

---

# Faster and Cheaper Energy Demand Forecasting at Scale

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Energy demand forecasting is one of the most challenging tasks for grids operators.  
2 Many approaches have been suggested over the years to tackle it. Yet, those still  
3 remain too expensive to train in terms of both time and computational resources,  
4 hindering their adoption as customers behaviors are continuously evolving.  
5 We introduce Transplit, a new lightweight transformer-based model, which sig-  
6 nificantly decreases this cost by exploiting the seasonality property and learning  
7 typical days of power demand. We show that Transplit can be run efficiently on  
8 CPU and is several hundred times faster than state-of-the-art predictive models,  
9 while performing as well.

## 10 1 Introduction

11 Energy demand forecasting is of utmost importance for grid operators and energy producers to ensure  
12 the stability of their service. Over the past decades, many approaches to provide the most accurate  
13 forecasting have been suggested [10]. In the recent years, approaches based on transformers models  
14 have started to be explored, but the cost to train them at scale, in a context where frequent retraining  
15 is a norm rather than an exception, hinders their adoption [4, 3]. Indeed recent events (e.g. climate  
16 change, remote working, energy prices) have highlighted that the way we consume electricity is in  
17 constant evolution and models need to be regularly retrained to adapt to new contexts. This makes  
18 previously suggested methods hardly applicable and therefore calls for cheaper and faster models.

19 Hence, we introduce Transplit, a transformer-based approach trained to recognize typical days  
20 of consumption from various profiles, that requires less computation time and resources, while  
21 conserving state-of-the-art performances. To validate our approach, we compare it against 4 related  
22 state-of-the-art transformers based approaches on two different datasets of electricity consumption: a  
23 publicly available one and another from our industrial partner (a national grid operator), that cannot  
24 be shared for privacy reasons. Overall, we observe that our approach performances are in line with  
25 the best approaches, while being **380** to **940** times faster to train on a single CPU, without relying on  
26 dedicated resources such as GPU, allowing it to run on significantly cheaper platforms.

## 27 2 Forecasting Energy Demand

28 Load forecasting is essential for electricity grid management and security, with applications ranging  
29 from peak prediction, incident prevention, to maintenance planning. It usually involves to forecast the  
30 energy demand of a household or a neighborhood in kWh over a given period of time (from minutes  
31 to months) based on historical data. Early on, Machine Learning approaches have been suggested [2],  
32 but their application remained limited by the amount of available data and computing power, which  
33 is not an issue anymore with the development of smartgrid and the recent increase in computing  
34 capability. Among ML-based approaches, we distinguish two categories.

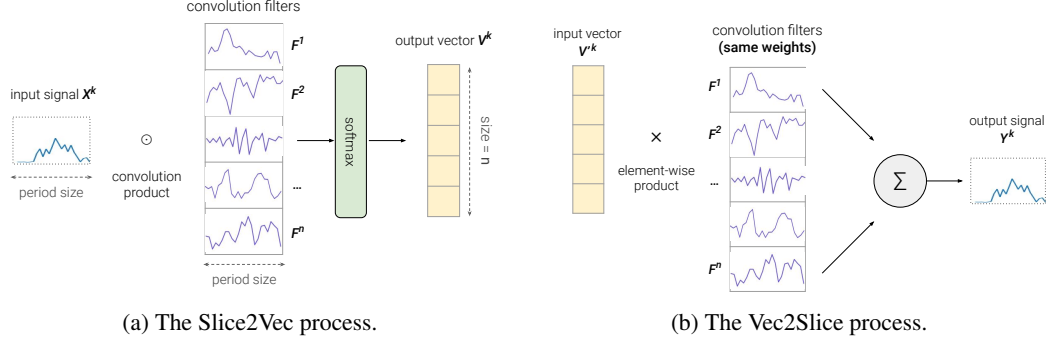


Figure 1: The SVS module is composed of two layers which share the same weights.

35 Approaches falling in the first one consist in training one model per household. This was usually  
 36 achieved by models such as random forests, SVMs, or linear regressions [11][3], but more recently  
 37 neural networks based approaches started to be suggested. Neural Prophet (2021) [15] for example  
 38 has been applied to electric load forecasting [14] with interesting results, as well as Deep Probabilistic  
 39 Koopman models [9]. Although these approaches perform well, the number of models to train for a  
 40 grid rapidly becomes too expensive to deploy.

41 On the opposite, approaches from the second category rely on a single model which generalizes  
 42 different behaviors of the demand time series, discerns consumers' profiles and prolong their demand  
 43 curve according to their shape and other factors. Those approaches are usually less costly, but also  
 44 less precise. LSTMs [5] have been widely used in this context [12][18] and are still the subject of  
 45 recent load forecasting research [7][18].

46 Yet, with the development of transformers [16], approaches such as LogTrans [8], Reformer [6],  
 47 Informer [19] and lastly Autoformer [17] started to appear with better accuracy and efficiency.  
 48 These approaches all derived from the original work differ in their inner attention mechanism that is  
 49 modify to lower the computation complexity. Still, these models are not entirely applicable to energy  
 50 demand forecasting as they require high computation power and considerable time to train, which is  
 51 problematic when models need to be retrained often. There is therefore a need for a much lighter  
 52 model, requiring less computational power and time, with at least similar performances.

### 53 3 Proposed approach

54 In this regards, we introduce Transplit, a lightweight load forecasting model. To design it, we started  
 55 from an observation on a fundamental difference between most of language models and time series  
 56 models, which also impacts their efficiency.

#### 57 3.1 Chunking time series

58 Whereas language models perform well by tokenizing input text [13] instead of splitting it character by  
 59 character, current time series models process every single input as a vector. As electrical consumption  
 60 data is seasonal, we propose to embed each season – in our case study, the time series is split into  
 61 days. The encoding and decoding parts are achieved by learning a vocabulary of daily consumptions.  
 62 Our solution then consists in reasoning in blocks of consumption days rather than single values.

##### 63 3.1.1 Slice2Vec

64 As described in Figure 1a, Slice2Vec takes a sequence  $X \in \mathbb{R}^L$ , where  $L = m \times T$  is the input  
 65 sequence length which is a multiple of a period  $T$ .  $X$  is split into  $m$  periods of size  $T$ , that we denote  
 66  $(X^k)_{k \in \{1, \dots, m\}}$ <sup>1</sup>.

67 For each  $X^k$ , we define a vector  $V^k$  with the following convolution product:

$$\forall i \in \{1, \dots, n\}, V_i^k = X^k \odot F^i \quad (1)$$

<sup>1</sup>Concretely, if the period  $T$  is one day, then  $X^k$  is the  $k^{th}$  day of the input.

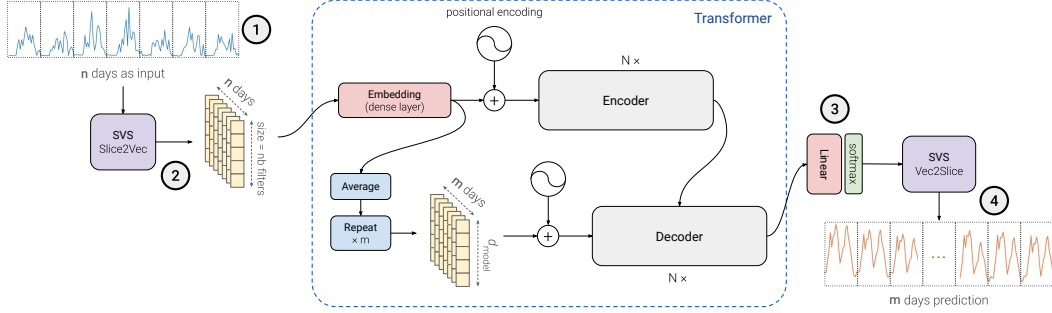


Figure 2: Model architecture of Transplit: (1) the input consumption is split into days; (2) Slice2Vec converts each day into a vector and passes it to the transformer (3) that forecasts vectors (4) which are converted to daily consumptions with Vec2Slice.

68 where  $n$  is the chosen size for the vectors,  $V_i^k$  is the  $i^{\text{th}}$  coefficient of  $V^k$ , and  $F^i$  is a convolution  
 69 filter of size  $T$ . Since  $X^k$  and  $F^i$  have the same size, the convolution product is also equivalent to the  
 70 sum of their term-wise product.

### 71 3.1.2 Vec2Slice

72 In order to transform a vector  $V^k$  back in a  $T$ -sized slice of time series called  $Y^k$ , the filters  $F^i$   
 73 are weighted according to the coefficients of  $V^k$ , then summed:

$$Y^k = \sum_i V_i^k F^i \quad (2)$$

74 as illustrated in Figure 1b. The slices  $Y^k$  are concatenated to form the output  $Y$ . It is important to  
 75 underline that Vec2Slice is not the inverse function of Slice2Vec:  $Vec2Slice(Slice2Vec(X)) \neq X$ .  
 76 Vec2Slice shares the filters weights with Slice2Vec. The filters are therefore used twice: to recognize  
 77 the input and to match the shape of the expected output.

### 78 3.2 Transplit

79 From these two components, we propose *Transplit*, a new load forecasting model based on the  
 80 original transformer [16] with simple full attention mechanisms. The architecture of the model is  
 81 presented in figure 2.

## 82 4 Evaluation

### 83 4.1 Criteria

84 Transplit is designed to (**req. 1**) maintain state-of-the-art performances for load forecasting, while  
 85 (**req. 2**) being lighter and faster, in terms of number of parameters and training time.

86 To validate these two requirements, we compare Transplit with 4 state-of-the-art approaches: Auto-  
 87 former (2021) [17], Informer (2021) [19], Reformer (2020) [6] and a vanilla Transformer architecture  
 88 (2017) [16]. Classical models such as LSTMs are disregarded as it has been shown that transformers  
 89 perform better for load forecasting [17, 19, 6]. The assessment is done on two forecasting ranges: 7  
 90 days and 30 days. We use two datasets and two performance metrics, and also measure training time  
 91 taken by using a GPU and only using a CPU (without multiprocessing).

### 92 4.2 Experimental setting

93 We evaluate those approaches on 2 datasets of household's hourly consumption expressed in kWh.  
 94 The first one, denoted **IND**, is provided by our industrial partner (a national grid operator) and  
 95 encompasses the consumption of 6010 households over 2 years (2020-2021). The second one **ECL**[1]  
 96 is an open dataset containing the consumption of 321 households over 3 years (2012-2014).

	ECL				IND			
	7 days		30 days		7 days		30 days	
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
<b>Transplit</b>	0.584	0.178	<b>0.688</b>	0.189	0.177	0.189	<b>0.202</b>	<b>0.205</b>
<b>Autoformer</b>	<b>0.577</b>	<b>0.153</b>	0.736	<b>0.174</b>	0.185	0.210	0.204	0.217
<b>Informer</b>	2.68	0.376	2.71	0.387	0.201	0.213	0.226	0.231
<b>Reformer</b>	1.28	0.262	1.79	0.313	<b>0.172</b>	<b>0.185</b>	0.205	0.207
<b>Transformer</b>	0.781	0.197	0.918	0.213	0.177	0.191	0.203	0.207

Table 1: Error metrics for the ECL and IND datasets

	<b>Transplit</b>	<b>Autoformer</b>	<b>Informer</b>	<b>Reformer</b>	<b>Transformer</b>
<b># of parameters</b>	<b>183,616</b>	10,505,217	11,306,497	5,782,529	10,518,529
<b>ECL: train. with GPU</b>	<b>3m02s</b>	4h16m	1h53m	4h13m	2h51m
<b>IND: train. with GPU</b>	<b>21m13s</b>	1d 12h15m	19h37m	20h56m	1d 9h00m
<b>ECL: train. with CPU</b>	<b>2m47s</b>	1d 18h01m	17h11m	1d 13h49m	1d 3h01m
<b>IND: train. with CPU</b>	<b>32m05s</b>	20d 21h59m	8d 13h21m	18d 19h47m	13d 10h42m

Table 2: Number of trainable parameters and training time for each model, dataset and infrastructure used, for a 720  $\rightarrow$  720 values forecast.

97 We set ratios of training / validation / testing to respectively 70 / 10 / 20%, based on the timeline. We  
98 feed the models with inputs of 720 values (30 days).

99 We configure all baseline approaches according to the recommendation from their authors for the  
100 ECL dataset. Regarding Transplit, we set  $T = 24$  hours and 512 filters for the SVS module and use 1  
101 encoder and 1 decoder, with a vector size ( $d_{model}$ ) of 64.

102 While the training data is standardized, we reverse it at the output in the testing phase to reflect the  
103 loss with the original scale in kWh. Further details are available on our repository<sup>2</sup>.

### 104 4.3 Results

105 Results for both datasets are presented in Table 1. Although there is no significant improvement,  
106 Transplit succeeds in keeping the performance of state-of-the-art time series transformer models.

107 The key element of Transplit is its lightness: what is important to distinguish is the time taken to  
108 train these models, shown in Table 2. We observe that Transplit took 21 minutes to learn the 6010  
109 consumers profiles of the IND dataset using a GPU, and is then 102 times faster than the Autoformer.

110 Another big advantage for Transplit is that using a CPU instead of a GPU doesn’t significantly slow  
111 down the model’s computation time, as processing smaller sequences implies smaller operations, and  
112 data still has to be cached in the GPU, which can take a lot of time in total. The non-necessity of  
113 GPU to efficiently run the model allows it to be used on less expensive infrastructures, while being  
114 faster and conserving good forecasting performances.

## 115 5 Conclusion

116 We have seen that Transplit remains in line with state-of-the-art time series deep learning models that  
117 are able to generalize consumption profiles and predict a detailed state of the grid in terms of electric  
118 load, while being several hundreds of times faster and lighter.

119 It opens up the possibility to run a deep-learning predictive model globally at a lower cost, as only  
120 one CPU is enough to keep Transplit efficient, making it possible to run it locally. This finally offers  
121 the possibility to update the model regularly with new data. Yet, Transplit might be ineffective to  
122 predict non-seasonal time series, but could be extended to seasonal signals in general.

<sup>2</sup><https://anonymous.4open.science/r/Transplit-BC41>

## References

- 123
- 124 [1] Electrical consumption load. [https://archive.ics.uci.edu/ml/datasets/](https://archive.ics.uci.edu/ml/datasets/ElectricityLoadDiagrams20112014)  
125 ElectricityLoadDiagrams20112014.
- 126 [2] Hesham Alfares and Nazeeruddin Mohammad. Electric load forecasting: Literature survey and  
127 classification of methods. *International Journal of Systems Science - IJSSc*, 33, 01 2002.
- 128 [3] Alfonso González-Briones, Guillermo Hernández, Juan M. Corchado, Sigeru Omatu, and  
129 Mohd Saberi Mohamad. Machine learning models for electricity consumption forecasting: A  
130 review. In *ICCAIS'19*, 2019.
- 131 [4] H.S. Hippert, C.E. Pedreira, and R.C. Souza. Neural networks for short-term load forecasting:  
132 a review and evaluation. *IEEE Transactions on Power Systems*, 16(1), 2001.
- 133 [5] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8),  
134 1997.
- 135 [6] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer, 2020.
- 136 [7] Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long- and short-term  
137 temporal patterns with deep neural networks. In *SIGIR '18*, 2018.
- 138 [8] Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyong Zhou, Wenhui Chen, Yu-Xiang Wang, and Xifeng  
139 Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series  
140 forecasting. In *NIPS'19*, volume 32, 2019.
- 141 [9] Alex Mallen, Henning Lange, and J. Nathan Kutz. Deep probabilistic koopman: Long-term  
142 time-series forecasting under periodic uncertainties, 2021.
- 143 [10] Isaac kofi Nti, Adebayo Adekoya, Owusu Nyarko-Boateng, and Moses Teimah. Electricity load  
144 forecasting: a systematic review. *Journal of Electrical Systems and Information Technology*, 7,  
145 09 2020.
- 146 [11] A. Parrado-Duque, S. Kelouwani, K. Agbossou, S. Hosseini, N. Henao, and F. Amara. A  
147 comparative analysis of machine learning methods for short-term load forecasting systems. In  
148 *SmartGridComm'21*, 2021.
- 149 [12] Andrew Pulver and Siwei Lyu. Lstm with working memory. In *IJCNN'17*, 2017.
- 150 [13] Phillip Rust, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder, and Iryna Gurevych. How good is  
151 your tokenizer? on the monolingual performance of multilingual language models. In *ACL'21*,  
152 August 2021.
- 153 [14] Md Jamal Ahmed Shohan, Md Omar Faruque, and Simon Y. Foo. Forecasting of electric load  
154 using a hybrid lstm-neural prophet model. *Energies*, 15(6), 2022.
- 155 [15] Oskar Triebe, Hansika Hewamalage, Polina Pilyugina, Nikolay Laptev, Christoph Bergmeir,  
156 and Ram Rajagopal. NeuralProphet: Explainable Forecasting at Scale, 2021.
- 157 [16] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
158 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS'17*, 2017.
- 159 [17] haixu wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition  
160 transformers with auto-correlation for long-term series forecasting. In *NIPS'21*, 2021.
- 161 [18] Bailin Yang, Shulin Sun, Jianyuan Li, Xianxuan Lin, and Yan Tian. Traffic flow prediction  
162 using lstm with feature enhancement. *Neurocomputing*, 332, 2019.
- 163 [19] Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai  
164 Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In  
165 *AAAI'21*, 2021.

166 **Checklist**

167 The checklist follows the references. Please read the checklist guidelines carefully for information on  
168 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
169 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
170 the appropriate section of your paper or providing a brief inline description. For example:

- 171 • Did you include the license to the code and datasets? **[Yes]** See Section 4.2.

172 Please do not modify the questions and only use the provided macros for your answers. Note that the  
173 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
174 block and only keep the Checklist section heading above along with the questions/answers below.

175 1. For all authors...

- 176 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
177 contributions and scope? **[Yes]**  
178 (b) Did you describe the limitations of your work? **[Yes]** See the discussion section.  
179 (c) Did you discuss any potential negative societal impacts of your work? **[No]**  
180 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
181 them? **[Yes]**

182 2. If you are including theoretical results...

- 183 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**  
184 (b) Did you include complete proofs of all theoretical results? **[N/A]**

185 3. If you ran experiments...

- 186 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
187 mental results (either in the supplemental material or as a URL)? **[Yes]** See Github link  
188 as a footnote in subsection 4.2.  
189 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
190 were chosen)? **[Yes]** Subsection 4.2.  
191 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
192 ments multiple times)? **[No]**  
193 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
194 of GPUs, internal cluster, or cloud provider)? **[Yes]** 4.2 and 4.3.

195 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 196 (a) If your work uses existing assets, did you cite the creators? **[Yes]** Link to UCI ML  
197 repository with no citation request  
198 (b) Did you mention the license of the assets? **[No]**  
199 (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**  
200 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
201 using/curating? **[Yes]** for the dataset from our industrial partner, **[No]** for ECL.  
202 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
203 information or offensive content? **[Yes]** with our partner.

204 5. If you used crowdsourcing or conducted research with human subjects... **[N/A]**