 and datasets are open sourced in MultiAttLLM GitHub Repository.

352 All models were encoded using PyTorch and trained on a computer with
353 a 14th Intel(R) Core(TM) i9-14900F CPU 3.20 GHz and 64 GB of working
354 memory (RAM). The models were solved and calculated using a GPU (NVIDIA
355 GeForce RTX 4090 24 GB).

Table 3: Hyperparameters configurations for benchmark models and our model.

Model	Hyperparameter	Value
LSTM	Hidden size	256
	Number of layers	2
Informer	Dimensions of model	512
	Dimensions of feed-forward	2048
	Number of encoder layers	4
	Number of decoder layers	2
Autoformer	Dimensions of model	512
	Dimensions of feed-forward	2048
	Number of encoder layers	2
	Number of decoder layers	1
iTTransformer	Dimensions of model	512
	Dimensions of feed-forward	512
	Number of encoder layers	3
	Number of decoder layers	1
TimesNet	Dimensions of model	32
	Dimensions of feed-forward	64
	Number of encoder layers	2
	Number of decoder layers	1
TimeLLM	Dimensions of patch embedding	32
	Text Prototype	1000
	Number of LLM layers	6
MultiAttLLM	Dimensions of model	64
	Dimensions of feed-forward	128
	Number of LLM layers	6
	Number of decoder layers	4

356 **3. Case study**

357 *3.1. Introduction of the dataset*

358 The datasets used for the case study include hourly electricity data (demand
 359 and generation), as well as weather data. The electricity data spans from 2016
 360 to 2023, covering four regions in Japan: Kyushu, Tokyo, Tohoku, and Hokkaido.
 361 For each region, the dataset includes hourly total electricity demand and gen-
 362 eration data from various sources. Electricity data comes from public data of
 363 power companies in various regions. Fig. 6 illustrates the electricity demand
 364 and generation by different energy sources for the four selected regions in Japan
 365 in 2023. As shown in Fig. 6, the Tokyo and Hokkaido regions exhibit similar
 366 patterns, with total generation and demand both around 30 MWh. In these re-
 367 gions, renewable energy generation primarily comes from solar and hydroelectric
 368 sources. In contrast, the Tohoku and Kyushu regions have much higher electric-
 369 ity demand and generation, both reaching approximately 100 MWh. Notably,
 370 in these regions, nuclear energy also constitutes a significant portion of the re-
 371 newable energy mix. This variation in energy profiles across regions reflects the
 372 diversity in energy infrastructure and resource availability, making them ideal
 373 for testing the generalization capability of the proposed model.

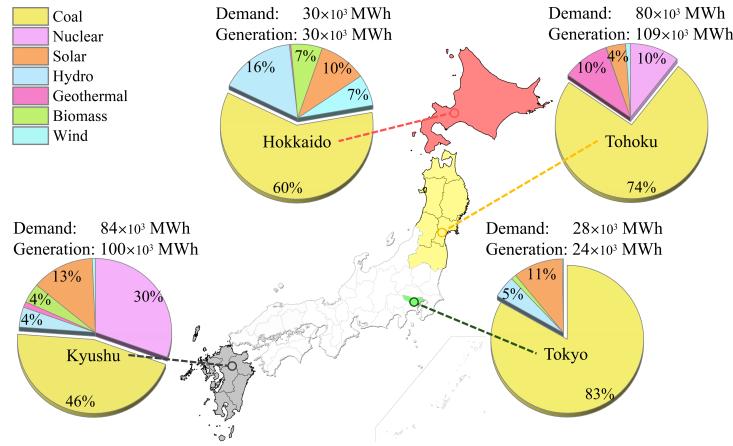


Figure 6: An Overview of electricity demand and generation data for selected regions of Japan in 2023.

374 In addition to electricity demand and generation data, the weather con-
375 ditions in the four selected regions—Fukuoka (Kyushu region), Tokyo (Tokyo
376 region), Sendai (Tohoku region), and Sapporo (Hokkaido region)—play a cru-
377 cial role in influencing both energy consumption and generation from renewable
378 sources. The weather data, sourced from the Japan Meteorological Agency, also
379 covers the period from 2016 to 2023. The weather data from 2023, as shown in
380 Fig. 7, captures key meteorological variables such as maximum and minimum
381 daily temperatures and solar radiation for each region, providing important
382 context for understanding the seasonal variations in energy patterns.

383 Hokkaido, located in the northernmost part of Japan, experiences long, cold
384 winters with significant snowfall and relatively cool summers. Solar radiation
385 is lower compared to other regions. Tohoku, in northeastern Japan, also has
386 cold winters, though less severe than Hokkaido. The region experiences distinct
387 seasons, with cooler temperatures in winter and moderate summers. Tokyo,
388 located in the central part of Japan, has a temperate climate with hot, humid
389 summers and mild winters. The relatively higher solar radiation in Tokyo,
390 particularly during summer months, contributes significantly to solar power
391 generation. Kyushu, in the southernmost part of Japan, enjoys a warm climate
392 year-round with hot, humid summers and mild winters. The region receives
393 the highest solar radiation among the four regions, making solar power a key
394 contributor to its renewable energy mix. Kyushu also benefits from a higher
395 proportion of nuclear energy generation, complementing its renewable energy
396 sources.

397 The details of all features in our datasets are summarized in Table 4. The se-
398 lection of these four regions was made to ensure a comprehensive representation
399 of Japan’s diverse climatic zones, ranging from the warmer southern regions to
400 the cooler northern areas. Additionally, these regions feature varying compo-
401 sitions of electricity generation sources, which allows for a thorough evaluation
402 of the proposed model’s generalization capability across different climatic and
403 energy production conditions.

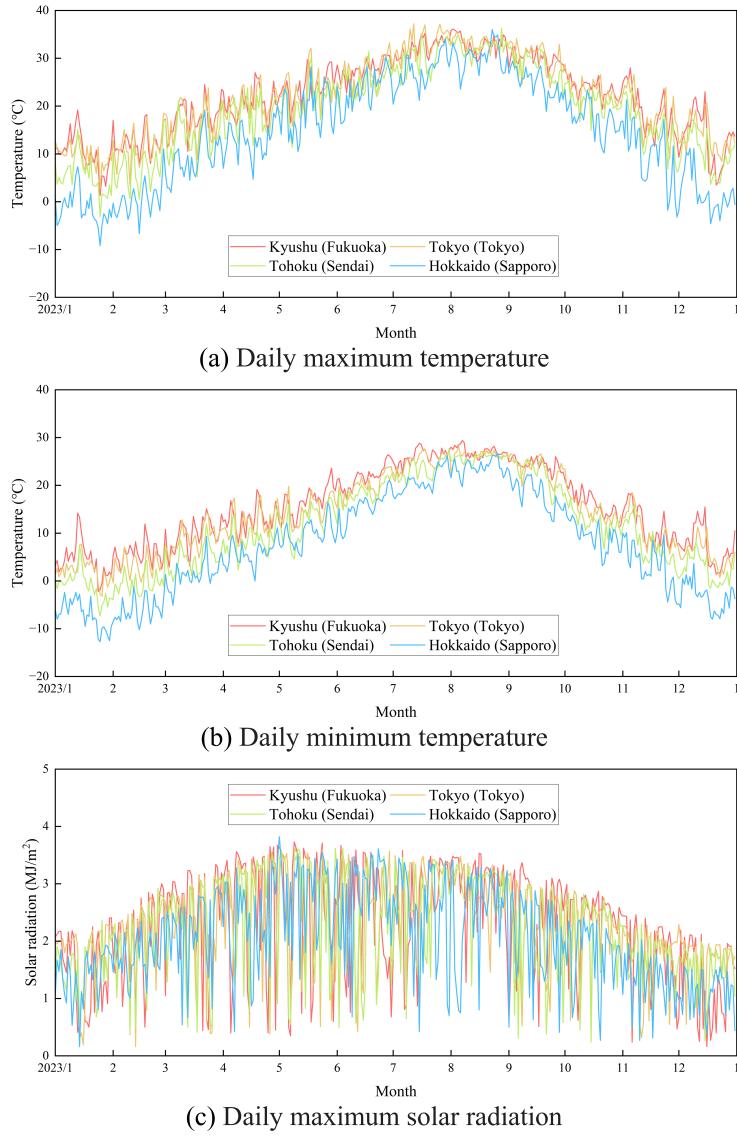


Figure 7: An Overview of daily weather data for selected regions of Japan in 2023.

Table 4: Description of the datasets(*:The value of renewable energy generation is equal to the sum of all other generation except fossil energy generation; data interval: one-hour).

Data	Features
Electricity demand	Total demand
	Fossil
	Nuclear
	Hydro
	Geothermal
Electricity generation	Biomass
	Solar
	Wind
	Renewable energy*
	Temperature
	Relative humidity
	Precipitation
	Dew point
Weather	Vapor pressure
	Wind speed
	Sunshine duration
	Snowfall
	Global horizontal irradiance

⁴⁰⁴ 3.2. Data standardization

⁴⁰⁵ Data standardization adjusts different data ranges to a common scale, which
⁴⁰⁶ helps to minimize regression errors while preserving correlations within the
⁴⁰⁷ dataset [61]. This study utilizes Z-score standardization, which normalizes the
⁴⁰⁸ data to have a mean of zero and a standard deviation of one [62]. The formula
⁴⁰⁹ for Z-score standardization is:

$$x' = \frac{x - \mu}{\delta} \quad (13)$$

⁴¹⁰ where x' and x represent the standardized data and original data, respectively;
⁴¹¹ and μ and δ represent the mean and the standard deviation of the original
⁴¹² data, respectively. Standardization of all feature data is essential before model
⁴¹³ integration. The formula for inverse standardization can be expressed as

$$x = x' \times \delta + \mu \quad (14)$$

⁴¹⁴ 3.3. Data splitting methodology

⁴¹⁵ To validate model effectiveness, we collected eight years of data from 2016
⁴¹⁶ to 2023. First, we assessed how training-set length impacts forecasting accuracy
⁴¹⁷ by training representative models on spans ranging from one to six years. The
⁴¹⁸ results, shown in Fig. 8, indicate that performance stabilizes once the training
⁴¹⁹ period exceeds four years. Balancing predictive gains against computational
⁴²⁰ cost, we therefore selected 2018–2021 as our training dataset. To evaluate gen-
⁴²¹ eralization to the most recent data, 2022 was used as the validation dataset and
⁴²² 2023 as the test dataset, with the checkpoint achieving the lowest validation
⁴²³ loss retained during training.

⁴²⁴ We employed a fixed-window, interval-output forecasting protocol: at each
⁴²⁵ prediction point, the model ingests the entire historical lookback window and
⁴²⁶ simultaneously outputs all future horizon values in a single forward pass. This
⁴²⁷ interval-output approach prevents any overlap between inputs and targets—unlike
⁴²⁸ rolling-window schemes—thereby avoiding data leakage.

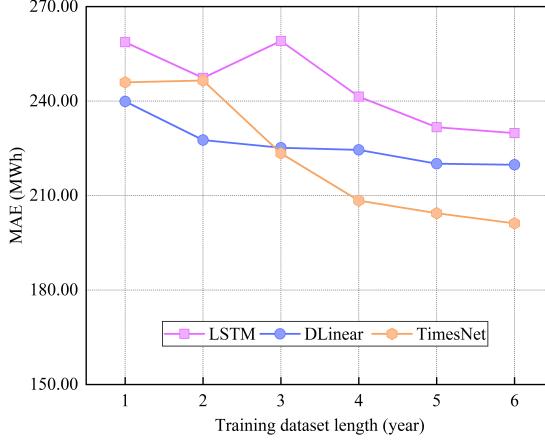


Figure 8: Impact of training dataset length (1–6 years) on MAE for LSTM, DLinear, and TimesNet models (2022 as the validation dataset and 2023 as the test dataset).

429 *3.4. Evaluation metrics*

430 We incorporated three widely used criteria to assess the predictive performance
 431 of the model from multiple perspectives: mean absolute error (MAE),
 432 mean absolute percentage error (MAPE), root mean squared error (RMSE) and
 433 correlation coefficient (R^2). These evaluation metrics enable us to gauge the predictive
 434 ability of the model from various angles. The formulas to calculate MAE,
 435 MAPE, RMSE and R^2 are

$$436 \quad \text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (15)$$

$$437 \quad \text{MAPE} = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \quad (16)$$

$$438 \quad \text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (17)$$

$$439 \quad R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (18)$$

440 where y_i , \hat{y}_i , and \bar{y} represent the measured value, predicted value, and mean of
 441 measured value, respectively; n is the length of sequence. In addition, in order
 to evaluate the versatility of the model on different data sets, we also used a

⁴⁴² dimensionless indicator relative absolute error (RAE), which is calculated as
⁴⁴³ follows:

$$RAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{\sum_{i=1}^n |y_i - \bar{y}|} \quad (19)$$

⁴⁴⁴ As all input data is standardized beforehand, the output of the model repre-
⁴⁴⁵ sents the predicted value after standardization. We de-standardize the predicted
⁴⁴⁶ value output by the model and compute the evaluation metrics based on the
⁴⁴⁷ measured values to more intuitively compare the performance of the model.

⁴⁴⁸ 3.5. Experimental Setup

⁴⁴⁹ Fossil fuel energy and renewable energy play critical roles in shaping energy
⁴⁵⁰ policies and grid management strategies [63]. Accurate forecasting of electricity
⁴⁵¹ demand, as well as generation from fossil fuels and renewable sources, is essential
⁴⁵² for optimizing grid operations, planning energy transitions, and ensuring energy
⁴⁵³ security [64][65]. Given the growing emphasis on carbon emission reduction and
⁴⁵⁴ renewable energy integration, accurate forecasting is essential. Our model aims
⁴⁵⁵ to provide reliable predictions for these key metrics. Therefore, the selected
⁴⁵⁶ targets for prediction include overall electricity demand, fossil fuel-based power
⁴⁵⁷ generation, and renewable energy generation, as shown in Fig. 9. By focusing on
⁴⁵⁸ these aspects, the proposed model can support more informed decision-making
⁴⁵⁹ for energy policy and grid management.

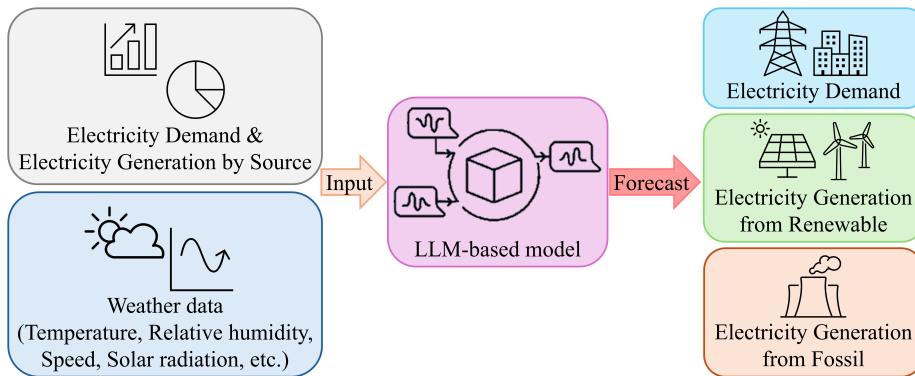


Figure 9: Electricity demand and generation forecasting task.

460 To further validate the generalization capability of the proposed model, we
461 employed a zero-shot learning approach. The model, trained solely on the Tokyo
462 dataset, was tested on the remaining three regions: Hokkaido, Tohoku, and
463 Kyushu. This approach allowed us to evaluate the model’s ability to adapt
464 and provide accurate forecasts without additional training on the new regions’
465 datasets. By comparing the zero-shot performance across these regions, we could
466 assess the robustness and generalization capacity of the proposed model under
467 different climatic and energy generation conditions. In all the above tasks, we
468 fixed the input of historical data for 3 days (72 hours) and predicted the data
469 for the next week (168 hours).

470 To select an appropriate LLM backbone, we evaluated four candidate mod-
471 els—BERT (420 MB) [66], GPT2 (522 MB) [67], LLAMA-3.2-1b (2.30 GB) [68],
472 and LLAMA-3.2-3b (5.97 GB) [68]. Their memory footprints and per-iteration
473 training times are plotted in Fig. 10, and the average MAE across three fore-
474 casting targets was 213, 202, 201, and 198, respectively. While LLAMA-3.2-3b
475 achieves the lowest error, it incurs a very large memory footprint and training
476 latency. GPT2, by contrast, delivers near-state-of-the-art accuracy with the
477 shortest training time, making it the most practical choice on our single-GPU
478 setup. Accordingly, we adopt GPT2 as the LLM component in all subsequent
479 experiments.

480 4. Results and discussions

481 4.1. Ablation study

482 The results of the ablation study, shown in Table 5, demonstrate that each
483 module in the proposed model positively contributes to its predictive perfor-
484 mance. The pre-trained LLM has the most significant impact on the model’s
485 accuracy, as removing this module leads to a 17.3% drop in MAE. The self-
486 attention module also plays a critical role in extracting relationships between
487 non-target and target features, with its removal resulting in a 16.3% decrease
488 in MAE.

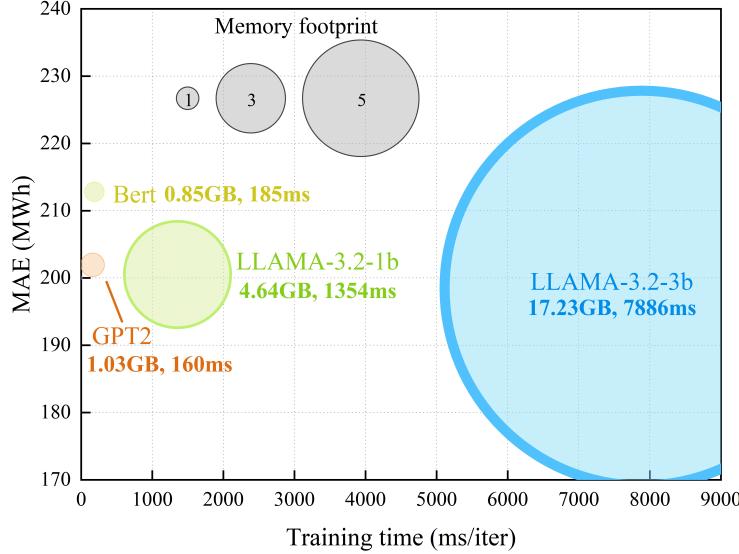


Figure 10: LLM efficiency comparison for proposed model (lookback window length: 72; forecast horizons: 168).

Table 5: Performance of the MultiAttLLM model after removing specific components (- indicates the full model without the corresponding part; * indicates that only the target variable is used as input; ①: Word Projection; ②: Cross-Attention; ③: Pre-trained LLM; ④: Feature Extractor; ⑤: Self-Attention; ⑥: Output Projection).

Descriptions	①	②	③	④	⑤	⑥	MAE
Full model	✓	✓	✓	✓	✓	✓	202 (0.0%)
- Cross-Attention	✗	✗	✓	✓	✓	✓	214 (5.9% ↓)
- Pre-trained LLM	✗	✗	✗	✓	✓	✓	237 (17.3% ↓)
- Feature Extractor	✗	✗	✗	✗	✓	✓	241 (19.3% ↓)
- Feature Extractor	✓	✓	✓	✗	✓	✓	223 (10.4% ↓)
- Self-Attention	✓	✓	✓	✓	✗	✓	235 (16.3% ↓)
Full model (CI)	✓	✓	✓	✓	✓	✓	241 (19.3% ↓)
- Feature Extractor* (CI)	✓	✓	✗	✓	✓	✓	232 (14.9% ↓)

489 Other modules, such as the feature extractor and the word projection, have
490 relatively smaller effects on the model’s performance. Removing the feature ex-
491 tractor decreases the MAE by 10.4%, while removing the word projection mod-
492 ule decreases the MAE by 5.9%. These findings indicate that, while all compo-
493 nents contribute positively to the model’s overall performance, the pre-trained
494 LLM and the self-attention mechanisms are particularly crucial for achieving
495 the model’s high accuracy.

496 Additionally, we also evaluated the model under the CI training strategy.
497 As shown in Table 5 last two row, adopting CI causes a substantial degradation
498 in accuracy: MAE increases by 19.3% compared to our standard setup. When
499 CI is applied with only the target variable as input (i.e., only use pre-trained
500 LLM without feature extractor), MAE still increases by 14.9%. This drop arises
501 because CI computes the loss for a single variable at each iteration. With all
502 variables supplied, the model optimizes to minimize the aggregate loss across
503 every channel, rather than focusing on the intended target, which harms its abil-
504 ity to predict that target accurately. Meantime, when only the target variable
505 is provided, the model cannot leverage cross-channel information, resulting in a
506 further reduction in predictive performance.

507 4.2. Performance under different prompt engineering

508 To evaluate the influence of prompt engineering on model performance, we
509 compared the LLM-based TimeLLM model with the proposed MultiAttLLM
510 model under different prompt conditions. The results, shown in Table 6, include
511 scenarios without a prompt and with three progressively more complex prompts
512 (Prompt1, Prompt2, and Prompt3). A detailed introduction to prompt is shown
513 in Appendix A.2.

514 The results indicate that designing more complex and accurate prompts
515 can improve the prediction accuracy of the TimeLLM model. Specifically, for
516 TimeLLM, using Prompt3 improves the MAE by up to 8.1% compared to the
517 no-prompt condition. However, the proposed MultiAttLLM model shows mini-
518 mal sensitivity to prompt complexity. Regardless of the prompt condition, the

Table 6: Performance of LLM-based models under different prompts (Prompt1: description of the prediction task only; Prompt2: description of the prediction task and dataset; Prompt3: description of the prediction task, dataset, and input sequence statistics; lookback window length: 72; forecast horizons: 168; the unit of MAE and RMSE is MWh, MAPE is %).

Model	TimeLLM				MultiAttLLM			
	Metrics	MAE	MAPE	RMSE	R ²	MAE	MAPE	RMSE
Without prompt	236	21.02	325	0.70	202	15.40	283	0.77
Prompt1	235	21.79	325	0.70	209	16.46	292	0.75
Prompt2	230	19.20	320	0.71	208	16.40	291	0.75
Prompt3	217	16.24	307	0.73	205	16.20	288	0.76

519 MultiAttLLM consistently outperforms TimeLLM, achieving the highest accuracy
 520 (MAE: 202, RMSE: 283, R²: 0.77) even without any prompt.

521 These findings suggest that the proposed MultiAttLLM model can effectively
 522 extract and process feature relationships through its multi-module architecture
 523 without relying on prompt design. This highlights its robustness and efficiency
 524 in utilizing LLMs, making it highly suitable for time series forecasting tasks
 525 without the added complexity of prompt engineering.

526 *4.3. Performance comparison of benchmark models and the proposed MultiAt-
 527 tLLM model*

528 The electricity demand and generation forecasting results for both the bench-
 529 mark models and the proposed model are presented in Table 7. The results
 530 demonstrate that the traditional LSTM model performs the worst across all
 531 prediction targets and metrics. The proposed MultiAttLLM model achieves
 532 the best results for all targets and metrics, outperforming the other models in
 533 terms of MAE, RMSE, and R² scores. The models that improve upon tra-
 534 ditional approaches, Informer, iTransformer, Autoformer and TimesNet, show
 535 similar performance. For some specific variables and metrics, these models can
 536 match the best results achieved by the proposed model.

Table 7: Electricity demand and generation forecasting results (lookback window length: 72; forecast horizons: 168; the unit of MAE and RMSE is MWh, MAPE is %; boldface indicates the best performance; underlined values denote the second-best performance).

Target variable	Metrics	Method							
		LSTM	DLinear	Informer	Autoformer	iTransformer	TimesNet	TimeLLM	MultiAttLLM
Electricity demand	MAE	288	259	235	265	<u>232</u>	251	250	228
	MAPE	8.88	<u>7.77</u>	7.15	8.23	7.01	7.58	7.42	6.89
	RMSE	384	355	317	356	<u>313</u>	337	346	311
	R ²	0.73	0.77	<u>0.80</u>	0.75	0.82	0.79	0.78	0.82
Generation (Renewable energy)	MAE	165	153	<u>135</u>	154	150	148	139	127
	MAPE	48.35	39.25	26.39	43.43	36.46	35.63	29.34	<u>27.47</u>
	RMSE	237	243	<u>222</u>	<u>222</u>	237	235	232	213
	R ²	0.75	0.73	0.73	0.73	0.75	0.75	<u>0.76</u>	0.80
Generation (Fossil energy)	MAE	278	257	271	258	<u>252</u>	251	262	251
	MAPE	13.19	11.97	13.05	11.86	11.85	11.49	11.96	<u>11.83</u>
	RMSE	360	336	343	335	<u>326</u>	332	342	325
	R ²	0.60	<u>0.66</u>	0.62	0.64	0.68	<u>0.66</u>	0.64	0.68
Mean	MAE	244	223	214	226	<u>211</u>	216	217	202
	MAPE	23.47	19.66	<u>15.53</u>	21.17	18.44	18.23	16.24	15.40
	RMSE	327	311	294	304	<u>292</u>	302	307	283
	R ²	0.69	0.72	0.72	0.71	<u>0.75</u>	<u>0.75</u>	0.73	0.77

Furthermore, we also compared each model’s per-iteration training time and memory footprint, as shown in Fig. 11. DLinear exhibits the smallest memory footprint ($\approx 0.02\text{GB}$) and the fastest speed ($\approx 2\text{ms/iter}$), though with only moderate forecasting accuracy. The transformer-based baselines (Informer, Autoformer, iTransformer, TimesNet) cluster around similar MAE values and training times, each outperforming the vanilla LSTM while maintaining comparable efficiency. Under nearly identical parameter counts (TimeLLM: 61.7M; MultiAttLLM: 62.4M), our MultiAttLLM reduces memory usage by about 49.3% (from 2.03GB down to 1.03GB) and cuts per-iteration training time by roughly 35.7% (from 249ms to 160ms). Meanwhile, the proposed model outperforms TimeLLM, with average improvements of 6.9%, 7.8%, and 5.5% in MAE, RMSE, and R², respectively.

4.4. Model performance for different target variables

The relative errors for predicting electricity demand, renewable energy generation, and fossil fuel energy generation across all models are shown in Fig. 12.

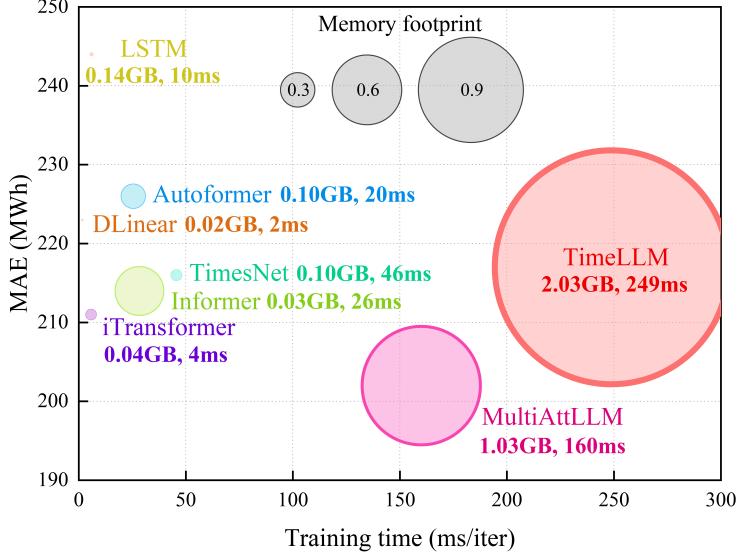


Figure 11: Model efficiency comparison (lookback window length: 72; forecast horizons: 168).

552 The calculation method for relative error (RE) is as follows:

$$\text{RE} = \frac{\hat{y}_i - y_i}{y_i} \quad (20)$$

553 where y_i and \hat{y}_i represent the measured value and predicted value.

554 The results indicate that all models perform best when predicting electricity
 555 demand, followed by fossil fuel energy generation. The prediction accuracy for
 556 renewable energy generation is the lowest. This is primarily because electric-
 557 ity demand and fossil fuel generation in a given region tend to be more stable
 558 and are less influenced by seasonal or weather changes. In contrast, renew-
 559 able energy generation, which includes solar, hydro, and wind power, is highly
 560 dependent on weather and seasonal variations, making it more challenging to
 561 predict accurately.

562 As illustrated in Fig. 13, a portion of the electricity demand forecasting
 563 results shows that the daily variation patterns are fairly consistent across all
 564 days, with peaks occurring during the day and troughs at night. When the peak
 565 demand and fluctuations between consecutive days are similar, all models show
 566 strong predictive performance. However, when there are significant deviations in

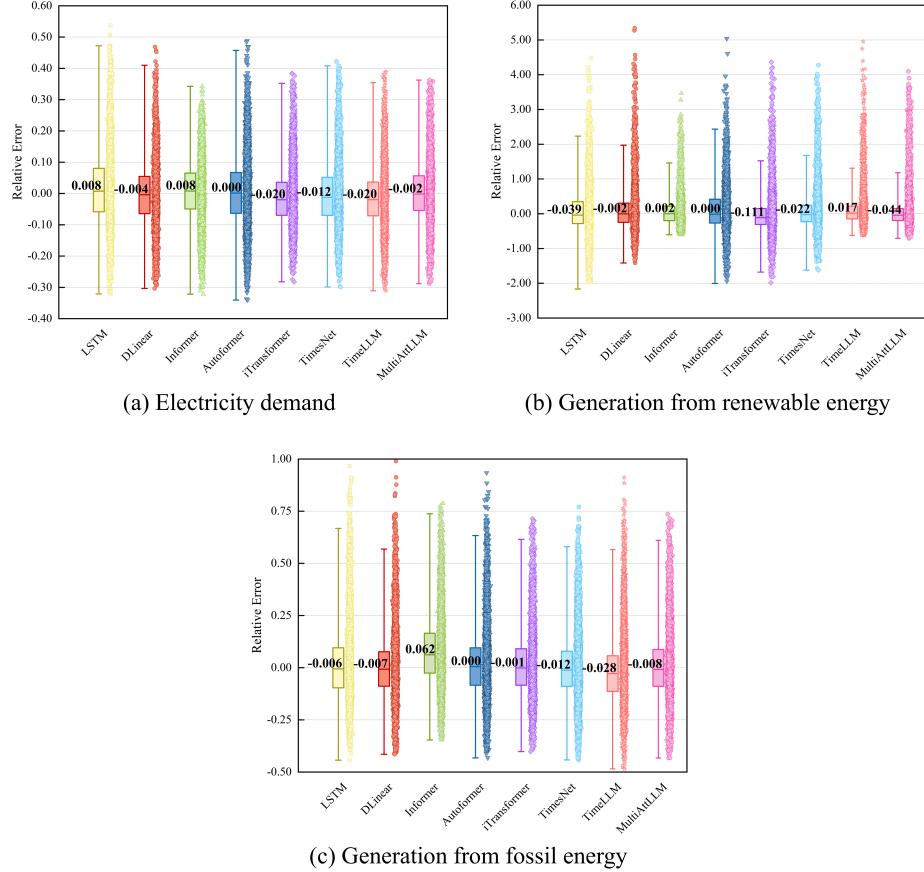


Figure 12: Relative error for electricity demand, renewable energy generation, and fossil energy generation predictions across different models(lookback window length: 72; forecast horizons: 168).

567 peak demand or fluctuations compared to the preceding and following days, the
 568 performance of all models decreases to some extent. Nonetheless, the proposed
 569 model continues to exhibit the best overall performance, particularly in scenarios
 570 with greater variability in demand.

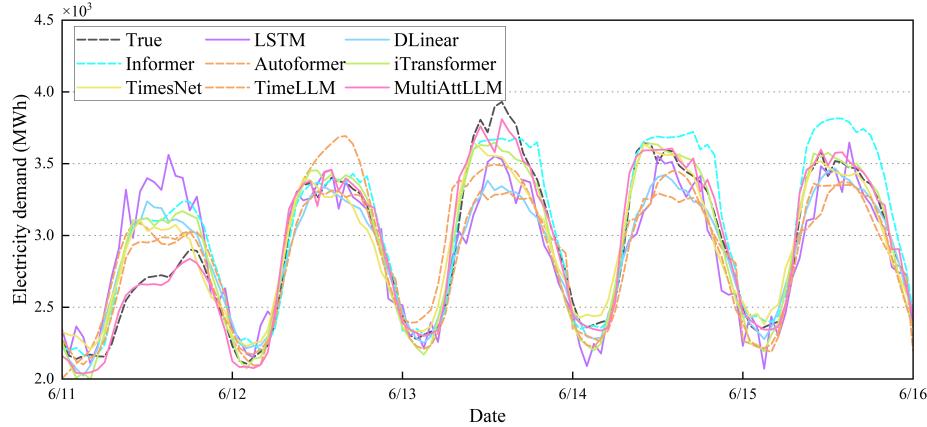


Figure 13: Comparison of electricity demand forecasting results across different models.

571 4.5. Performance across different forecasting horizons

572 To assess how horizon length impacts predictive accuracy, we evaluated every
 573 model on horizons ranging from 1 step (1 h) to 720 steps (30 days) for electricity
 574 demand forecasting, with the average results of three target variables plotted
 575 in Fig. 14. When forecasting just one step ahead, all models achieve very
 576 high accuracy due to the strong correlation between immediate history and
 577 the next value. As the horizon extends from 1 to 24 steps, performance for
 578 every method declines sharply. Beyond 24 steps, however, the rate of accuracy
 579 degradation slows markedly—longer-term targets bear weaker links to input
 580 history, so error growth tapers off. Among the baselines, Informer and TimesNet
 581 show the slowest drop in precision over long horizons, but our MultiAttLLM
 582 model consistently maintains the highest accuracy across both short-term and
 583 long-term forecasts.

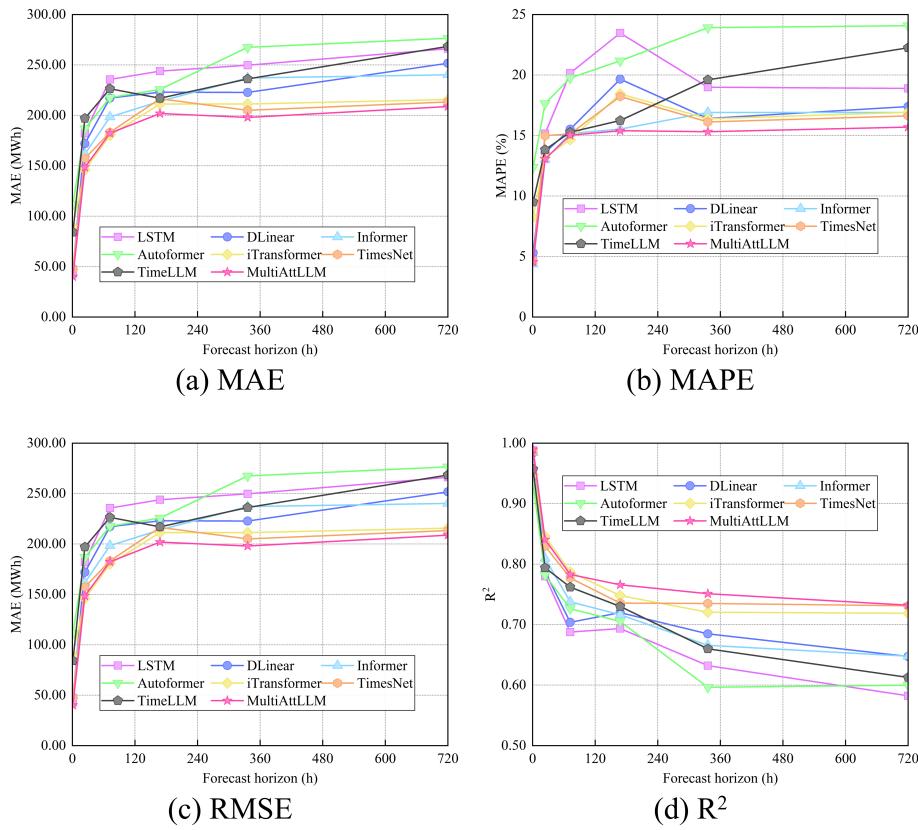


Figure 14: Performance across different forecasting horizons (from 1 steps to 720 steps)

584 *4.6. Generalization performance through zero-shot learning*

585 To evaluate the generalization capability of the proposed model, we con-
 586 ducted zero-shot learning experiments. For each experiment, data from one of
 587 the four regions—Tokyo, Hokkaido, Tohoku, or Kyushu—was used as the source
 588 domain to train the model. The trained model was then tested on all four tar-
 589 get domains without any fine-tuning or adjustments. The average results for
 590 each target domain are presented in Table 8. These results demonstrate that
 591 the proposed MultiAttLLM model consistently achieves the best generaliza-
 592 tion performance across all target domains. The TimeLLM model achieves the
 593 second-best accuracy, highlighting the superior generalization ability of LLM-
 594 based models. The DLinear and TimesNet model, which leverages time series
 595 decomposition, demonstrates a moderate level of generalization. In contrast,
 596 the self-attention based models, which perform reasonably well in non-transfer
 597 learning settings, show significant drops in performance when tested on unseen
 598 datasets. Detailed results for using each region as the source domain and the
 599 corresponding predictions on different target domains are provided in Appendix
 600 A.3.

Table 8: Zero-shot learning performance across target domains using different source regions
 (Results averaged across target domains for each source region; lookback window length: 72;
 forecast horizons: 168; the unit of MAE and RMSE is MWh).

Target domain	Tokyo				Hokkaido				Tohoku				Kyushu			
	Metrics	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE
LSTM	253	351	0.52	0.64	250	330	0.48	0.74	1097	1426	0.57	0.64	880	1210	0.52	0.66
DLinear	227	316	0.47	0.71	241	318	0.46	0.75	928	1245	0.48	0.73	824	1159	0.49	0.69
Informer	250	341	0.53	0.63	307	395	0.62	0.57	928	1240	0.54	0.64	906	1218	0.56	0.62
Autoformer	276	367	0.59	0.58	277	359	0.56	0.65	1018	1339	0.60	0.58	928	1238	0.58	0.60
iTransformer	288	380	0.60	0.56	352	474	0.60	0.57	1164	1528	0.59	0.60	1133	1501	0.67	0.48
TimesNet	237	328	0.49	0.69	244	323	0.46	0.75	1201	1576	0.62	0.54	841	1170	0.50	0.68
TimeLLM	229	318	0.47	0.71	238	314	0.45	0.76	943	1261	0.49	0.72	849	1185	0.51	0.68
MultiAttLLM	211	304	0.43	0.73	234	310	0.44	0.77	891	1220	0.46	0.74	796	1131	0.47	0.71

601 When using the Tokyo dataset as the source domain, the zero-shot learning
 602 results for Hokkaido, Tohoku, and Kyushu are shown in Fig. 15. The dimension-
 603 less metrics reveal that transformer-based models and the LSTM model exhibit
 604 the poorest generalization capability. While these models maintain reasonable

accuracy on the Tohoku dataset, which shares geographic and climatic similarities with Tokyo, their performance significantly deteriorates on the Hokkaido and Kyushu datasets, which feature distinct climatic and energy generation characteristics. In contrast, the proposed MultiAttLLM model demonstrates the best generalization performance across all test regions, achieving the lowest relative errors and maintaining robust predictions even under varying conditions.

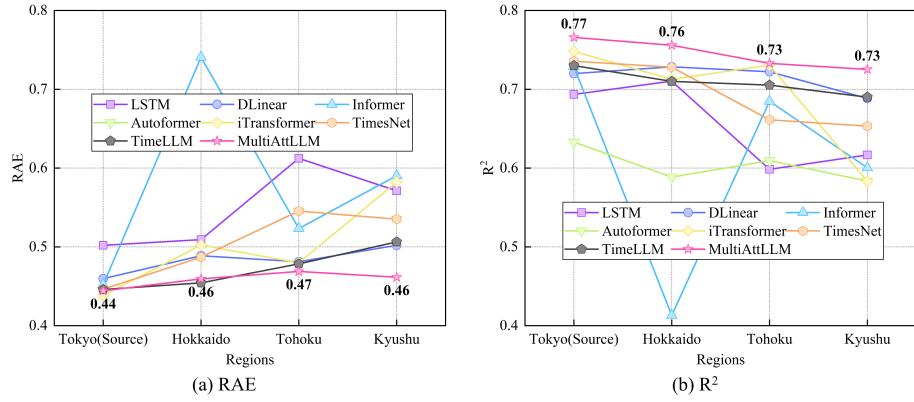


Figure 15: Zero-shot learning results of two dimensionless metrics (Source domain: Tokyo; Target domains: Hokkaido, Tohoku, Kyushu; lookback window length: 72; forecast horizons: 168).

5. Conclusion

Large language models (LLMs) have garnered significant attention in time series forecasting due to their strong performance across diverse tasks. However, their application to multivariate forecasting remains challenged by channel-independent training strategies that impede learning of cross-channel interactions and by the complexity of prompt engineering. To overcome these obstacles, we introduced a novel multi-attention LLM framework that integrates cross-attention, self-attention, and a frozen pre-trained LLM to extract rich interdependencies among target and auxiliary variables without any manual prompt design.

622 Extensive ablation studies demonstrate that every component of our ar-
623 chitecture contributes positively to forecasting accuracy, with the pre-trained
624 LLM backbone and the attention modules yielding the largest gains. In bench-
625 mark comparisons, our model outperforms state-of-the-art transformer-based
626 and LLM-based approaches, reducing mean absolute error by up to 20.8% rel-
627 ative to standard LSTM models. Compared with the current best time-series
628 LLM, the proposed model reduces memory usage by 49.3% and shortens train-
629 ing time by 35.7%. Moreover, it maintains superior trend-following capabil-
630 ities under volatile conditions, achieving the highest accuracy across electric-
631 ity demand, renewable generation, and fossil generation forecasts. Crucially,
632 in zero-shot evaluations on four geographically diverse Japanese regions, our
633 framework improves MAE by an average of 16.6% over LSTM, demonstrating
634 exceptional generalization ability. These findings underscore the effectiveness
635 of multi-attention reprogramming in unlocking the full potential of LLMs for
636 accurate, robust, and generalizable multivariate time series forecasting.

637 However, this study has some limitations. First, the model was not fine-
638 tuned for specific time series forecasting tasks, which could further improve
639 its performance. Additionally, the computational complexity of LLMs remains
640 a challenge, particularly for real-time applications. In future work, we aim
641 to explore further improvements by fine-tuning the LLMs module for specific
642 time series forecasting tasks and incorporating more complex data sources, such
643 as real-time data streams. Additionally, extending the model to handle multi-
644 horizon forecasting and integrating domain-specific knowledge into the attention
645 mechanisms could further enhance the model’s accuracy and applicability in
646 various industries.

647 **Acknowledgment**

648 This work was supported by project No. 23KJ0766 funded by the Japan
649 Society for The Promotion of Science.

650 **Appendix A.1 Long short-term memory network**

651 The core functionality of the long short-term memory (LSTM) unit lies in
 652 its three gating mechanisms, which are controlled by sigmoid functions. These
 653 gates regulate the proportion of past and current input information used in
 654 the unit's calculations, ultimately determining the output for the current time
 655 step through the output gate (denoted as $\mathbf{h}^{(t)}$ in Fig. 1). The cell state, $\mathbf{c}^{(t)}$,
 656 remains largely unaffected by the gates and stays relatively stable throughout
 657 the computation process. This "conveyor belt" system plays a critical role in
 658 maintaining the long-term memory capabilities of the LSTM network.

659 LSTM is calculated using:

$$\mathbf{f}^{(t)} = \sigma(\mathbf{W}_f[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}_f) \quad (1)$$

$$\mathbf{i}^{(t)} = \sigma(\mathbf{W}_i[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}^{(t)}) \quad (2)$$

$$\tilde{\mathbf{c}}^{(t)} = \tanh(\mathbf{W}_c[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}_c) \quad (3)$$

$$\mathbf{c}^{(t)} = \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \tilde{\mathbf{c}}^{(t)} \quad (4)$$

$$\mathbf{o}^{(t)} = \sigma(\mathbf{W}_o[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}_o) \quad (5)$$

$$\mathbf{h}^{(t)} = \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)}) \quad (6)$$

665 where $\mathbf{f}^{(t)}$, $\mathbf{i}^{(t)}$, and $\mathbf{o}^{(t)}$ represents the calculation results of the forget gate,
 666 input gate, and output gate, respectively; $\tilde{\mathbf{c}}^{(t)}$ is the cell state update; $\mathbf{h}^{(t)}$ is
 667 the hidden state; \mathbf{W} , \mathbf{U} , and \mathbf{b} refer to the weights and biases in the model; \odot
 668 is the Hadamard product; σ is the sigmoid function.

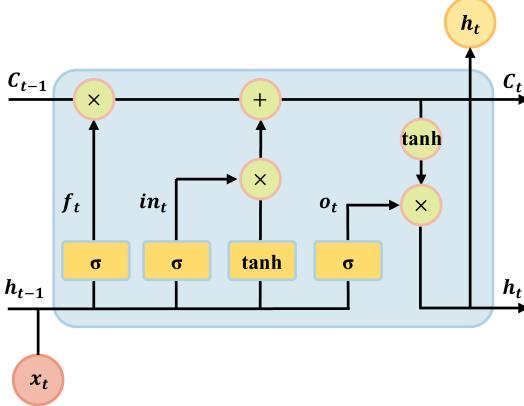


Figure 1: Unfold structure of the LSTM unit.

669 Appendix A.2 Prompt engineering

670 Based on the prompt engineering described in [36], we designed three types
 671 of prompts tailored for time series forecasting tasks. These prompts are detailed
 672 as follows ($\{\}$ is task-specific configurations or calculated input statistics):

- 673 1) Prompt1 (prediction task description only): $|-\text{start_prompt}-|$ Task de-
 674 scription: forecast the next $\{\text{forecast sequence length}\}$ steps given the
 675 previous $\{\text{input sequence length}\}$ steps information $|-\text{end_prompt}-|$
- 676 2) Prompt2 (prediction task and dataset description): $|-\text{start_prompt}-|$
 677 Task description: forecast the next $\{\text{forecast sequence length}\}$ steps given
 678 the previous $\{\text{input sequence length}\}$ steps information; Dataset descrip-
 679 tion: Electricity dataset is recorded every 1 hour, which contains meteo-
 680 rological indicators and power source composition of Japan , such as air
 681 temperature, humidity, solar, coal nuclear, etc. $|-\text{end_prompt}-|$
- 682 3) Prompt3 (prediction task, dataset description, and input sequence statis-
 683 tics): $|-\text{start_prompt}-|$ Task description: forecast the next $\{\text{forecast}$
 $\text{sequence length}\}$ steps given the previous $\{\text{input sequence length}\}$ steps
 685 information; Dataset description: Electricity dataset is recorded every 1

hour, which contains meteorological indicators and power source composition of Japan, such as air temperature, humidity, solar, coal nuclear, etc.;
 Input statistics: min value {min values}, max value {max values}, median value {median values}, the trend of input is {upward or downward}, top
 {top_k}, lags are : {lags values} i—end_prompt—i

691 Appendix A.3 Results of zero-shot learning

Table 9: Zero-shot learning performance across target domains using Kyushu region as source domain (The unit of MAE and RMSE is MWh).

Target domain	Tokyo				Hokkaido				Tohoku				Kyushu			
	Metrics	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE
LSTM	268	379	0.55	0.59	264	349	0.50	0.71	1036	1345	0.54	0.68	804	1126	0.48	0.71
DLinear	231	329	0.48	0.68	248	327	0.47	0.74	917	1232	0.48	0.73	784	1109	0.47	0.72
Informer	238	321	0.51	0.68	302	393	0.61	0.58	882	1203	0.52	0.67	718	985	0.45	0.75
Autoformer	284	383	0.61	0.55	274	355	0.55	0.65	1053	1400	0.63	0.54	844	1152	0.53	0.66
iTransformer	282	375	0.59	0.58	314	418	0.53	0.67	1114	1470	0.57	0.63	1056	1417	0.62	0.54
TimesNet	235	328	0.49	0.68	258	350	0.49	0.70	1554	2054	0.79	0.26	761	1072	0.45	0.74
TimeLLM	229	318	0.47	0.71	239	315	0.46	0.76	934	1251	0.49	0.72	849	1184	0.50	0.68
MultiAttLLM	211	304	0.43	0.73	237	314	0.45	0.76	891	1223	0.46	0.74	796	1130	0.47	0.71

Table 10: Zero-shot learning performance across target domains using Hokkaido region as source domain (The unit of MAE and RMSE is MWh).

Target domain	Tokyo				Hokkaido				Tohoku				Kyushu			
	Metrics	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE
LSTM	261	359	0.54	0.63	238	317	0.45	0.76	1177	1521	0.61	0.59	907	1237	0.54	0.65
DLinear	234	321	0.48	0.70	226	299	0.43	0.78	938	1256	0.49	0.72	808	1136	0.48	0.70
Informer	285	380	0.61	0.56	266	343	0.54	0.68	1182	1560	0.70	0.44	1054	1398	0.65	0.51
Autoformer	307	405	0.65	0.50	261	341	0.53	0.68	1073	1408	0.63	0.54	958	1273	0.60	0.58
iTransformer	295	387	0.62	0.54	371	503	0.63	0.52	1187	1555	0.59	0.59	1159	1533	0.68	0.45
TimesNet	255	349	0.52	0.66	231	305	0.44	0.77	1053	1384	0.55	0.66	862	1196	0.51	0.67
TimeLLM	231	321	0.47	0.70	236	312	0.45	0.76	950	1272	0.49	0.71	853	1192	0.51	0.67
MultiAttLLM	210	303	0.43	0.73	226	302	0.43	0.78	888	1216	0.46	0.74	791	1130	0.47	0.71

692 References

- [1] Y. Oswald, A. Owen, J. K. Steinberger, Large inequality in international and intranational energy footprints between income groups and across consumption categories, Nature Energy 5 (3) (2020) 231–239.

Table 11: Zero-shot learning performance across target domains using Tokyo region as source domain (The unit of MAE and RMSE is MWh).

Target domain	Tokyo				Hokkaido				Tohoku				Kyushu			
	Metrics	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE
LSTM	222	305	0.46	0.73	261	338	0.50	0.72	999	1318	0.52	0.69	902	1241	0.54	0.64
DLinear	208	293	0.43	0.75	264	347	0.50	0.71	919	1236	0.47	0.73	894	1258	0.53	0.63
Informer	208	288	0.45	0.73	370	470	0.74	0.41	898	1190	0.52	0.68	949	1252	0.59	0.60
Autoformer	256	341	0.55	0.63	307	391	0.62	0.59	1035	1335	0.61	0.59	975	1270	0.61	0.58
iTransformer	280	370	0.60	0.57	354	473	0.60	0.57	1168	1534	0.59	0.60	1160	1523	0.68	0.47
TimesNet	203	287	0.42	0.76	254	331	0.48	0.73	1142	1483	0.60	0.59	880	1216	0.52	0.66
TimeLLM	226	313	0.46	0.72	240	315	0.46	0.76	938	1250	0.49	0.72	843	1171	0.50	0.68
MultiAttLLM	202	283	0.42	0.77	246	323	0.46	0.76	898	1227	0.47	0.73	806	1135	0.46	0.73

Table 12: Zero-shot learning performance across target domains using Tohoku region as source domain (The unit of MAE and RMSE is MWh).

Target domain	Tokyo				Hokkaido				Tohoku				Kyushu			
	Metrics	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE	R ²	MAE	RMSE	RAE
LSTM	289	375	0.60	0.59	287	369	0.55	0.67	961	1267	0.50	0.72	1255	1594	0.73	0.42
DLinear	222	304	0.46	0.73	236	308	0.45	0.77	892	1195	0.46	0.75	828	1177	0.49	0.68
Informer	267	374	0.56	0.57	290	374	0.59	0.62	748	1009	0.44	0.77	902	1237	0.57	0.60
Autoformer	256	340	0.55	0.64	265	349	0.54	0.67	909	1215	0.53	0.66	934	1257	0.58	0.60
iTransformer	292	384	0.61	0.55	360	482	0.61	0.56	1168	1535	0.59	0.60	1159	1529	0.68	0.46
TimesNet	543	696	1.08	-0.33	396	515	0.75	0.36	928	1262	0.48	0.72	1489	1981	0.89	0.06
TimeLLM	234	325	0.48	0.69	238	315	0.45	0.76	961	1284	0.50	0.71	873	1211	0.52	0.67
MultiAttLLM	210	304	0.43	0.73	230	307	0.44	0.77	886	1215	0.46	0.74	791	1129	0.47	0.71

- 696 [2] C. Dingbang, C. Cang, C. Qing, S. Lili, C. Caiyun, Does new energy con-
697 sumption conducive to controlling fossil energy consumption and carbon
698 emissions?-evidence from china, Resources policy 74 (2021) 102427.
- 699 [3] R. Khalili, A. Khaledi, M. Marzband, A. F. Nematollahi, B. Vahidi,
700 P. Siano, Robust multi-objective optimization for the iranian electricity
701 market considering green hydrogen and analyzing the performance of dif-
702 ferent demand response programs, Applied Energy 334 (2023) 120737.
- 703 [4] M. S. S. Danish, T. Senju, T. Funabashia, M. Ahmadi, A. M. Ibrahimi,
704 R. Ohta, H. O. R. Howlader, H. Zaheb, N. R. Sabory, M. M. Sediqi, A sus-
705 tainable microgrid: A sustainability and management-oriented approach,
706 Energy Procedia 159 (2019) 160–167.
- 707 [5] C.-T. Hsiao, C.-S. Liu, D.-S. Chang, C.-C. Chen, Dynamic modeling of the
708 policy effect and development of electric power systems: A case in taiwan,
709 Energy policy 122 (2018) 377–387.
- 710 [6] M. Ghiasi, T. Niknam, Z. Wang, M. Mehrandezh, M. Dehghani,
711 N. Ghadimi, A comprehensive review of cyber-attacks and defense mecha-
712 nisms for improving security in smart grid energy systems: Past, present
713 and future, Electric Power Systems Research 215 (2023) 108975.
- 714 [7] T. Hong, P. Pinson, Y. Wang, R. Weron, D. Yang, H. Zareipour, Energy
715 forecasting: A review and outlook, IEEE Open Access Journal of Power
716 and Energy 7 (2020) 376–388.
- 717 [8] P. Akbary, M. Ghiasi, M. R. R. Pourkheranjani, H. Alipour, N. Ghadimi,
718 Extracting appropriate nodal marginal prices for all types of committed
719 reserve, Computational Economics 53 (2019) 1–26.
- 720 [9] M. Sharma, N. Mittal, A. Mishra, A. Gupta, Survey of electricity demand
721 forecasting and demand side management techniques in different sectors to
722 identify scope for improvement, Smart Grids and Sustainable Energy 8 (2)
723 (2023) 9.

- 724 [10] M. Ghiasi, Z. Wang, M. Mehrandezh, S. Jalilian, N. Ghadimi, Evolution of
725 smart grids towards the internet of energy: Concept and essential compo-
726 nents for deep decarbonisation, *IET Smart Grid* 6 (1) (2023) 86–102.
- 727 [11] A. O. Aderibigbe, E. C. Ani, P. E. Ohenehen, N. C. Ohalete, D. O. Daraor-
728 jimba, Enhancing energy efficiency with ai: a review of machine learning
729 models in electricity demand forecasting, *Engineering Science & Technol-*
730 *ogy Journal* 4 (6) (2023) 341–356.
- 731 [12] A. Román-Portabales, M. López-Nores, J. J. Pazos-Arias, Systematic re-
732 view of electricity demand forecast using ann-based machine learning algo-
733 rithms, *Sensors* 21 (13) (2021) 4544.
- 734 [13] N. Sultana, S. Z. Hossain, S. H. Almuhamini, D. Düştégör, Bayesian opti-
735 mization algorithm-based statistical and machine learning approaches for
736 forecasting short-term electricity demand, *Energies* 15 (9) (2022) 3425.
- 737 [14] C. E. Velasquez, M. Zocatelli, F. B. Estanislau, V. F. Castro, Analysis of
738 time series models for brazilian electricity demand forecasting, *Energy* 247
739 (2022) 123483.
- 740 [15] W. Jiang, X. Wang, H. Huang, D. Zhang, N. Ghadimi, Optimal economic
741 scheduling of microgrids considering renewable energy sources based on
742 energy hub model using demand response and improved water wave opti-
743 mization algoritm, *Journal of Energy Storage* 55 (2022) 105311.
- 744 [16] J. F. Torres, F. Martínez-Álvarez, A. Troncoso, A deep lstm network for the
745 spanish electricity consumption forecasting, *Neural Computing and Appli-*
746 *cations* 34 (13) (2022) 10533–10545.
- 747 [17] H. Wu, T. Hu, Y. Liu, H. Zhou, J. Wang, M. Long, Timesnet: Temporal
748 2d-variation modeling for general time series analysis, arXiv preprint
749 arXiv:2210.02186 (2022).
- 750 [18] H. Zhou, S. Zhang, J. Peng, S. Zhang, J. Li, H. Xiong, W. Zhang, Informer:
751 Beyond efficient transformer for long sequence time-series forecasting, in:

- 752 Proceedings of the AAAI conference on artificial intelligence, Vol. 35, 2021,
753 pp. 11106–11115.
- 754 [19] H. Wu, J. Xu, J. Wang, M. Long, Autoformer: Decomposition transformers
755 with auto-correlation for long-term series forecasting, *Advances in neural*
756 information processing systems
- 34 (2021) 22419–22430.
- 757 [20] M. E. Günay, Forecasting annual gross electricity demand by artificial neu-
758 ral networks using predicted values of socio-economic indicators and cli-
759 matic conditions: Case of turkey, *Energy Policy* 90 (2016) 92–101.
- 760 [21] H. Iftikhar, J. E. Turpo-Chaparro, P. Canas Rodrigues, J. L. López-
761 Gonzales, Day-ahead electricity demand forecasting using a novel decom-
762 position combination method, *Energies* 16 (18) (2023) 6675.
- 763 [22] F. Pallonetto, C. Jin, E. Mangina, Forecast electricity demand in com-
764 mercial building with machine learning models to enable demand response
765 programs, *Energy and AI* 7 (2022) 100121.
- 766 [23] T. G. Grandón, J. Schwenzer, T. Steens, J. Breuing, Electricity demand
767 forecasting with hybrid classical statistical and machine learning algo-
768 rithms: Case study of ukraine, *Applied Energy* 355 (2024) 122249.
- 769 [24] E. Cebekhulu, A. J. Onumanyi, S. J. Isaac, Performance analysis of machine
770 learning algorithms for energy demand-supply prediction in smart grids,
771 *Sustainability* 14 (5) (2022) 2546.
- 772 [25] Z. Wang, Z. Chen, Y. Yang, C. Liu, X. Li, J. Wu, A hybrid autoformer
773 framework for electricity demand forecasting, *Energy Reports* 9 (2023)
774 3800–3812.
- 775 [26] C. Wu, J. Li, W. Liu, Y. He, S. Nourmohammadi, Short-term electricity
776 demand forecasting using a hybrid anfis–elm network optimised by an im-
777 proved parasitism–predation algorithm, *Applied Energy* 345 (2023) 121316.

- 778 [27] C. Sekhar, R. Dahiya, Robust framework based on hybrid deep learning
779 approach for short term load forecasting of building electricity demand,
780 Energy 268 (2023) 126660.
- 781 [28] E. C. May, A. Bassam, L. J. Ricalde, M. E. Soberanis, O. Oubram, O. M.
782 Tzuc, A. Y. Alanis, A. Livas-García, Global sensitivity analysis for a real-
783 time electricity market forecast by a machine learning approach: A case
784 study of mexico, International Journal of Electrical Power & Energy Sys-
785 tems 135 (2022) 107505.
- 786 [29] Y. Jiang, T. Gao, Y. Dai, R. Si, J. Hao, J. Zhang, D. W. Gao, Very short-
787 term residential load forecasting based on deep-autoformer, Applied Energy
788 328 (2022) 120120.
- 789 [30] M. Alhussein, K. Aurangzeb, S. I. Haider, Hybrid cnn-lstm model for short-
790 term individual household load forecasting, Ieee Access 8 (2020) 180544–
791 180557.
- 792 [31] R.-x. Nie, Z.-p. Tian, R.-y. Long, W. Dong, Forecasting household electric-
793 ity demand with hybrid machine learning-based methods: Effects of res-
794 idents' psychological preferences and calendar variables, Expert Systems
795 with Applications 206 (2022) 117854.
- 796 [32] T. B. Brown, Language models are few-shot learners, arXiv preprint
797 arXiv:2005.14165 (2020).
- 798 [33] A. Vaswani, Attention is all you need, Advances in Neural Information
799 Processing Systems (2017).
- 800 [34] M. U. Hadi, R. Qureshi, A. Shah, M. Irfan, A. Zafar, M. B. Shaikh,
801 N. Akhtar, J. Wu, S. Mirjalili, et al., A survey on large language mod-
802 els: Applications, challenges, limitations, and practical usage, Authorea
803 Preprints (2023).

- 804 [35] J. Kaddour, J. Harris, M. Mozes, H. Bradley, R. Raileanu, R. McHardy,
805 Challenges and applications of large language models, arXiv preprint
806 arXiv:2307.10169 (2023).
- 807 [36] M. Jin, S. Wang, L. Ma, Z. Chu, J. Y. Zhang, X. Shi, P.-Y. Chen, Y. Liang,
808 Y.-F. Li, S. Pan, et al., Time-lm: Time series forecasting by reprogram-
809 ming large language models, arXiv preprint arXiv:2310.01728 (2023).
- 810 [37] A. Nazir, A. K. Shaikh, A. S. Shah, A. Khalil, Forecasting energy consump-
811 tion demand of customers in smart grid using temporal fusion transformer
812 (tft), Results in Engineering 17 (2023) 100888.
- 813 [38] J. Su, C. Jiang, X. Jin, Y. Qiao, T. Xiao, H. Ma, R. Wei, Z. Jing, J. Xu,
814 J. Lin, Large language models for forecasting and anomaly detection: A
815 systematic literature review, arXiv preprint arXiv:2402.10350 (2024).
- 816 [39] N. Gruver, M. Finzi, S. Qiu, A. G. Wilson, Large language models are zero-
817 shot time series forecasters, Advances in Neural Information Processing
818 Systems 36 (2023) 19622–19635.
- 819 [40] H. Xue, F. D. Salim, Promptcast: A new prompt-based learning paradigm
820 for time series forecasting, IEEE Transactions on Knowledge and Data
821 Engineering 36 (11) (2023) 6851–6864.
- 822 [41] D. Cao, F. Jia, S. O. Arik, T. Pfister, Y. Zheng, W. Ye, Y. Liu, Tempo:
823 Prompt-based generative pre-trained transformer for time series forecast-
824 ing, arXiv preprint arXiv:2310.04948 (2023).
- 825 [42] M. Tan, M. Merrill, V. Gupta, T. Althoff, T. Hartvigsen, Are language
826 models actually useful for time series forecasting?, Advances in Neural In-
827 formation Processing Systems 37 (2024) 60162–60191.
- 828 [43] T. Zhou, P. Niu, L. Sun, R. Jin, et al., One fits all: Power general time
829 series analysis by pretrained lm, Advances in neural information processing
830 systems 36 (2023) 43322–43355.

- 831 [44] C. Chang, W.-Y. Wang, W.-C. Peng, T.-F. Chen, Llm4ts: Aligning
832 pre-trained llms as data-efficient time-series forecasters, arXiv preprint
833 arXiv:2308.08469 (2023).
- 834 [45] C. Sun, H. Li, Y. Li, S. Hong, Test: Text prototype aligned embedding
835 to activate llm's ability for time series, arXiv preprint arXiv:2308.08241
836 (2023).
- 837 [46] H. Xue, F. D. Salim, Utilizing language models for energy load forecasting,
838 in: Proceedings of the 10th ACM International Conference on Systems for
839 Energy-Efficient Buildings, Cities, and Transportation, 2023, pp. 224–227.
- 840 [47] T. Wu, Q. Ling, Stellm: Spatio-temporal enhanced pre-trained large lan-
841 guage model for wind speed forecasting, Applied Energy 375 (2024) 124034.
- 842 [48] D. Han, W. Guo, H. Chen, B. Wang, Z. Guo, Lest: Large language models
843 and spatio-temporal data analysis for enhanced sino-us exchange rate fore-
844 casting, International Review of Economics & Finance 96 (2024) 103508.
- 845 [49] Z. Lai, T. Wu, X. Fei, Q. Ling, Bert4st:: Fine-tuning pre-trained large
846 language model for wind power forecasting, Energy Conversion and Man-
847 agement 307 (2024) 118331.
- 848 [50] X. Yu, Z. Chen, Y. Ling, S. Dong, Z. Liu, Y. Lu, Temporal data meets
849 llm-explainable financial time series forecasting. arxiv 2023, arXiv preprint
850 arXiv:2306.11025.
- 851 [51] G. Liu, Y. Bai, K. Wen, X. Wang, Y. Liu, G. Liang, J. Zhao, Z. Y. Dong,
852 Lflm: A large language model for load forecasting, Authorea Preprints
853 (2024).
- 854 [52] Y. Wang, H. A. Karimi, Exploring large language models for climate fore-
855 casting, arXiv preprint arXiv:2411.13724 (2024).
- 856 [53] Z. Duan, C. Bian, S. Yang, C. Li, Prompting large language model for
857 multi-location multi-step zero-shot wind power forecasting, Expert Systems
858 with Applications (2025) 127436.

- 859 [54] W. Wang, Y. Luo, M. Ma, J. Wang, C. Sui, A novel forecasting framework
860 leveraging large language model and machine learning for methanol price,
861 Energy 320 (2025) 135123.
- 862 [55] T. Guo, E. Hauptmann, Fine-tuning large language models for stock return
863 prediction using newsflow, arXiv preprint arXiv:2407.18103 (2024).
- 864 [56] E. Spiliotis, Time series forecasting with statistical, machine learning, and
865 deep learning methods: Past, present, and future, in: Forecasting with
866 Artificial Intelligence: Theory and Applications, Springer, 2023, pp. 49–75.
- 867 [57] L. Han, H.-J. Ye, D.-C. Zhan, The capacity and robustness trade-off: Re-
868 visiting the channel independent strategy for multivariate time series fore-
869 casting, IEEE Transactions on Knowledge and Data Engineering (2024).
- 870 [58] S. Hochreiter, Long short-term memory, Neural Computation MIT-Press
871 (1997).
- 872 [59] A. Zeng, M. Chen, L. Zhang, Q. Xu, Are transformers effective for time
873 series forecasting?, in: Proceedings of the AAAI conference on artificial
874 intelligence, Vol. 37, 2023, pp. 11121–11128.
- 875 [60] Y. Liu, T. Hu, H. Zhang, H. Wu, S. Wang, L. Ma, M. Long, itransformer:
876 Inverted transformers are effective for time series forecasting, arXiv preprint
877 arXiv:2310.06625 (2023).
- 878 [61] M. Shanker, M. Y. Hu, M. S. Hung, Effect of data standardization on neural
879 network training, Omega 24 (4) (1996) 385–397.
- 880 [62] N. Fei, Y. Gao, Z. Lu, T. Xiang, Z-score normalization, hubness, and few-
881 shot learning, in: Proceedings of the IEEE/CVF International Conference
882 on Computer Vision, 2021, pp. 142–151.
- 883 [63] IEA, World Energy Outlook 2023, IEA, Paris, 2023, licence: CC BY 4.0
884 (report); CC BY NC SA 4.0 (Annex A).
885 URL <https://www.iea.org/reports/world-energy-outlook-2023>

- 886 [64] X. Fang, S. Misra, G. Xue, D. Yang, Smart grid—the new and improved
887 power grid: A survey, *IEEE communications surveys & tutorials* 14 (4)
888 (2011) 944–980.
- 889 [65] A. Olabi, M. A. Abdelkareem, Renewable energy and climate change, Re-
890 newable and Sustainable Energy Reviews 158 (2022) 112111.
- 891 [66] J. Devlin, Bert: Pre-training of deep bidirectional transformers for language
892 understanding, arXiv preprint arXiv:1810.04805 (2018).
- 893 [67] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al.,
894 Language models are unsupervised multitask learners, OpenAI blog 1 (8)
895 (2019) 9.
- 896 [68] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman,
897 A. Mathur, A. Schelten, A. Yang, A. Fan, et al., The llama 3 herd of
898 models, arXiv preprint arXiv:2407.21783 (2024).