

Paper Review: Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecasting

1

Abdullah Mamun

May 10, 2021

About this paper

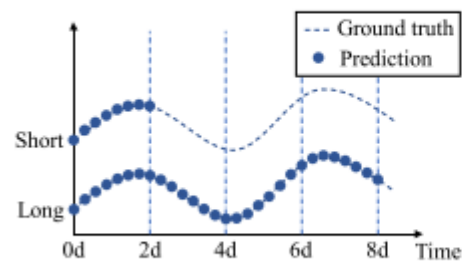
Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecasting

Haoyi Zhou,¹ Shanghang Zhang,² Jieqi Peng,¹ Shuai Zhang,¹ Jianxin Li,¹
Hui Xiong,³ Wancai Zhang⁴

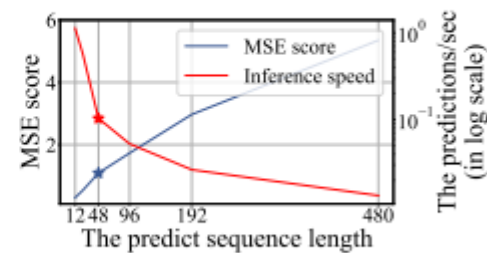
¹ Beihang University ² UC Berkeley ³ Rutgers University ⁴ SEDD Company
{zhouhy, pengjq, zhangs, lijx}@act.buaa.edu.cn, shz@eecs.berkeley.edu, {xionghui, zhangwancaibuaa}@gmail.com

Abstract

Many real-world applications require the prediction of long sequence time-series, such as electricity consumption planning. Long sequence time-series forecasting (LSTF) demands a high prediction capacity of the model, which is the ability to capture precise long-range dependency coupling between output and input efficiently. Recent studies have shown the



(a) Sequence Forecasting.



(b) Run LSTM on sequences.

About this paper

- ▶ Published in the **35th AAAI** conference in 2021.
- ▶ Cited by **7** as of May 10, 2021
- ▶ Received the best paper award in 35th AAAI 2021.
- ▶ Main goal is building a **forecasting model** for multi-dimensional time series data.

Why Informer?

Transformer cell:

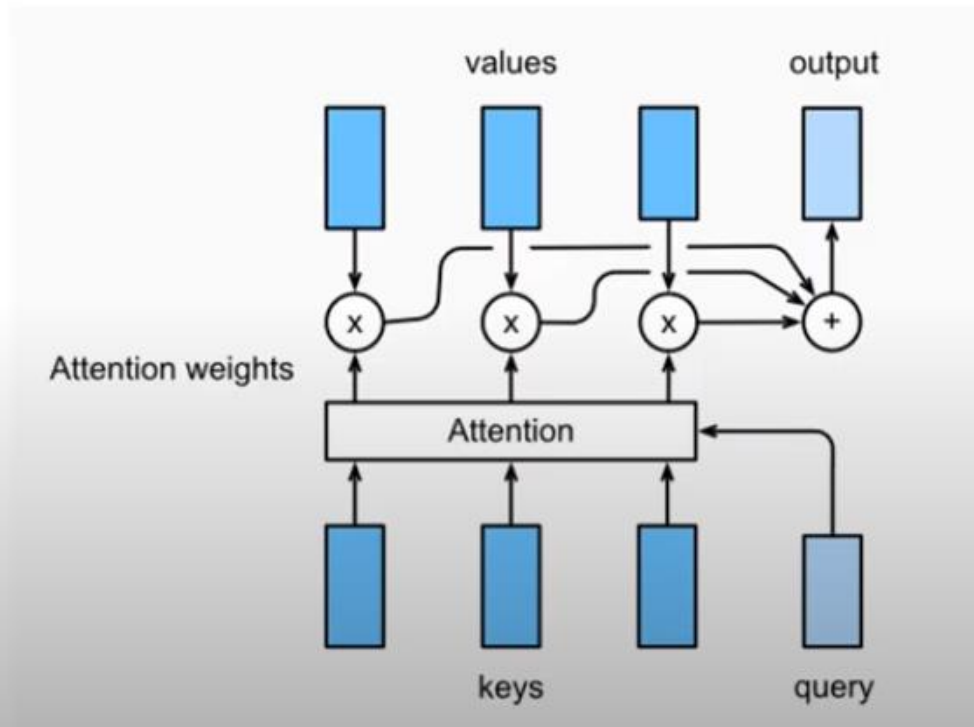
High 1. Space-2. Time complexity

$$\text{softmax} \left(\frac{\begin{matrix} \text{Q} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{matrix} \square \\ \square \\ \square \end{matrix} \end{matrix}}{\sqrt{d_k}} \right) \begin{matrix} \text{V} \\ \begin{matrix} \square & \square \\ \square & \square \end{matrix} \end{matrix}$$
$$= \begin{matrix} \text{Z} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix} \quad O(L_K L_Q)$$

Why Informer?

Attention mechanism:

3. Predicts one output at a time

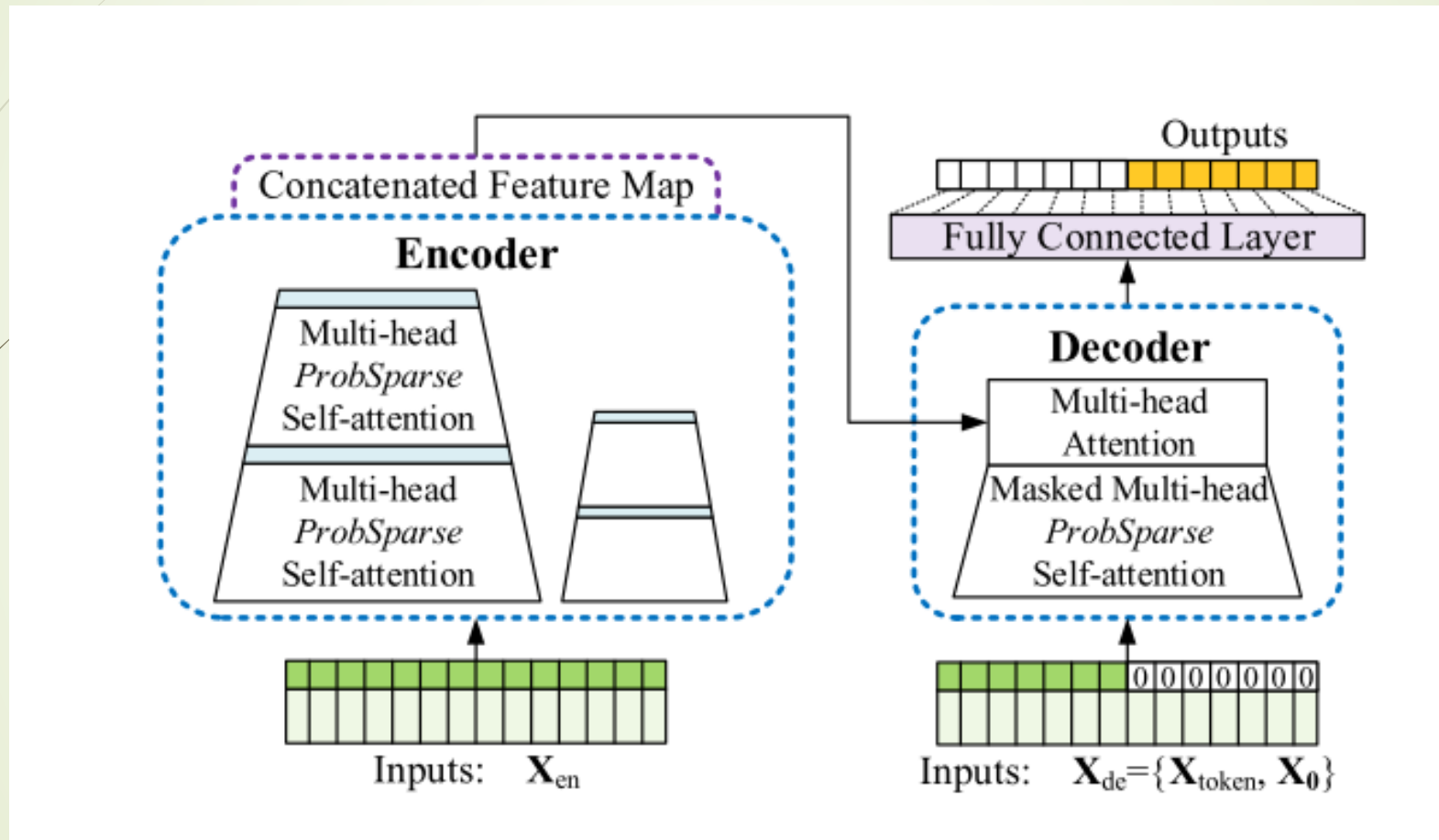


Why Informer?

Informer:

1. **Query sparsity** Measurement lowers **Space** complexity
2. **ProbSparse** lowers **Time** Complexity
3. Predicts **sequence** in one batch (Generative Style Decoder)

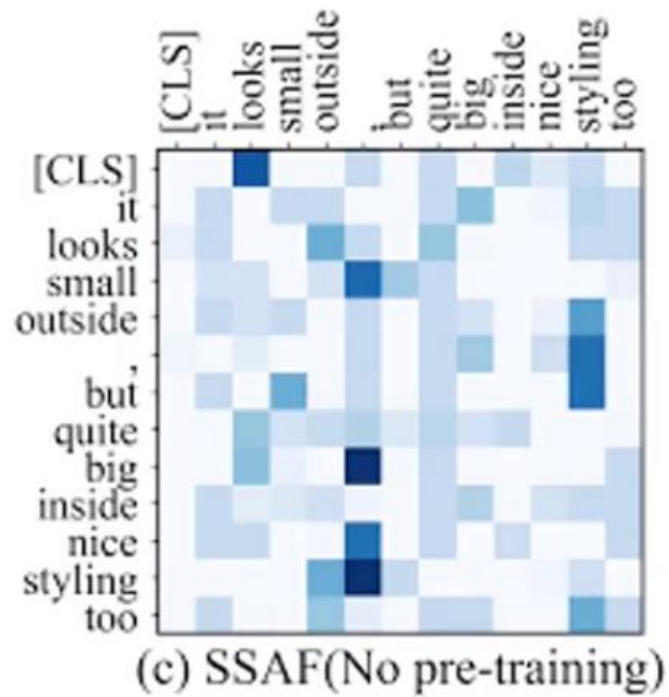
Informer Model Overview



Main Improvements

1

Attention values:



long-tail distributed

Therefore:

Prune needed.

Main Improvements

$$\mathcal{A}(\mathbf{q}_i, \mathbf{K}, \mathbf{V}) = \sum_j \frac{k(\mathbf{q}_i, \mathbf{k}_j)}{\sum_l k(\mathbf{q}_i, \mathbf{k}_l)} \mathbf{v}_j = \mathbb{E}_{p(\mathbf{k}_j | \mathbf{q}_i)} [\mathbf{v}_j]$$

1

Therefore: thresholding the queries, so that queries inducing top-u KLD-values are kept (rests set to zero)

$u = c \ln L_Q$ (not proven!)

Main Improvements

2

Attention values:

$$\text{softmax} \left(\frac{Q \times K^T}{\sqrt{d_k}} \right) V = Z$$

The diagram shows the attention mechanism. A purple 2x3 matrix Q is multiplied by an orange 3x2 matrix K^T (highlighted with a red box and a red checkmark). The result is divided by the square root of the key dimension $\sqrt{d_k}$. This is followed by a softmax operation and multiplication by a blue 2x2 matrix V to produce a pink 2x2 matrix Z .

Main Improvements

However, the traversing of all the queries for the measurement $M(\mathbf{q}_i, \mathbf{K})$ requires calculating each dot-product pairs, i.e., quadratically $\mathcal{O}(L_Q L_K)$, besides the LSE operation has the potential numerical stability issue. Motivated by this, we propose an empirical approximation for the efficient acquisition of the query sparsity measurement.

Lemma 1. *For each query $\mathbf{q}_i \in \mathbb{R}^d$ and $\mathbf{k}_j \in \mathbb{R}^d$ in the keys set \mathbf{K} , we have the bound as $\ln L_K \leq M(\mathbf{q}_i, \mathbf{K}) \leq \max_j \{\mathbf{q}_i \mathbf{k}_j^\top / \sqrt{d}\} - \frac{1}{L_K} \sum_{j=1}^{L_K} \{\mathbf{q}_i \mathbf{k}_j^\top / \sqrt{d}\} + \ln L_K$. When $\mathbf{q}_i \in \mathbf{K}$, it also holds.*

From the Lemma 1 (proof is given in Appendix D.1), we propose the max-mean measurement as

$$\bar{M}(\mathbf{q}_i, \mathbf{K}) = \max_j \left\{ \frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d}} \right\} - \frac{1}{L_K} \sum_{j=1}^{L_K} \frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d}}. \quad (4)$$

Main Improvements

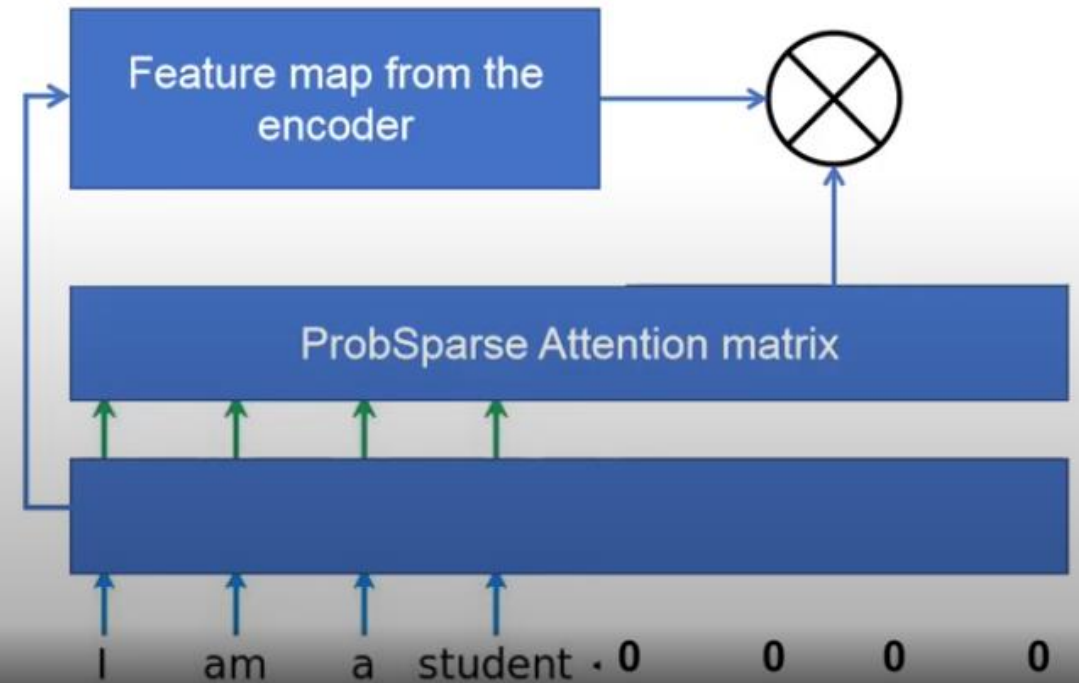
12

3

Generative Style Decoder:

Informer:
Predicts the target in one batch

Je suis étudiant </s>



Experiments and Results

13

4 Experiment

Datasets

We extensively perform experiments on four datasets, including 2 collected real-world datasets for LSTF and 2 public benchmark datasets.

ETT (Electricity Transformer Temperature)^[2]: The ETT is a crucial indicator in the electric power long-term deployment. We collected 2-year data from two separated counties in China. To explore the granularity on the LSTF problem, we create separate datasets as $\{ETTh_1, ETTh_2\}$ for 1-hour-level and $ETTm_1$ for 15-minute-level. Each data point consists of the target value "oil temperature" and 6 power load features. The train/val/test is 12/4/4 months.

ECL (Electricity Consuming Load)^[3]: It collects the electricity consumption (Kwh) of 321 clients. Due to the missing data (Li et al. 2019), we convert the dataset into hourly consumption of 2 years and set 'MT_320' as the target value. The train/val/test is 15/3/4 months.

Weather^[4]: This dataset contains local climatological data for nearly 1,600 U.S. locations, 4 years from 2010 to 2013, where data points are collected every 1 hour. Each data point

Experiments and Results

14

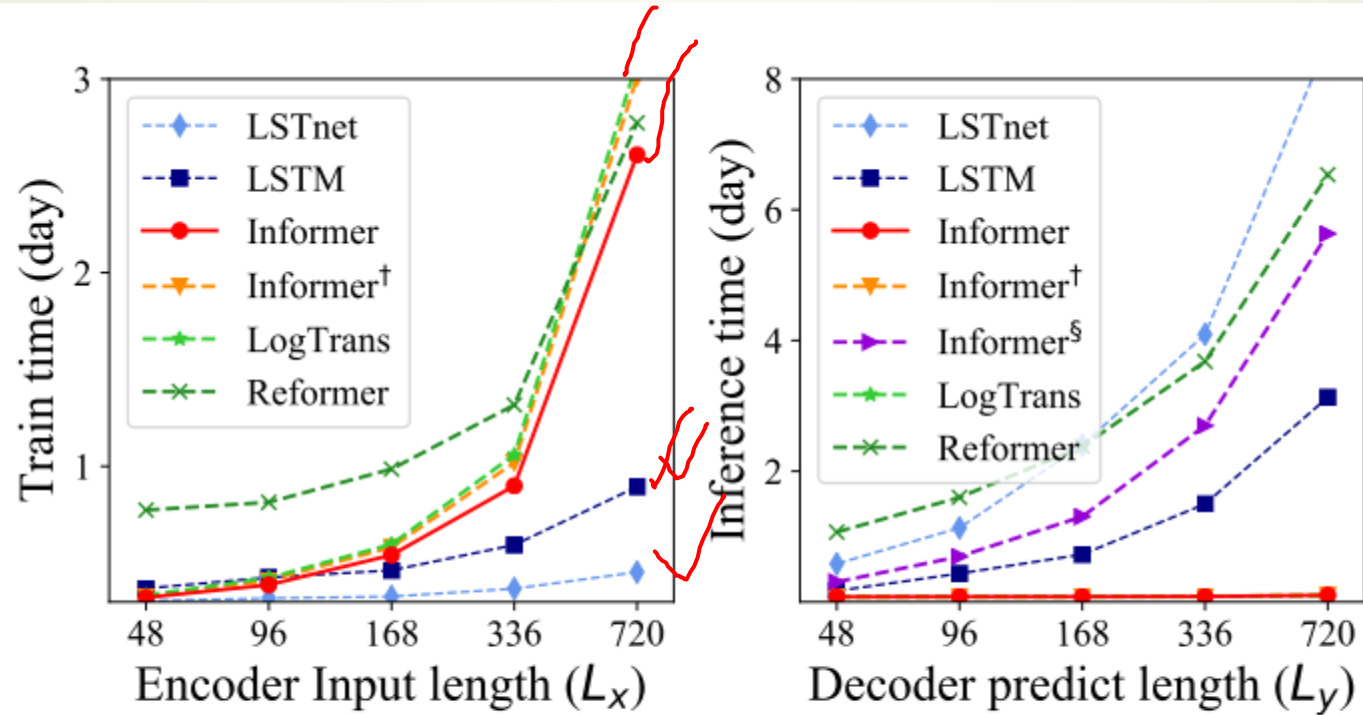
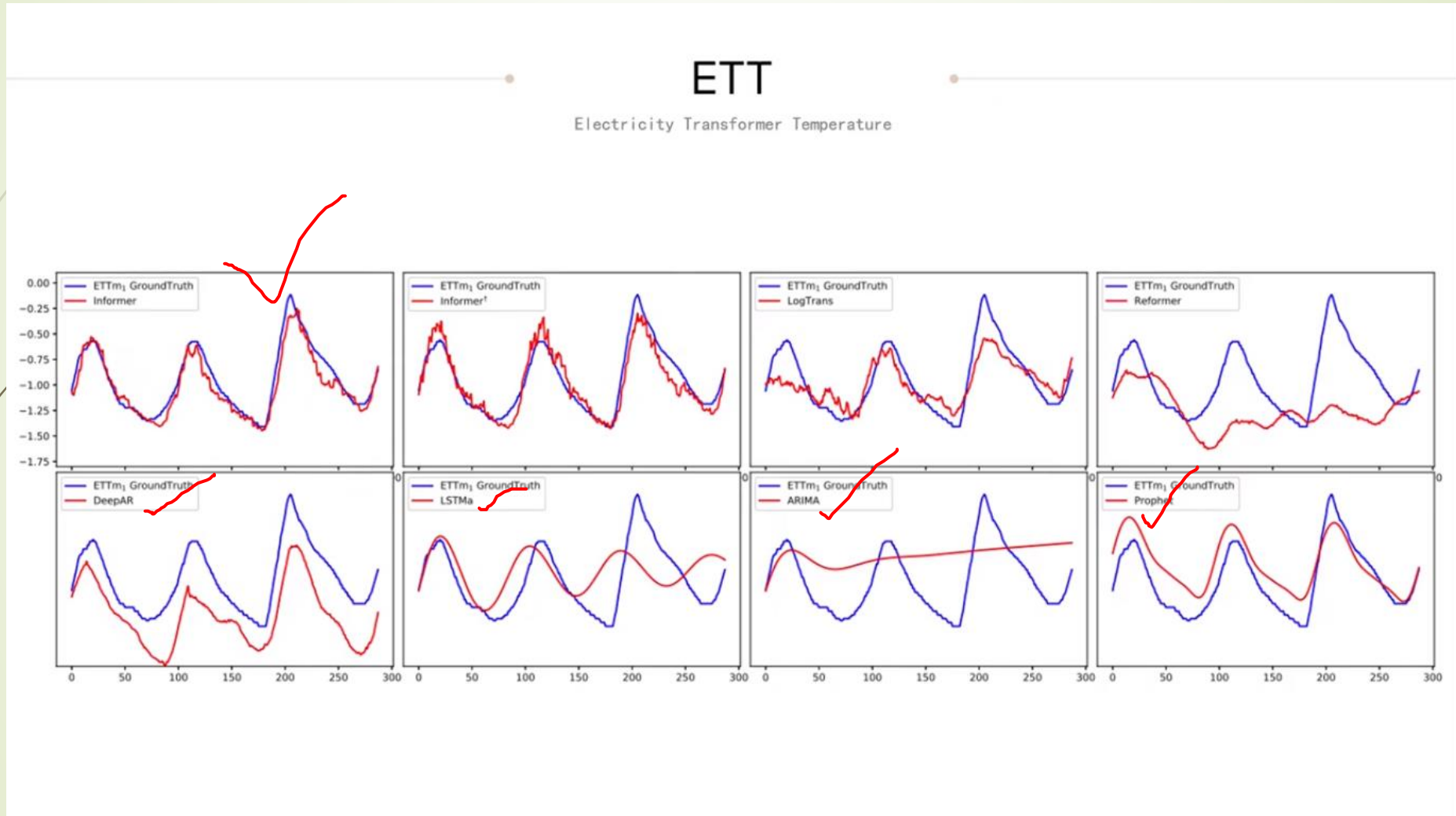


Figure 5: The total runtime of training/testing phase.

Experiments and Results

15



Limitations

16

1. assume all keys are IID following some Gaussian distributions **(basic assumption)**
2. It holds only for queries resulting in top-u KLDs **(condition)**.
3. The attention values are long-tail distributed **(verified, but on one dataset only)**.
4. The variance of query-key dot products also decreases along with the KLD **(Alert, comes from no where)**.

Thank You!

Contact: abdullahal.mamun1@wsu.edu