国际人工智能会议 AAAI 2021 论文北京预讲会

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

Authors: Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang Jianxin Li, Xiong Hui, Wancai Zhang









## CONTENT

#### Background

Motivation: why self-attention?
Methods: the details of Informer
Experiments
Summary

#### What's the main topic of this paper?

Sequence prediction is a problem that involves using historical sequence information to predict the next value or values in the sequence. It is a basic but important research problem.



Stock Market Prediction



**Robots Action Prediction** 



Human Position Prediction



Weather Prediction



Supply Chain prediction



Society Event Prediction

#### An Example

Electricity Transformer (over 100kv) is a valuable asset for the energy company. The cost of a single unit can reach 30 million US dollars, which makes only conservative adjustments are allowed, i.e., one adjustment lasts 2 weeks, for stability.





Ground Truth

Time

Predictions

6d

#### Long Sequence Predictions

Problem: Long Sequence Time-series Forecasting (LSTF) is an important problem that has not been well addressed in the time-series analysis field for a long time. On the other hand, because of the "huge data, high frequency, long prediction" features of the real-world time-series data, how to solve this problem is one of the key limitations in big data applications.



As the length of the data sequence increases, the inference speed of LSTM decreases rapidly, resulting in the inability to predict long sequences with limited computing power and time successfully. Meanwhile, the continuous accumulation of error causes the MSE score to increase rapidly, making the results unusable.

#### Another Similar Problem ...

**Discrimination: Long Sequence Input Learning (LSIL)** is another important problem that remains to be solved in the machine learning field. It is not necessary related to the LSTF problem but the literature always refers to LSIL when talking about the long sequence.

#### Example:



**Classification:** All-age Novels

□ **Title Summary:** A boy wizard begins training and must battle for his life with the Dark Lord who murdered his parents.

Output

Short Summary: A boy who learns on his eleventh birthday that he is the orphaned son of two powerful wizards and possesses unique magical powers of his own. He is summoned from his life as an unwanted child to become a student at Hogwarts, an English boarding school for wizards. There, he meets several friends who become his closest allies and help him discover the truth about his parents' mysterious deaths.

#### The previous research on LSIL

Literature Review of Long Sequence Input Learning Problem (LSIL) We capture the long term dependencies with gradient descent, however, is difficult because the gradients computed by BPTT tend to vanish or explode (Hochreiter et al., 2001).





#### The previous research on LSIL

#### Literature Review of Long Sequence Input Learning Problem (LSIL)

We capture the long term dependencies with gradient descent, however, is difficult because the gradients computed by BPTT tend to vanish or explode (Hochreiter et al., 2001).

Practice: Truncate Sequences / Summarize Sequences / Random Sampling. These methods reduce the length of the input data to make it tractable.
 Regression: fast F2F '14 / additive F2F '18. These methods can perform the large-scale regression.

- RNN-based model: Periodically States '16 / Truncated BPTT '17 / Auxiliary Losses '18 / Recurrent Highway Networks '19 / Bootstrapping Regularizer '20. They try to only calculate <u>selected hidden states</u> or storage hidden states in the memory.
- Attention-based model: Naive Attention '14 / End-to-end Memory '15 / ED network '17. They improve the gradient flow by <u>shortening the path</u> that relevant information needs to traverse in the computation graph.
- CNN-based model: Wavenet '16 / TCNnet 18' / Seq-U-Net 19'. They use the convolutional neural networks with dilated filters to capture the long term dependencies. Their receptive fields grow exponentially with the number of layers.
- Self-attention model: Transformer '17 / Bert '19 / Transformer-xl '19. The calculation of self-attention is quadratic in the length of the sequence inputs.

## problem



## CONTENT

Background
Motivation: why self-attention?
Methods: the details of Informer
Experiments
Summary

## problem?

Problem Setting: Long Sequence Forecasting predicts the long sequence output from massive long sequence inputs.

How do we answer a question?



How do we plan gifts for Christmas day?



We will have some feelings of what you have seen when it comes to the prediction.

#### The intuitive introduction of attention

Attention Mechanism is arguably one of the most powerful concepts in the deep learning field nowadays. It is based on a common intuition that we "attend to" a certain part when processing a large amount of information.



#### Example:

We can explain the relationship between words in one sentence or close context. When we see "eating", we expect to encounter a food word very soon. The color term describes the food, but probably not so much with "eating" directly.



Acknowledgement from https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html

### The self-attention in NLP/CV field









Acknowledgement to Figure from http://jalammar.github.io/illustrated-bert/

#### The previous research on efficient self-attention



Tay Y, Dehghani M, Bahri D, et al. Efficient transformers: A survey[J]. arXiv preprint arXiv:2009.06732, 2020.

#### If we apply the Transformer in LSTF problem ...

	ARIMA	Prophet	LSTM	Conv	Transformer
Long Input	1	1	×	1	×
Long Output	×	×	×	×	×
Complexity / layer	O(L)	Not clear	O(L*d <sup>2</sup> )	O(k*L*d)	O(L <sup>2</sup> *d)
Max Path			O(L)	O(log <sub>k</sub> L)	O(1)

L is sequence length, D is the representation dimension, k is the kernel size of Convs.



## CONTENT

Background
Motivation: why self-attention?
Methods: the details of Informer
Experiments
Summary

#### Challenges

#### 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 O(L2).

[Complexity/layer]

#### The memory bottleneck in stacking layers.

The stack of J encoder/decoder layer makes total memory usage to be  $O(J \cdot L2)$ , which limits the model scalability on receiving long sequence inputs. [Long Input]

#### The speed plunge in predicting long outputs.

The dynamic decoding of vanilla Transformer makes the inference speed as slow as RNN-based model.

[Long Output]



An overall graph of the Informer model.

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







#### **Previous work:**







Rewrite the original self-attention into the probability formulation

$$\mathcal{A}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d}})\mathbf{V} \qquad \square \qquad \mathcal{A}(\mathbf{q}_i, \mathbf{K}, \mathbf{V}) = \sum_j \frac{k(\mathbf{q}_i, \mathbf{k}_j)}{\sum_l k(\mathbf{q}_i, \mathbf{k}_l)} \mathbf{v}_j = \mathbb{E}_{p(\mathbf{k}_j | \mathbf{q}_i)}[\mathbf{v}_j]$$

Establish the measurement

Attention probability 
$$p(\mathbf{k}_{j}|\mathbf{q}_{i})$$
? Uniform probability  $q(\mathbf{k}_{j}|\mathbf{q}_{i}) = \frac{1}{L_{K}}$   
 $KL(q||p) = \sum_{j=1}^{L_{K}} \frac{1}{L_{K}} \ln \frac{1/L_{K}}{k(\mathbf{q}_{i},\mathbf{k}_{j})/\sum_{l} k(\mathbf{q}_{i},\mathbf{k}_{l})}$   
 $= \ln \sum_{l=1}^{L_{K}} e^{\frac{\mathbf{q}_{i}\mathbf{k}_{l}^{\mathsf{T}}}{\sqrt{d}}} - \frac{1}{L_{K}} \sum_{i=1}^{L_{K}} \frac{\mathbf{q}_{i}\mathbf{k}_{j}^{\mathsf{T}}}{\sqrt{d}} - \ln L_{K}$ 

Dropping the constant, we define the *i*-th query's sparsity measurement as

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

Rewrite the original self-attention into the probability formulation

$$\mathcal{A}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d}})\mathbf{V} \qquad \square \qquad \mathcal{A}(\mathbf{q}_i, \mathbf{K}, \mathbf{V}) = \sum_j \frac{k(\mathbf{q}_i, \mathbf{k}_j)}{\sum_l k(\mathbf{q}_i, \mathbf{k}_l)} \mathbf{v}_j = \mathbb{E}_{p(\mathbf{k}_j | \mathbf{q}_i)}[\mathbf{v}_j]$$

Establish the measurement

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

Define ProbSparse self-attention

Based on proposed measurement, we have the *ProbSparse* Self-attention by allowing each key only to attend to the *u* dominate queries

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

where  $\overline{\mathbf{Q}}$  is a sparse matrix of the same size of **q** and it only contains the Topu queries under the sparsity measurement  $M(\mathbf{q}, \mathbf{K})$ .

- Rewrite the original self-attention into the probability formulation
- Establish the measurement
- Define ProbSparse self-attention

The measurement have to calculate each dot-product pairs, which requires O(L<sup>2</sup>) computation. And log-sum-exp fun has the potential numerical stability issue.

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

**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 following lower and upper bounds

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

*Especially when*  $q_i \in K$ *, it also holds.* 

- Rewrite the original self-attention into the probability formulation
- Establish the measurement
- Define ProbSparse self-attention

Use an approximation to the measurement



**PROPOSITION** 1. Assuming  $\mathbf{k}_{j} \sim \mathcal{N}(\mu, \Sigma)$  and let  $\mathbf{q}\mathbf{k}_{i} = \{\frac{\mathbf{q}_{i}\mathbf{k}_{j}^{\top}}{\sqrt{d}} | j = 1, \ldots, L_{K}\}, \forall M_{m} = \max_{i} M(\mathbf{q}_{i}, \mathbf{K}) \text{ there exist } \kappa > 0 \text{ such that in that interval } \forall \mathbf{q}_{1}, \mathbf{q}_{2} \in \{\mathbf{q} | M(\mathbf{q}, \mathbf{K}) \in [M_{m}, M_{m} - \kappa)\} \text{ if } \overline{M}(\mathbf{q}_{1}, \mathbf{K}) > \overline{M}(\mathbf{q}_{2}, \mathbf{K}) \text{ and } \operatorname{Var}(\mathbf{q}\mathbf{k}_{1}) > \operatorname{Var}(\mathbf{q}\mathbf{k}_{2}), \text{ we have high probability that } M(\mathbf{q}_{1}, \mathbf{K}) > M(\mathbf{q}_{2}, \mathbf{K}).$ 

Algorithm 1 ProbSparse self-attentionInput: Tensor  $Q \in \mathbb{R}^{m \times d}$ ,  $K \in \mathbb{R}^{n \times d}$ ,  $V \in \mathbb{R}^{n \times d}$ 1: initialize: set hyperparameter  $c, u = c \ln m$  and  $U = m \ln n$ 2: randomly select U dot-product pairs from K as  $\bar{K}$ 3: set the sample score  $\bar{S} = Q\bar{K}^{\top}$ 4: compute the measurement  $M = \max(\bar{S}) - \max(\bar{S})$  by row5: set Top-u queries under M as  $\bar{Q}$ 6: set  $S_1 = \operatorname{softmax}(\bar{Q}K^{\top}/\sqrt{d}) \cdot V$ 7: set  $S_0 = \operatorname{mean}(V)$ 8: set  $S = \{S_1, S_0\}$  by their original rows accordinglyOutput: self-attention feature map S.

Sample L log L dot-product pairs  $\rightarrow$  O(L log L) O(L<sup>2</sup>) \checkmark

#### Challenge 2: Self-attention Distilling Operation





#### **Challenge 3: Generative Style Decoder**



Start token is an efficient technique in NLP "dynamic decoding", especially for per-training model, and we extend it into a generative way.

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.

Take predicting 480 points as an example (5-day prediction in the ETT dataset), we will take the known 5 days before this time sequence as "start-token", and feed the generative-style inference decoder with them.

$$\mathbf{X}_{\text{feed}\_de}^{t} = \text{Concat}(\mathbf{X}_{\text{token}}^{t}, \mathbf{X}_{0}^{t}) \in \mathbb{R}^{(L_{\text{token}}+L_{y}) \times d_{\text{model}}}$$

$$\underbrace{\mathbf{X}_{\text{feed}\_de}^{t} = \text{Concat}(\mathbf{X}_{\text{token}}^{t}, \mathbf{X}_{0}^{t}) \in \mathbb{R}^{(L_{\text{token}}+L_{y}) \times d_{\text{model}}}$$

$$\underbrace{\mathbf{Long Output}}_{t}$$

## model

#### **Challenges:** [Complexity/layer ✓ ][Long Input ✓ ][Long Output ✓ ]



An overall graph of the Informer model.

The Informer network components in details

		Ν							
1x3 Conv1d	Embedding ( $d = 512$ )								
Multi-head ProbSparse Attention ( $h = 16, d = 32$ )									
Add, LayerNorm, Dropout ( $p = 0.1$ )									
Pos-wise FFN ( $d_{inner} = 2048$ ), GELU									
Add, LayerN	orm, Dropout $(p = 0.1)$	1							
1x3	conv1d, ELU	1							
Max pooling (stride = $2$ )									
		N							
1x3 Conv1d Embedding ( $d = 512$ )									
add Mask on Attention Block									
Multi-head Attention ( $h = 8, d = 64$ )									
Add, LayerNorm, Dropout $(p = 0.1)$									
Pos-wise FFN ( $d_{inner} = 2048$ ), GELU									
Add, LayerNorm, Dropout ( $p = 0.1$ )									
Outputs FCN $(d = d_{out})$									
	1x3 Conv1d Multi-head ProbSp Add, LayerN Pos-wise FFN Add, LayerN 1x3 Max po 1x3 Conv1d add Mask Multi-head A Add, LayerN Pos-wise FFN Add, LayerN	1x3 Conv1dEmbedding $(d = 512)$ Multi-head ProbSparse Attention $(h = 16, d = 32)$ Add, LayerNorm, Dropout $(p = 0.1)$ Pos-wise FFN $(d_{inner} = 2048)$ , GELUAdd, LayerNorm, Dropout $(p = 0.1)$ 1x3 conv1d, ELUMax pooling (stride = 2)1x3 Conv1dEmbedding $(d = 512)$ add Mask on Attention BlockMulti-head Attention $(h = 8, d = 64)$ Add, LayerNorm, Dropout $(p = 0.1)$ Pos-wise FFN $(d_{inner} = 2048)$ , GELUAdd, LayerNorm, Dropout $(p = 0.1)$ Pos-wise FFN $(d_{inner} = 2048)$ , GELUAdd, LayerNorm, Dropout $(p = 0.1)$ FCN $(d = d_{out})$							

## CONTENT

Background
Motivation: why self-attention?
Methods: the details of Informer
Experiments
Summary

#### **Experiments settings**

#### Datasets

- ETT (Electricity Transformer Temperature)
  - 2 stations, 2 years, every 15 minutes.
- ECL (Electricity Consuming Load)
  - 1 country, 2 years, every 1 hour.
- Weather

◆ 1600 US locations, 4 years, every 1 hour. Train/Val/Test is 28/10/10 months.

#### **Baselines**

- Time-series methods: ARIMA, DeepAR, Prophet, LSTMa, LSTnet
- □ Transformer-based methods: Vanilla Transformer, Reformer, LogSparse Transformer

#### Metrics

Mean Absolute Error (MAE), Mean Squared Error (MSE)

#### Platform

All methods are run on one single Nvidia V100 GPU.

```
Train/Val/Test is 12/4/8 months.
```

```
Train/Val/Test is 15/3/6 months.
```

#### **Overall Results: Univariate Time-series Forecasting**

Mathada	Matria			ETTh <sub>1</sub>				2507	ETTh <sub>2</sub>	ē				ETTm					Weathe	r				ECL	2		
Methods	Wietric	24	48	168	336	720	24	48	168	336	720	24	48	96	288	672	24	48	168	336	720	48	168	336	720	960	count
Informer	MSE MAE	0.062	<b>0.108</b> 0.245	0.146 0.294	0.208 0.363	<b>0.193</b> 0.365	0.079	0.103 0.240	0.143 0.296	0.171 0.327	0.184 0.339	0.051 0.153	0.092 0.217	0.119 0.249	0.181 0.320	0.204 0.345	0.107	0.164 0.282	0.226 0.338	0.241 0.352	0.259 0.367	0.335 0.423	0.408 0.466	0.451 0.488	0.466 0.499	0.470 0.520	28
Informer <sup>†</sup>	MSE MAE	0.046	0.129 0.274	0.183 0.337	0.189 0.346	0.201 0.357	0.083	0.111 0.249	0.154 0.306	0.166 0.323	0.181 0.338	0.054	0.087 0.210	0.115 0.248	0.182 0.323	0.207 0.353	0.107 0.220	0.167 0.284	0.237 0.352	0.252 0.366	0.263 0.374	0.304 0.404	0.416 0.478	0.479 0.508	0.482 0.515	0.538 0.560	14
LogTrans	MSE MAE	0.059	0.111 0.263	0.155 0.309	0.196 0.370	0.217 0.379	0.080	0.107 0.262	0.176 0.344	0.175 0.345	0.185 0.349	0.061	0.156 0.322	0.229 0.397	0.362 0.512	0.450 0.582	0.120	0.182 0.312	0.267 0.387	0.299 0.416	0.274 0.387	0.360 0.455	0.410 0.481	0.482 0.521	0.522 0.551	0.546 0.563	0
Reformer	MSE MAE	0.172 0.319	0.228 0.395	1.460 1.089	1.728 0.978	1.948 1.226	0.235	0.434 0.505	0.961 0.797	1.532 1.060	1.862 1.543	0.055	0.229 0.340	0.854 0.675	0.962 1.107	1.605 1.312	0.197	0.268 0.381	0.590 0.552	1.692 0.945	1.887 1.352	0.917 0.840	1.635 1.515	3.448 2.088	4.745 3.913	6.841 4.913	0
LSTMa	MSE MAE	0.094	0.175 0.322	0.210 0.352	0.556 0.644	0.635 0.704	0.135	0.172 0.318	0.359 0.470	0.516 0.548	0.562 0.613	0.099	0.289 0.371	0.255 0.370	0.480 0.528	0.988 0.805	0.107 0.222	0.166 0.298	0.305 0.404	0.404 0.476	0.784 0.709	0.475 0.509	0.703 0.617	1.186 0.854	1.473 0.910	1.493 0.9260	1
DeepAR	MSE MAE	0.089	0.126 0.291	0.213 0.382	0.403 0.496	0.614 0.643	0.080	0.125 0.283	0.179 0.346	0.568 0.555	0.367 0.488	0.075	0.197 0.332	0.336 0.450	0.908 0.739	2.371 1.256	0.108	0.177 0.313	0.259 0.397	0.535 0.580	0.407 0.506	0.188 0.317	0.295 0.398	0.388 0.471	0.471 0.507	0.583 0.583	6
ARIMA	MSE MAE	0.086	0.133 0.242	0.364 0.456	0.428 0.537	0.613 0.684	3.538 0.407	3.168 0.440	2.768 0.555	2.717 0.680	2.822 0.952	0.074	0.157 0.274	0.242 0.357	0.424 0.500	0.565 0.605	0.199	0.247 0.375	0.471 0.541	0.678 0.666	0.996 0.853	0.861 0.726	1.014 0.797	1.102 0.834	1.213 0.883	1.322 0.908	1
Prophet	MSE MAE	0.093	0.150 0.300	1.194 0.721	1.509 1.766	2.685 3.155	0.179	0.284 0.428	2.113 1.018	2.052 2.487	3.287 4.592	0.102	0.117 0.273	0.146 0.304	0.414 0.482	2.671 1.112	0.280	0.421 0.492	2.409 1.092	1.931 2.406	3.759 1.030	0.506 0.557	2.711 1.239	2.220 3.029	4.201 1.363	6.827 4.184	0

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

- The *Informer* greatly improves the inference performance.
- The *Informer* and its canonical degradation *Informer+* show comparable performance and Informer beats *Informer+* mostly in wining-counts, i.e. 28>14.
- The *Informer* model shows significantly better results than time-series models.
- Our proposed method achieves better results than vTrans, Reformer and LogSparseTrans.

#### **Overall Results: Multivariate Time-series Forecasting**

			lable	2. r	viuit	Ivan	aten	ong	sequ	ence	unne	2-501	les i	orec	asum	gre	suits	0111	our	uatas	sets (	inve	case	5)			
Mathada	Matria	ETTh <sub>1</sub>					ETTh <sub>2</sub>	5			ETTm <sub>1</sub>			Weather					ECL				agunt				
Methods	Metric	24	48	168	336	720	24	48	168	336	720	24	48	96	288	672	24	48	168	336	720	48	168	336	720	960	count
Informer	MSE MAE	0.509	0.551	0.878 0.722	0.884	0.941 0.768	0.446	0.934	1.512 0.996	1.665	2.340	0.325	0.472	<b>0.642</b> 0.626	1.219 0.871	1.651 1.002	0.353	0.464	0.592 0.531	0.623	0.685	0.269	0.300	0.311	0.308	0.328	32
Informer <sup>†</sup>	MSE MAE	0.550	0.602 0.581	0.893 0.733	0.836 0.729	0.981 0.779	0.683	0.977 0.793	1.873 1.094	1.374 0.935	2.493 1.253	0.324 0.440	0.446 0.508	0.651 <b>0.616</b>	1.342 0.927	1.661 1.001	0.355 0.383	0.471 0.456	0.613 0.544	0.626 0.548	0.680 0.569	0.269 0.351	0.281 0.366	0.309 0.383	0.314 0.388	0.356 0.394	12
LogTrans	MSE MAE	0.656	0.670 0.611	$\begin{array}{c} 0.888\\ 0.766\end{array}$	0.942 0.766	1.109 0.843	0.726	$\begin{array}{c} 1.728\\ 0.944\end{array}$	3.944 1.573	3.711 1.587	2.817 1.356	0.341 0.495	0.495 0.527	0.674 0.674	1.728 1.656	1.865 1.721	0.365 0.405	0.496 0.485	0.649 0.573	$0.666 \\ 0.584$	$\begin{array}{c} 0.741 \\ 0.611 \end{array}$	<b>0.267</b> 0.366	0.290 0.382	<b>0.305</b> 0.395	0.311 0.397	0.333 0.413	2
Reformer	MSE MAE	0.887 0.630	1.159 0.750	1.686 0.996	1.919 1.090	2.177 1.218	1.381 1.475	1.715 1.585	4.484 1.650	3.798 1.508	5.111 1.793	0.598 0.489	0.952 0.645	1.267 0.795	$\begin{array}{c} 1.632\\ 0.886 \end{array}$	1.943 1.006	0.583 0.497	0.633 0.556	1.228 0.763	1.770 0.997	$\begin{array}{c} 2.548 \\ 1.407 \end{array}$	$\begin{array}{c} 1.312\\ 0.911\end{array}$	1.453 0.975	1.507 0.978	1.883 1.002	1.973 1.185	0
LSTMa	MSE MAE	0.536	0.616 0.577	1.058 0.725	1.152 0.794	1.682 1.018	1.049 0.689	$\begin{array}{c} 1.331\\ 0.805 \end{array}$	3.987 1.560	3.276 1.375	3.711 1.520	0.511 0.517	1.280 0.819	1.195 0.785	1.598 0.952	2.530 1.259	0.476 0.464	0.763 0.589	0.948 0.713	1.497 0.889	1.314 0.875	0.388 0.444	0.492 0.498	0.778 0.629	1.528 0.945	1.343 0.886	0
LSTnet	MSE MAE	1.175 0.793	1.344 0.864	1.865 1.092	2.477 1.193	1.925 1.084	2.632	3.487 1.577	1.442 2.389	<b>1.372</b> 2.429	2.403 3.403	1.856 1.058	1.909 1.085	2.654 1.378	<b>1.009</b> 1.902	1.681 2.701	0.575 0.507	0.622 0.553	0.676 0.585	0.714 0.607	0.773 0.643	0.279 <b>0.337</b>	0.318 0.368	0.357 0.391	0.442 0.433	0.473 0.443	4

Table 2: Multivariate lang acquance time caries forecasting results on four detects (fue acces)

- The *Informer* greatly improves the inference performance.
- The Informer and its canonical degradation Informer+ show comparable performance and Informer beats Informer+ mostly in wining-counts, i.e. 32>12.
- The previous conclusion in multivariate case holds.

#### **Overall Results: LSTF with Granularity Consideration**



#### Parameter Sensitivity



#### **Ablation study**

Prediction ler	ngth		336		720					
Encoder's inp	out	336	720	1440	720	1440	2880			
Informer	MSE	0.243	0.225	0.212	0.258	0.238	0.224			
	MAE	0.487	0.404	0.381	0.503	0.399	0.387			
T	MSE	0.214	0.205	-	0.235	-	-			
mormer	MAE	0.369	0.364	-	0.401	-	-			
LogTrans	MSE	0.256	0.233	-	0.264	-	-			
Log Hans	MAE	0.496	0.412	-	0.523	-	-			
Pafarmar	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			

Table 3: Ablation of *ProbSparse* mechanism

<sup>1</sup> Informer<sup>†</sup> uses the canonical self-attention mechanism. <sup>2</sup> The '-' indicates failure for out-of-memory.

Prediction	length			336					480		
Encoder's input		336	480	720	960	1200	336	480	720	960	1200
T.C. t	MSE	0.201	0.175	0.215	0.185	0.172	0.136	0.213	0.178	0.146	0.121
Informer	MAE	0.360	0.335	0.366	0.355	0.321	0.282	0.382	0.345	0.296	0.272
Informer <sup>‡</sup>	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	-	-

#### Table 4: Ablation of Self-attention Distilling

<sup>1</sup> Informer<sup> $\ddagger$ </sup> removes the self-attention distilling from Informer<sup> $\dagger$ </sup>.

<sup>2</sup> The '-' indicates failure for out-of-memory.

#### **Ablation study**

Ta	ble 5:	Abla	tion o	of Gen	nerativ	ve Sty	le De	coder				
Prediction le	ength		33	36		480						
Prediction of	offset	+0	+12	+24	+48	+0	+48	+96	+168			
t	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			
Informar§	MSE	0.152	-	-	-	0.462	-	-	-			
mormer	MAE	0.294	1. <del></del>	-	-	0.595	-	-	-			

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

#### **Computation efficiency**



The total runtime of training/testing phase.

## CONTENT

Background
Motivation: why self-attention?
Methods: the details of Informer
Experiments
Summary

#### Things to take

Long Sequence Forecasting problem is a long-standing problem.
The Self-attention model is a tractable solution in time-series problem.
ProbSparse Self-attention -> an efficient self-attention, also probabilistic/analyzable
Self-attention Distilling Operation -> reduces overheads, allows for longer inputs
Generative Style Decoder -> allows for longer outputs, even arbitrary-step predictions

Sparsity hypothesis works on Self-attention remarkably.

Our proposed methods can bring substantial benefits to other domains, such as long sequence generation of text, music, image, and video.

## problem

- Paper: https://arxiv.org/abs/2012.07436
  - Code: <a href="https://github.com/zhouhaoyi/Informer2020">https://github.com/zhouhaoyi/Informer2020</a>
  - Dataset: https://github.com/zhouhaoyi/ETDataset

#### **Electricity Transformer Dataset (ETDataset)**



In this Github repo, we provide several datasets could be used for the long sequence time-series problem. All datasets have been preprocessed and they were stored as .csv files. The dataset ranges from 2016/07 to 2018/07, and we will update to 2019 soon. 中文版本 | ChineseVersion

#### Dataset list

- ETT-small: The data of 2 Electricity Transformers at 2 stations, including load, oil temperature.
- ETT-large: The data of 39 Electricity Transformers at 39 stations, including load, oil temperature.
- ETT-full: The data of 69 Transformer station at 39 stations, including load, oil temperature, location, climate, demand.





北航ACTBD公众号



Figure 1. The overall view of "OT" in the ETT-small. Figure 2. The autocorrelation graph of all variables.

国际人工智能会议 AAAI 2021论文北京预讲会

# THANKS

2020.12.19