# 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], Xiong Hui[3], Wancai Zhang[4]
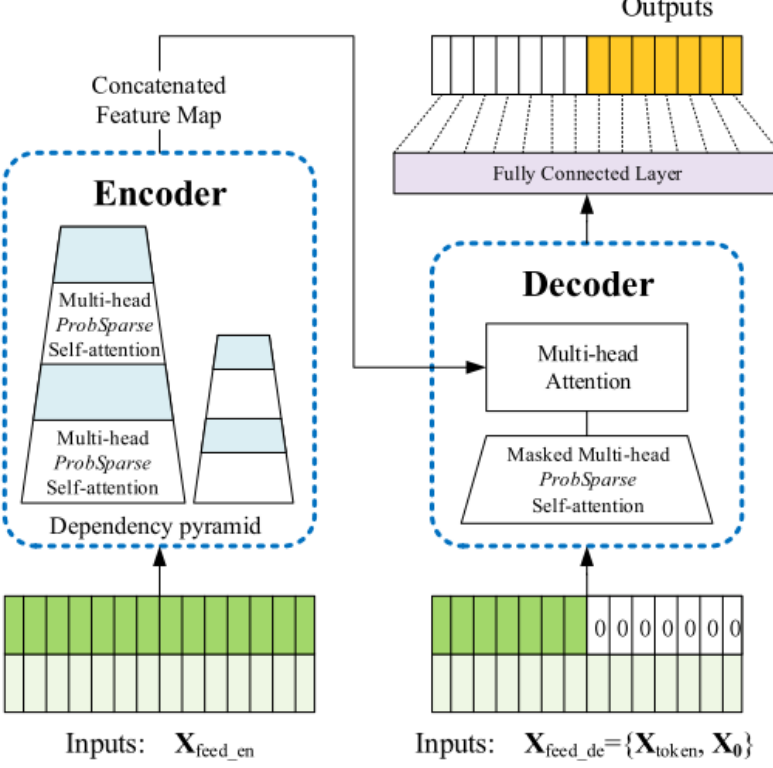
[1] Beihang University, [2] UC Berkeley, [3] Rutgers University,
[4] Beijing Guowang Fuda Science & Technology Development Company

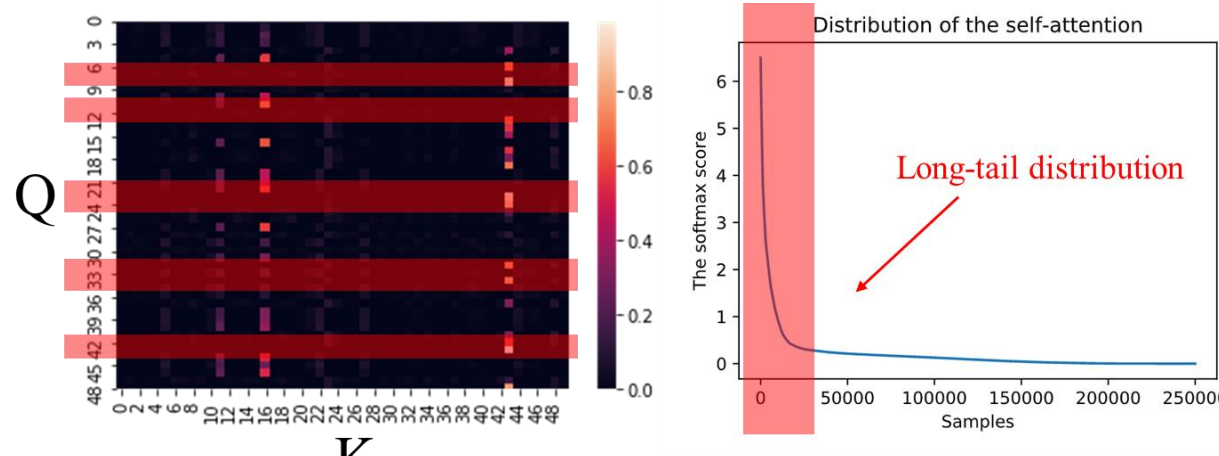## Highlights

❖ Propose Informer to successfully enhance the prediction capacity in the Long Sequence Time-series Forecasting (LSTF) problem.

❖ Propose **ProbSparse Self-Attention** mechanism that achieves $\mathcal{O}(L\log L)$ in time complexity and memory usage.

❖ Propose **Self-attention Distilling** to sharply reduce space complexity to $\mathcal{O}((2-\epsilon)L\log L)$.

❖ Propose **Generative Style Decoder** to acquire long sequence output with only one forward step needed.

## Introduction

❖ **Long Sequence Time-series Forecasting (LSTF)**

- Near future predictions
- **Limited** adjustment

- Coarse predictions
- **Inefficient** adjustment

- Long sequence predictions
- **Proper** adjustment

The major challenge for LSTF in **enhancing the prediction capacity to meet the increasingly long sequences demand**

- Extraordinarily **long-range alignment ability**
- **Efficient operations** on long sequence inputs and outputs

❖ **Applying Transformer models in LSTF problem**

**Limitations of LSTM in LSTF**

- The inference speed of LSTM decrease rapidly.
- Continuous accumulation of errors causes the MSE score increase rapidly.

**Advances of Transformer**

- The self-attention mechanism **reduce max path of network signals to $\mathcal{O}(1)$ and avoids recurrent structure**.

**Limitations of Transformer**

- *The quadratic computation of self-attention.* The atom operation of self-attention mechanism, namely canonical dot-product, causes the time complexity and memory usage per layer to be $\mathcal{O}(L^2)$.
- *The memory bottleneck in stacking layers for long inputs.* The stack of $J$ encoder/decoder layer makes total memory usage to be $\mathcal{O}(J \cdot L^2)$ which limits the model scalability on receiving long sequence inputs.
- *The speed plunge in predicting long outputs.* The dynamic decoding of vanilla Transformer makes the step-by-step inference as slow as RNN-based model.

## Methodology

❖ *ProbSparse* **Self-attention**

Rewrite the self-attention into the probability formulation

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

Establish the measurement

Attention probability $p(\mathbf{k}_j|\mathbf{q}_i)$ **?** Uniform probability $q(\mathbf{k}_j|\mathbf{q}_i) = \frac{1}{L_K}$

$$KL(q\|p) = \sum_{j=1}^{L_K} \frac{1}{L_K} \ln \frac{1/L_K}{k(\mathbf{q}_i, \mathbf{k}_j)/\sum_l k(\mathbf{q}_i, \mathbf{k}_l)}$$

$$= \ln \sum_{l=1}^{L_K} e^{\frac{\mathbf{q}_i \mathbf{k}_l^\top}{\sqrt{d}}} - \frac{1}{L_K} \sum_{j=1}^{L_K} \frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d}} - \ln L_K$$

Dropping the constant, and **use an approximation** to the measurement for computation simplicity:

$$M(\mathbf{q}_i, \mathbf{K}) = \ln \sum_{j=1}^{L_K} e^{\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d}}} - \frac{1}{L_K} \sum_{j=1}^{L_K} \frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d}} \quad \textbf{Replace} \quad \max_j\{\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d}}\}$$

Define *ProbSparse* self-attention

$$\mathcal{A}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}(\frac{\overline{\mathbf{Q}}\mathbf{K}^\top}{\sqrt{d}})\mathbf{V}$$

Where $\overline{Q}$ is a sparse matrix of the same size of $q$ and it only contains the Top-$u$ queries under the sparsity measurement $M(q, K)$.
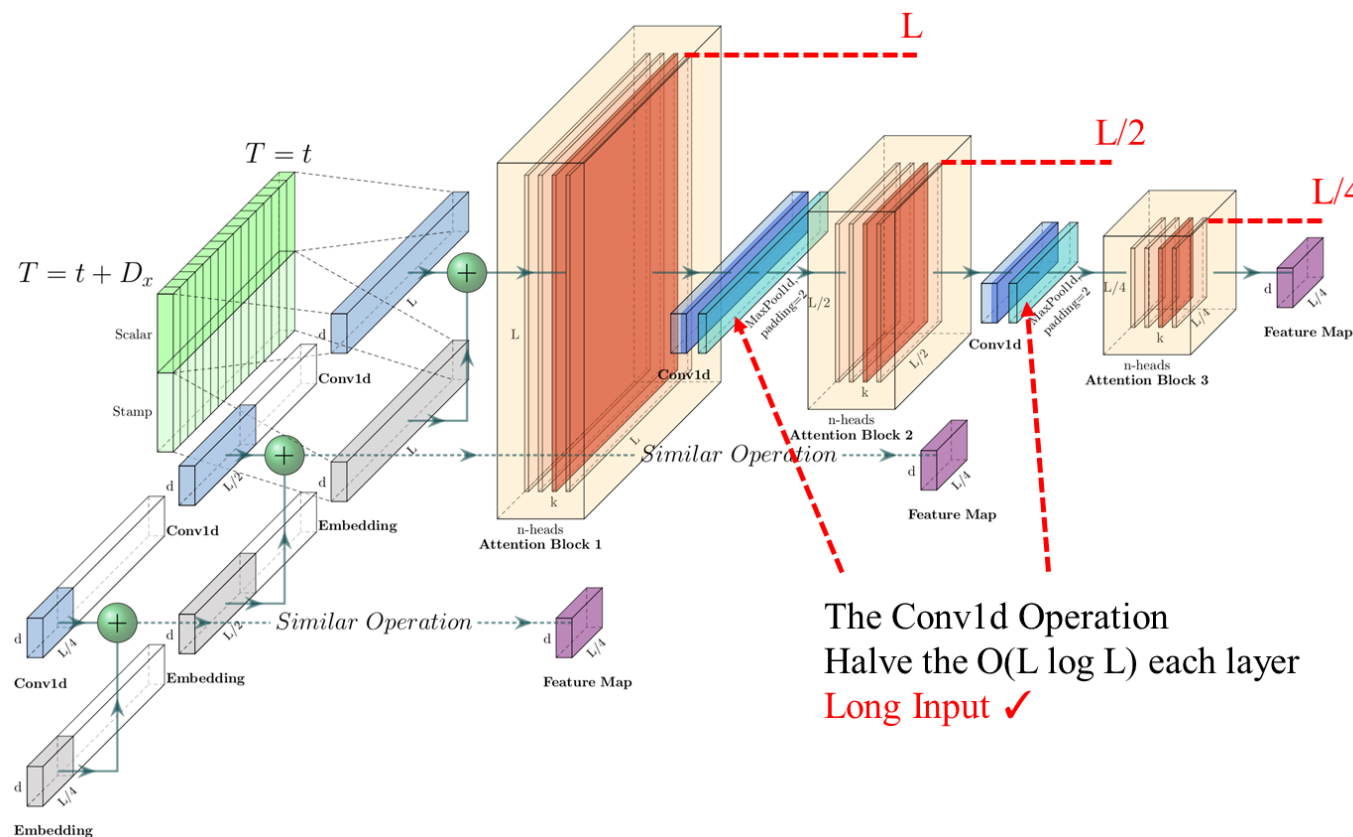
**Algorithm 1** ProbSparse self-attention

**Input:** Tensor $\mathbf{Q} \in \mathbb{R}^{m \times d}$, $\mathbf{K} \in \mathbb{R}^{n \times d}$, $\mathbf{V} \in \mathbb{R}^{n \times d}$
1: **initialize:** set hyperparameter $c$, $u = c \ln m$ and $U = m \ln n$
2: randomly select $U$ dot-product pairs from $\mathbf{K}$ as $\bar{\mathbf{K}}$
3: set the sample score $\bar{\mathbf{S}} = \mathbf{Q}\bar{\mathbf{K}}^\top$
4: compute the measurement $M = \max(\bar{\mathbf{S}}) - \text{mean}(\bar{\mathbf{S}})$ by row
5: set Top-$u$ queries under $M$ as $\bar{\mathbf{Q}}$
6: set $\mathbf{S}_1 = \text{softmax}(\bar{\mathbf{Q}}\mathbf{K}^\top/\sqrt{d}) \cdot \mathbf{V}$
7: set $\mathbf{S}_0 = \text{mean}(\mathbf{V})$
8: set $\mathbf{S} = \{\mathbf{S}_1, \mathbf{S}_0\}$ by their original rows accordingly
**Output:** self-attention feature map S.

Sample $L\log L$ dot-product pairs → $\mathcal{O}(L\log L)$

❖ **Self-attention Distilling**

The Conv1d Operation
Halve the O(L log L) each layer
Long Input ✓

❖ **Generative-style Decoder**

$$\mathbf{X}_{\text{feed\_de}}^t = \text{Concat}(\mathbf{X}_{\text{token}}^t, \mathbf{X}_0^t) \in \mathbb{R}^{(L_{\text{token}}+L_y)\times d_{\text{model}}}$$

**Predicts all the outputs by one forward procedure.**

- $\mathbf{X}_{\text{token}}^t$: **Generative start token**, instead of choosing a specific flag as the token, we sample a "shorter" long sequence in input sequence, which is an earlier slice before output sequence.

## Experiment

❖ **Univariate Time-series Forecasting**

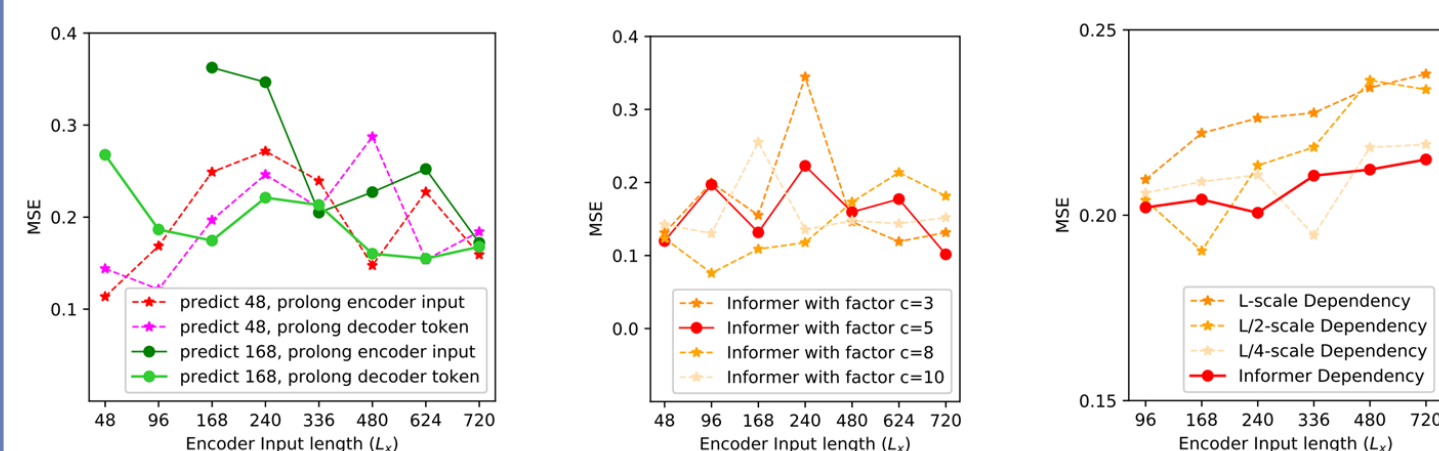Table 1: Univariate long sequence time-series forecasting results on four datasets (five cases)

❖ **Multivariate Time-series Forecasting**

Table 2: Multivariate long sequence time-series forecasting results on four datasets (five cases)

❖ **Parameter Sensitivity**

(a) Input length.  (b) Sampling.  (c) Stacking.

❖ **Ablation study**

Table 3: Ablation of *ProbSparse* mechanism

| Prediction length | 336 | | | 720 | | |
|---|---|---|---|---|---|---|
| Encoder's input | 336 | 720 | 1440 | 720 | 1440 | 2880 |
| Informer MSE | 0.243 | 0.225 | 0.212 | 0.258 | 0.238 | 0.224 |
| Informer MAE | 0.487 | 0.404 | 0.381 | 0.503 | 0.399 | 0.387 |
| Informer† MSE | 0.214 | 0.205 | - | 0.235 | - | - |
| Informer† MAE | 0.369 | 0.364 | - | 0.401 | - | - |
| LogTrans MSE | 0.256 | 0.233 | - | 0.264 | - | - |
| LogTrans MAE | 0.496 | 0.412 | - | 0.523 | - | - |
| Reformer MSE | 1.848 | 1.832 | 1.817 | 2.094 | 2.055 | 2.032 |
| Reformer MAE | 1.054 | 1.027 | 1.010 | 1.363 | 1.306 | 1.334 |

[1] Informer† uses the canonical self-attention mechanism.
[2] The '-' indicates failure for out-of-memory.

Table 4: Ablation of Self-attention Distilling

| Prediction length | 336 | | | | | 480 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Encoder's input | 336 | 480 | 720 | 960 | 1200 | 336 | 480 | 720 | 960 | 1200 |
| Informer† MSE | 0.201 | 0.175 | 0.215 | 0.185 | 0.172 | 0.136 | 0.213 | 0.178 | 0.164 | 0.171 |
| Informer† MAE | 0.360 | 0.335 | 0.366 | 0.355 | 0.321 | 0.282 | 0.382 | 0.345 | 0.296 | 0.272 |
| Informer‡ MSE | 0.187 | 0.182 | 0.177 | - | - | 0.208 | 0.182 | 0.168 | - | - |
| Informer‡ MAE | 0.330 | 0.341 | 0.329 | - | - | 0.384 | 0.337 | 0.304 | - | - |

[1] Informer‡ removes the self-attention distilling from Informer†.
[2] The '-' indicates failure for out-of-memory.

Table 5: Ablation of Generative Style Decoder

| Prediction length | 336 | | | | 480 | | | |
|---|---|---|---|---|---|---|---|---|
| Prediction offset | +0 | +12 | +24 | +48 | +0 | +48 | +96 | +168 |
| Informer‡ MSE | 0.101 | 0.102 | 0.103 | 0.103 | 0.155 | 0.158 | 0.160 | 0.165 |
| Informer‡ MAE | 0.215 | 0.218 | 0.223 | 0.227 | 0.317 | 0.397 | 0.399 | 0.406 |
| Informer§ MSE | 0.152 | - | - | - | 0.462 | - | - | - |
| Informer§ MAE | 0.294 | - | - | - | 0.595 | - | - | - |

[1] Informer§ replaces our decoder with dynamic decoding one in Informer‡.
[2] The '-' indicates failure for the unacceptable metric results.

❖ **Computation Efficiency**

论文预印版本&补充材料：https://arxiv.org/abs/2012.07436

论文项目地址　　开源数据集地址　　北航ACTBD公众号