

# The sequential architecture & AI for scientific computing

Haoyi Zhou

Thursday, Feb 27, 2025



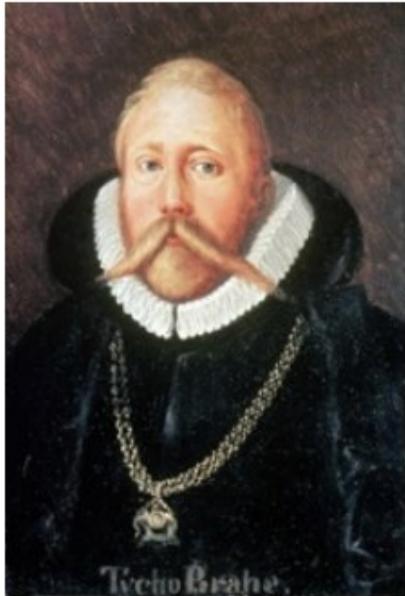
Beihang University

# - Content -

- **Sequential Modeling: from Science**
  - The Transformer models
  - The modern architecture
  - Some applications on scientific sequences
  - Future direction

# Tycho → Kepler → Newton

Three stages of transformation including data collection, model learning and model interpretation, which align closely with the development trajectory of AI.



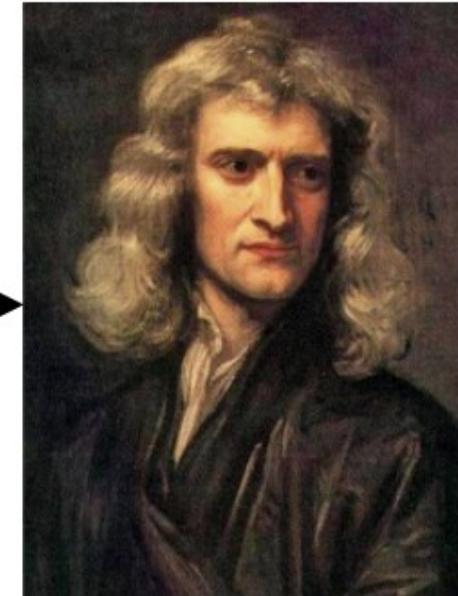
Tycho Brahe (1546-1610)

Observation  
Data Collection



Johannes Kepler (1571-1630)

Analyzation  
Model Learning



Isaac Newton (1643-1727)

Explanation  
Model Interpretation

# Outlook on Scientific Research Paradigm in 2050



Envisioning Science in 2050

Report of a Community of Interest Workshop  
on Future Scientific Methodologies



US Department of Energy National  
Laboratory Advanced Scientific  
Computing Research Group  
Reported in 2022, looking forward to  
the future form of scientific research in  
2050

## The Theory Machine Transforms Science

Intelligently complete experimental data analysis, theoretical derivation, and simulation calculations

## Delphi: The universal knowledge map

Intelligently summarize and extract scientific literature and experimental reports, help scientists analyze and predict research directions, and propose new scientific hypotheses

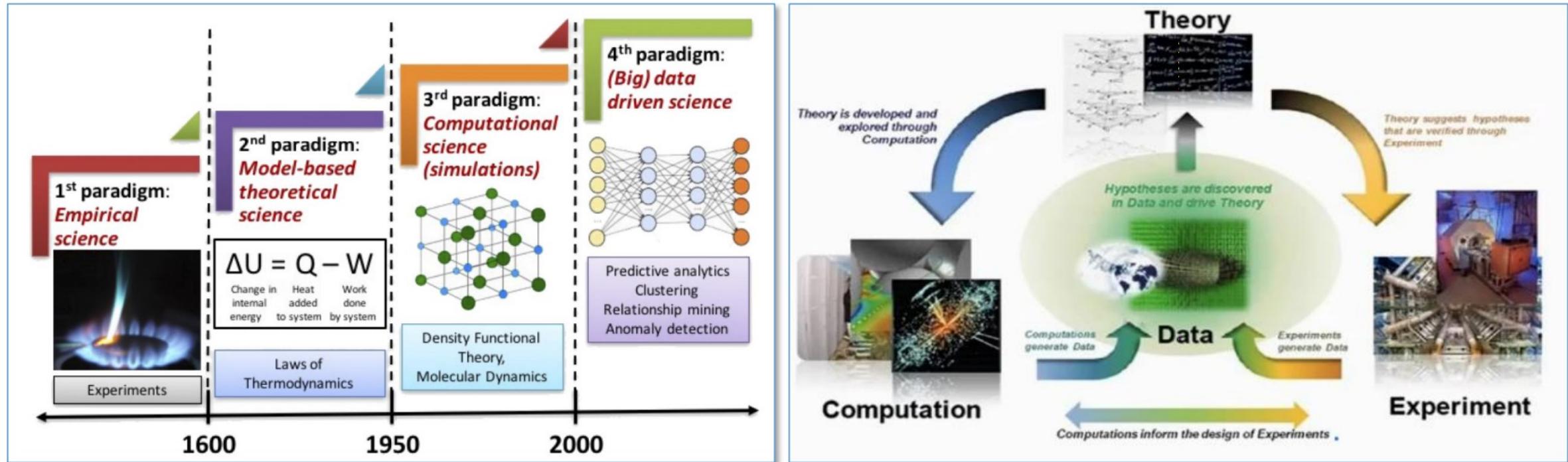
## Smart self-driving facilities automate experimentation

Like autonomous driving, independently completing scientific experiments and becoming an intelligent assistant for experimental scientists

## Universal Data Service Make All Data Easily Accessible

## Linked Facilities Form a Discovery Cloud

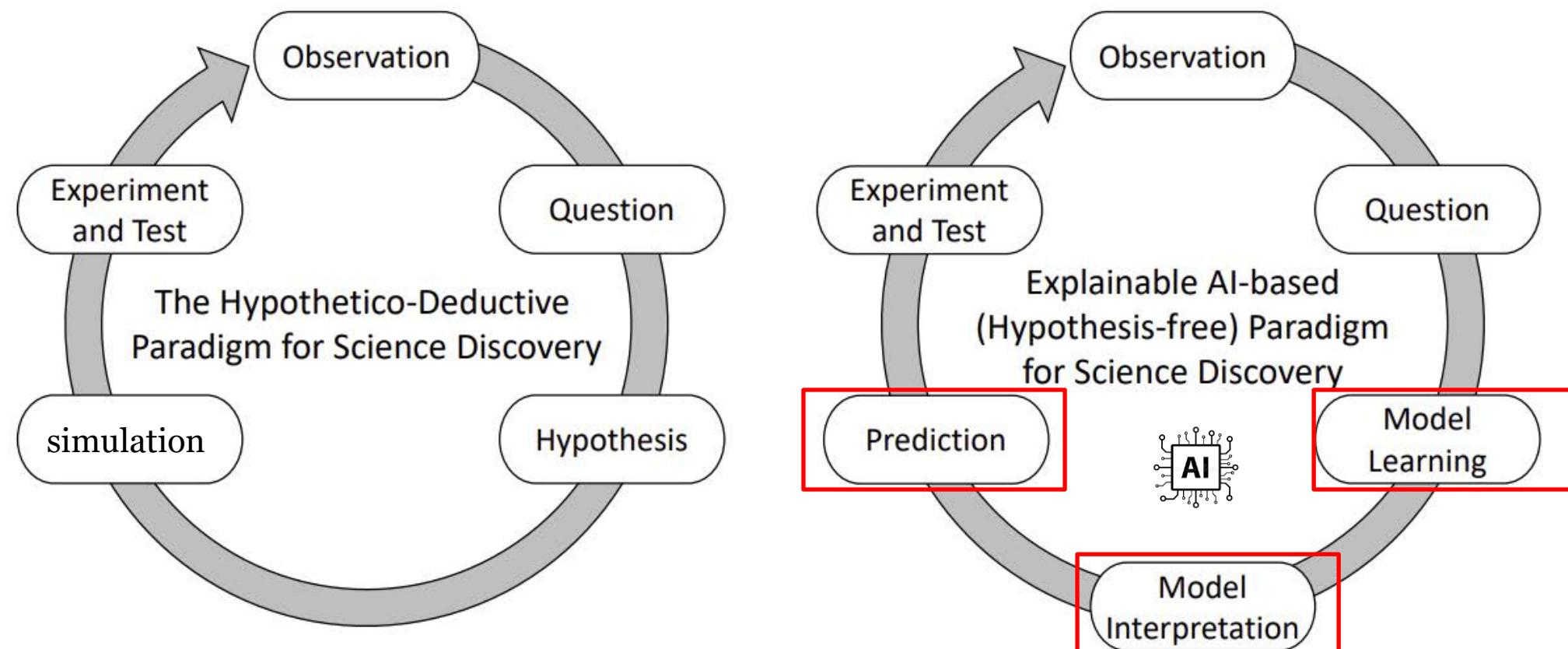
# AI is becoming a New Tool for Scientific Research



Experimental Science → Theoretical Science → Computing Science → AI pushes scientific research into the paradigm 4

# Hypothesis-free paradigm for Science Discovery

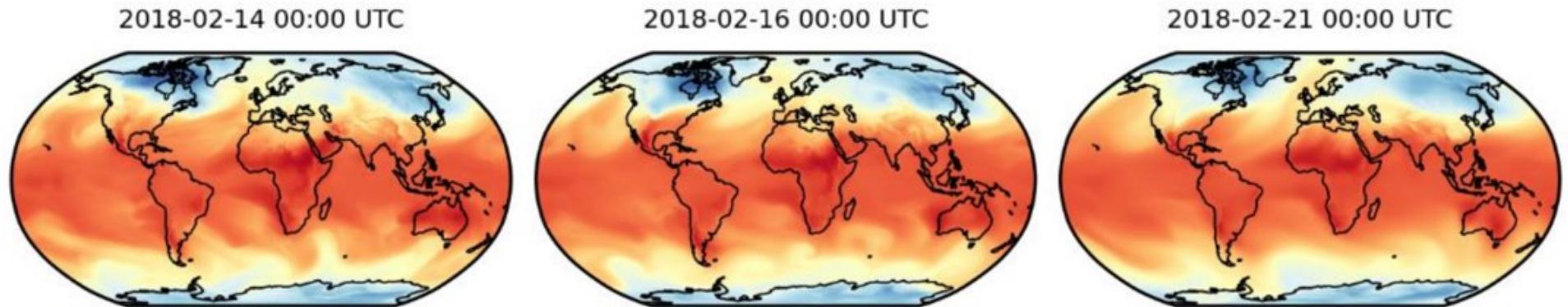
The hypothesis and simulation stage are replaced by AI models, which can learn underlying pattern in the observed data and directly output possible physical interpretation and prediction.



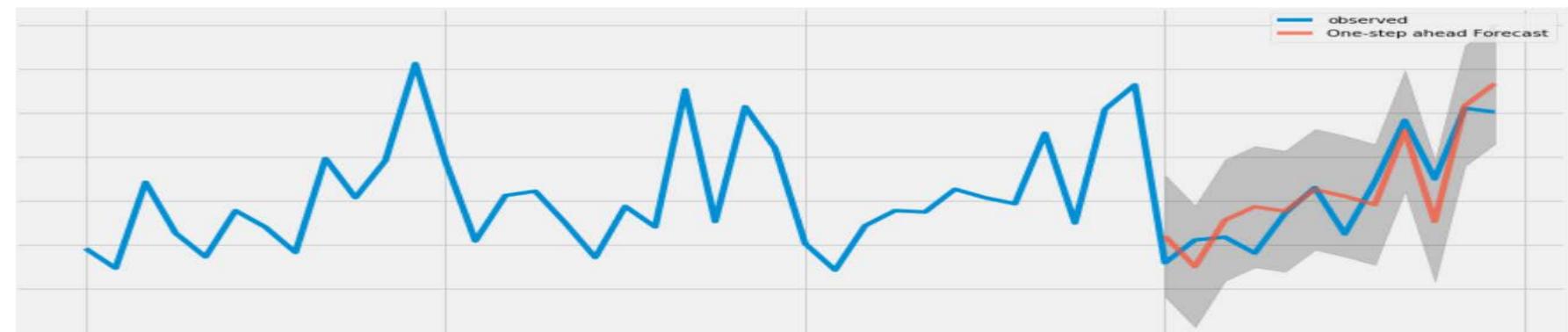
# Time-series and dynamic system

Physical dynamical systems are essentially time series controlled by explicit equations

Dynamic system



Time-series



# Sequential Data: The Requirement for Computation

Intelligent computing of **sequential data** is a crucial for the support of industrial development.

A320: 12,000 parameters, A350: nearly 600,000 parameters.

Boeing 787: One flight generates 0.5TB of data (including only cabin pressure, etc.).

## Social Activities



Behavioral data



Smart cities

Complex sensor data

## Intelligent Manufacturing



Aerospace



Robotic automation

Continuous high-frequency data

## Infrastructure



Power and energy industry



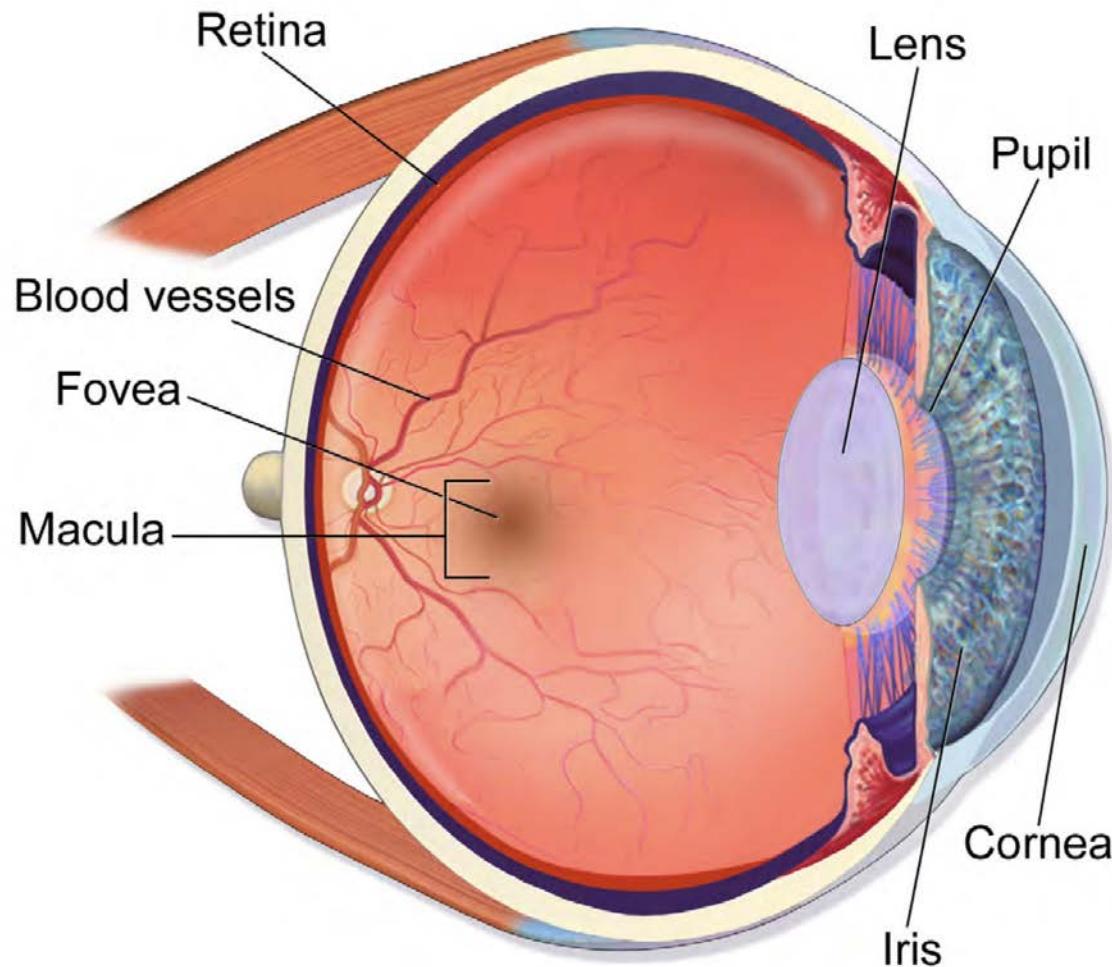
UAV inspection

Working condition data

# - Content -

- Sequential Modeling: from Science
- **The Transformer Model**
- The modern architecture
- Some applications on scientific sequences
- Future direction

# What is the attention?



\* Images are from google search engine.

# How do we tackle the forecasting problem?

**Problem Setting:** LSTF predicts the long sequence output from massive long sequence inputs.

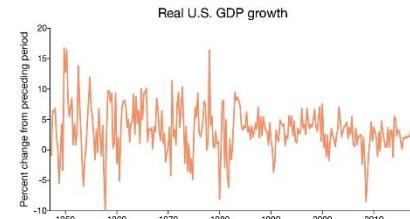
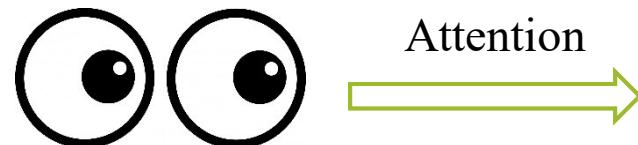
How do we plan gifts for Christmas day?



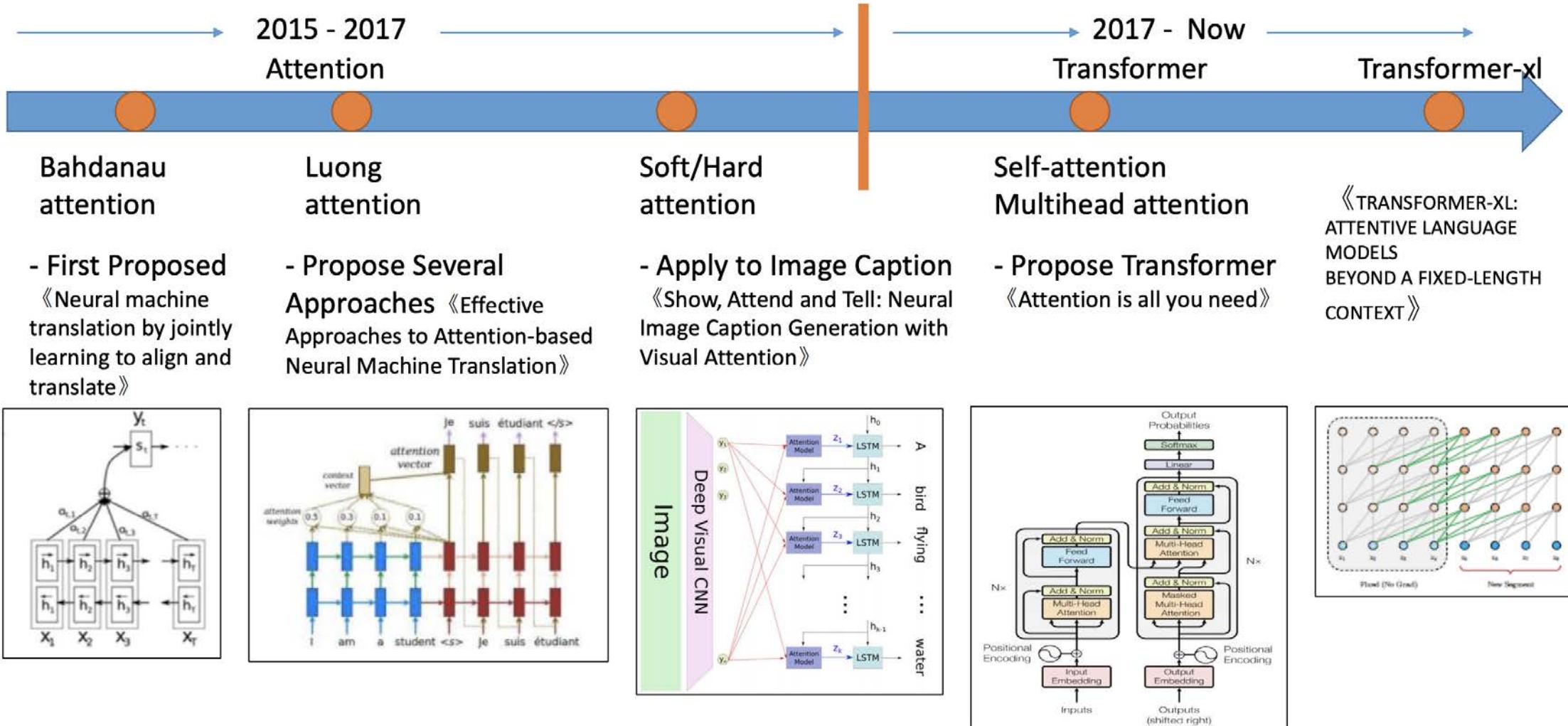
We will have some feelings of what you have seen when it comes to the prediction.

---

**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.



# The development of attention models



Figures are from Aliyun slides.

# Take the NLP as an example

<https://jalammar.github.io/illustrated-transformer/>

# ① Embeddings: Numbers more than Words

Token embeddings

---

# ① Embeddings: Numbers more than Words

Position embeddings

---

# ① Embeddings: Numbers more than Words

$$\begin{array}{c} \text{Token embedding matrix} \\ T \times C \\ \begin{matrix} \text{the} & \begin{bmatrix} 0.07 & 0.13 & 0.63 & \dots & 0.23 \end{bmatrix} \\ \text{robots} & \begin{bmatrix} 0.81 & 0.51 & 0.44 & \dots & 0.98 \end{bmatrix} \\ \text{will} & \begin{bmatrix} 0.62 & 0.29 & 0.80 & \dots & 0.34 \end{bmatrix} \\ \text{bring} & \begin{bmatrix} 0.50 & 0.16 & 0.93 & \dots & 0.19 \end{bmatrix} \end{matrix} \end{array} + \begin{array}{c} \text{Position embedding matrix} \\ T \times C \\ \begin{matrix} 0 & \begin{bmatrix} 0.32 & 0.12 & 0.88 & \dots & 0.75 \end{bmatrix} \\ 1 & \begin{bmatrix} 0.14 & 0.07 & 0.41 & \dots & 0.35 \end{bmatrix} \\ 2 & \begin{bmatrix} 0.07 & 0.47 & 0.23 & \dots & 0.86 \end{bmatrix} \\ 3 & \begin{bmatrix} 0.81 & 0.12 & 0.06 & \dots & 0.24 \end{bmatrix} \end{matrix} \end{array}$$

## ② Queries, keys and values

Query, key, value vectors for each word

---

$$\begin{array}{l} \text{the} \\ \text{robots} \\ \text{will} \\ \text{bring} \end{array} \begin{matrix} & T \times C \\ \left[ \begin{array}{cccc} 0.39 & 0.25 & \dots & 0.98 \\ 0.95 & 0.58 & \dots & 1.33 \\ 0.69 & 0.77 & \dots & 1.20 \\ 1.31 & 0.28 & \dots & 0.43 \end{array} \right] & \times & \left[ \begin{array}{cccc} 0.02 & 0.21 & \dots & 0.37 \\ 0.03 & 0.02 & \dots & 0.05 \\ 0.29 & 0.07 & \dots & 0.99 \\ 0.38 & 0.37 & \dots & 0.28 \\ \vdots & \vdots & \vdots & \vdots \end{array} \right] & = & W_q (C \times C) \end{matrix}$$

# ③ “Pay” Attention

## Self Attention

$$\begin{array}{ll} & Q_1 (T \times H) \\ \text{the} & \begin{bmatrix} 0.07 & 0.04 & \dots & 0.64 \end{bmatrix} \\ \text{robots} & \begin{bmatrix} 0.07 & 0.07 & \dots & 0.36 \end{bmatrix} \\ \text{will} & \begin{bmatrix} 0.08 & 0.06 & \dots & 0.44 \end{bmatrix} \\ \text{bring} & \begin{bmatrix} 0.04 & 0.04 & \dots & 0.34 \end{bmatrix} \end{array} \times \begin{array}{ll} & K_1 (T \times H) \\ \text{the} & \begin{bmatrix} 0.22 & 0.13 & \dots & 0.10 \end{bmatrix} \\ \text{robots} & \begin{bmatrix} 0.22 & 0.03 & \dots & 0.10 \end{bmatrix} \\ \text{will} & \begin{bmatrix} 0.20 & 0.04 & \dots & 0.11 \end{bmatrix} \\ \text{bring} & \begin{bmatrix} 0.19 & 0.06 & \dots & 0.24 \end{bmatrix} \end{array}$$

# ④ “Apply” Attention

## Self Attention

$$A_1 \ (T \times T)$$

the	0.18	0.22	0.23	0.37
robots	0.18	0.29	0.30	0.24
will	0.18	0.26	0.28	0.28
bring	0.17	0.30	0.31	0.22

the      robots      will      bring

# ⑤ Two heads are better than one

Split Q, K and V matrices into “heads”

---

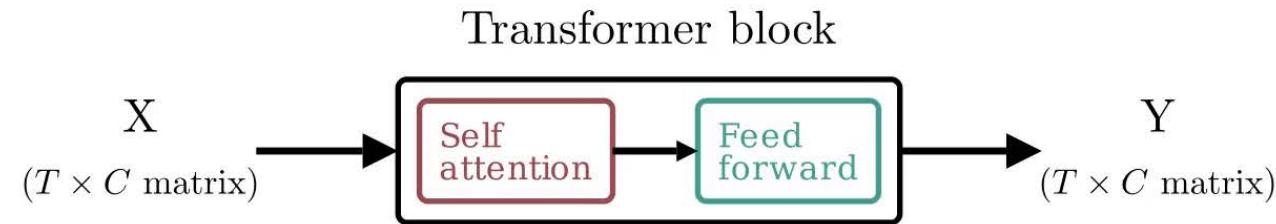
$Q (T \times C)$

$$\begin{bmatrix} 0.07 & 0.04 & \dots & 0.64 \\ 0.07 & 0.07 & \dots & 0.36 \\ 0.08 & 0.06 & \dots & 0.44 \\ 0.04 & 0.04 & \dots & 0.34 \end{bmatrix}$$

# ⑥ Putting all heads together

# ⑦ Stacking the attention block

A Transformer block



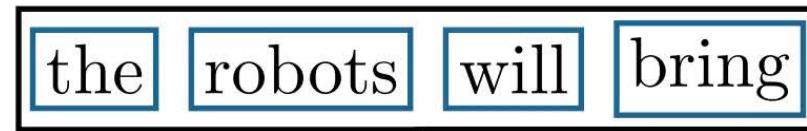
# ⑧ Making the prediction

# ⑨ Text Generator & Window

Generating text

---

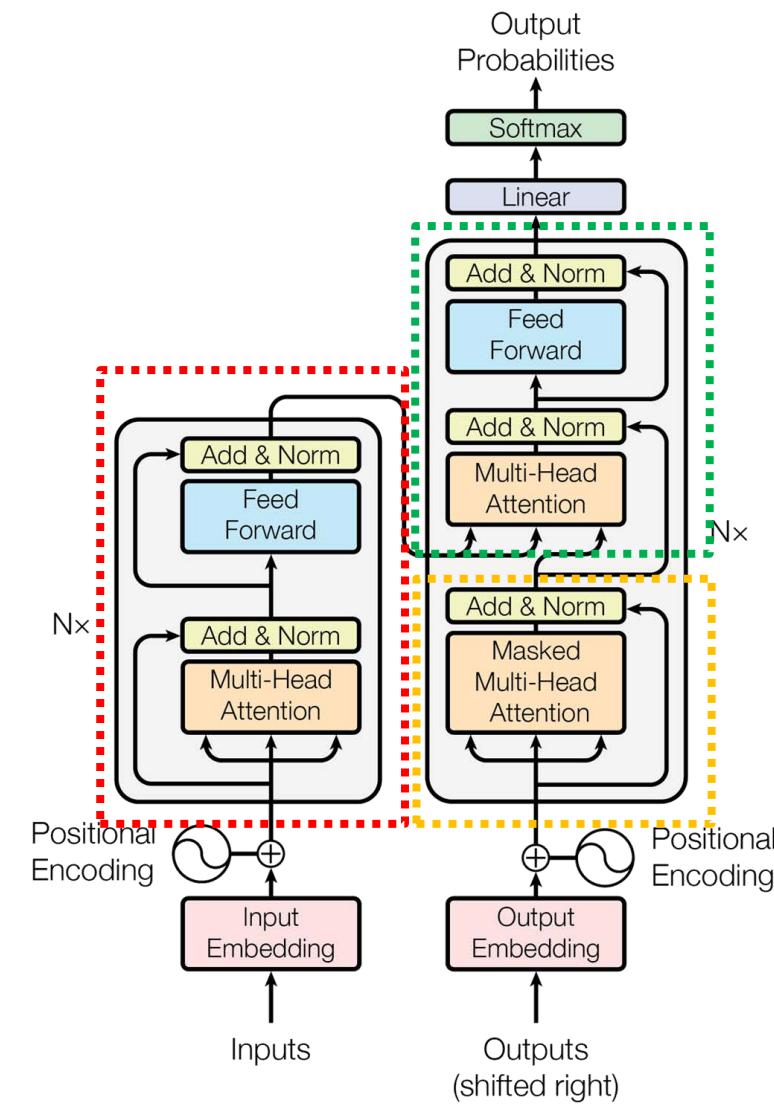
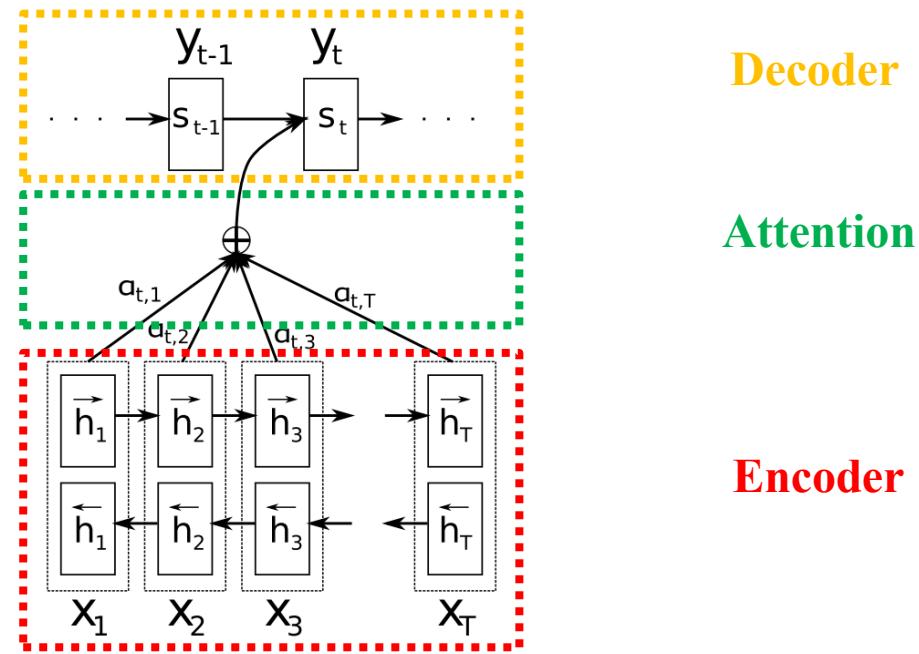
Max prompt length = 5



Transformer

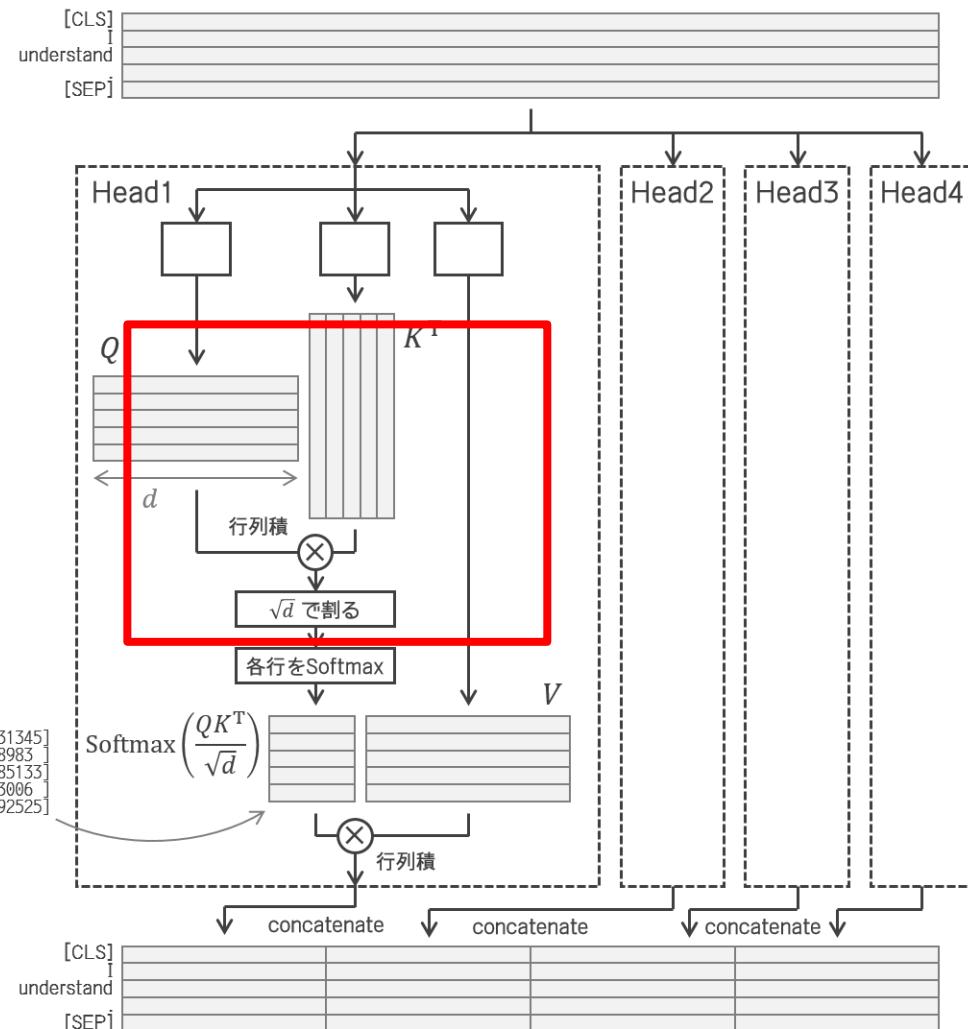
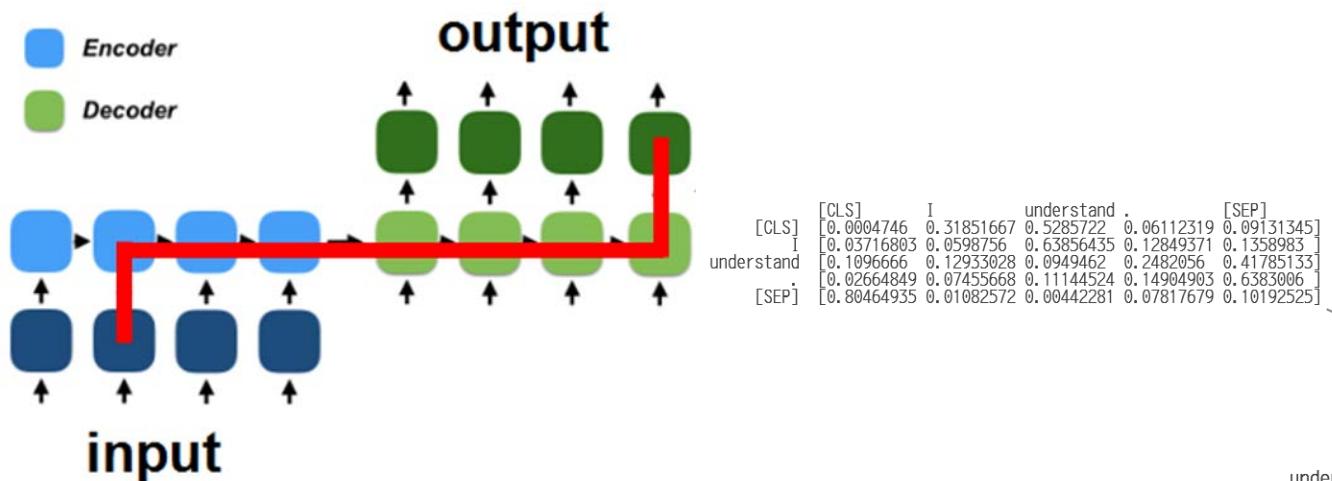
# The Transformer model

1. Encoder-Decoder
2. Attention
  - Scaled Dot-Product Attention
  - Multi-Head Attention
3. Position-wise Feed-Forward Networks
4. Positional Encoding



# The Transformer model

1. Encoder-Decoder
2. Attention
  - Scaled Dot-Product Attention
  - Multi-Head Attention
3. Position-wise Feed-Forward Networks
4. Positional Encoding



<https://cookie-box.hatenablog.com/entry/2021/02/11/195822>

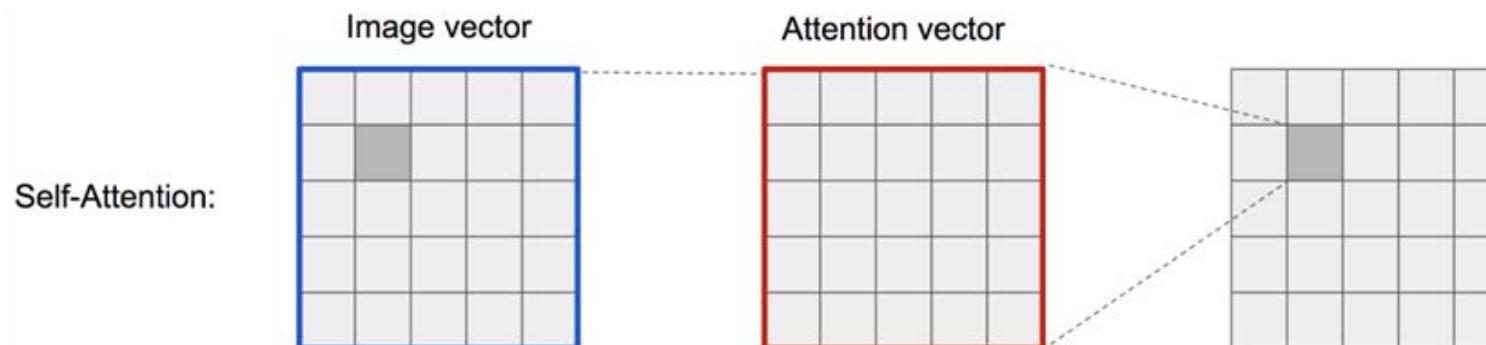
Multi-Headed Self-Attention

# The self-attention in NLP/CV field

The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .  
The FBI is chasing a criminal on the run .

$$\begin{array}{l} \mathbf{x} \times \mathbf{W}^Q = \mathbf{Q} \\ \mathbf{x} \times \mathbf{W}^K = \mathbf{K} \\ \mathbf{x} \times \mathbf{W}^V = \mathbf{V} \end{array}$$

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

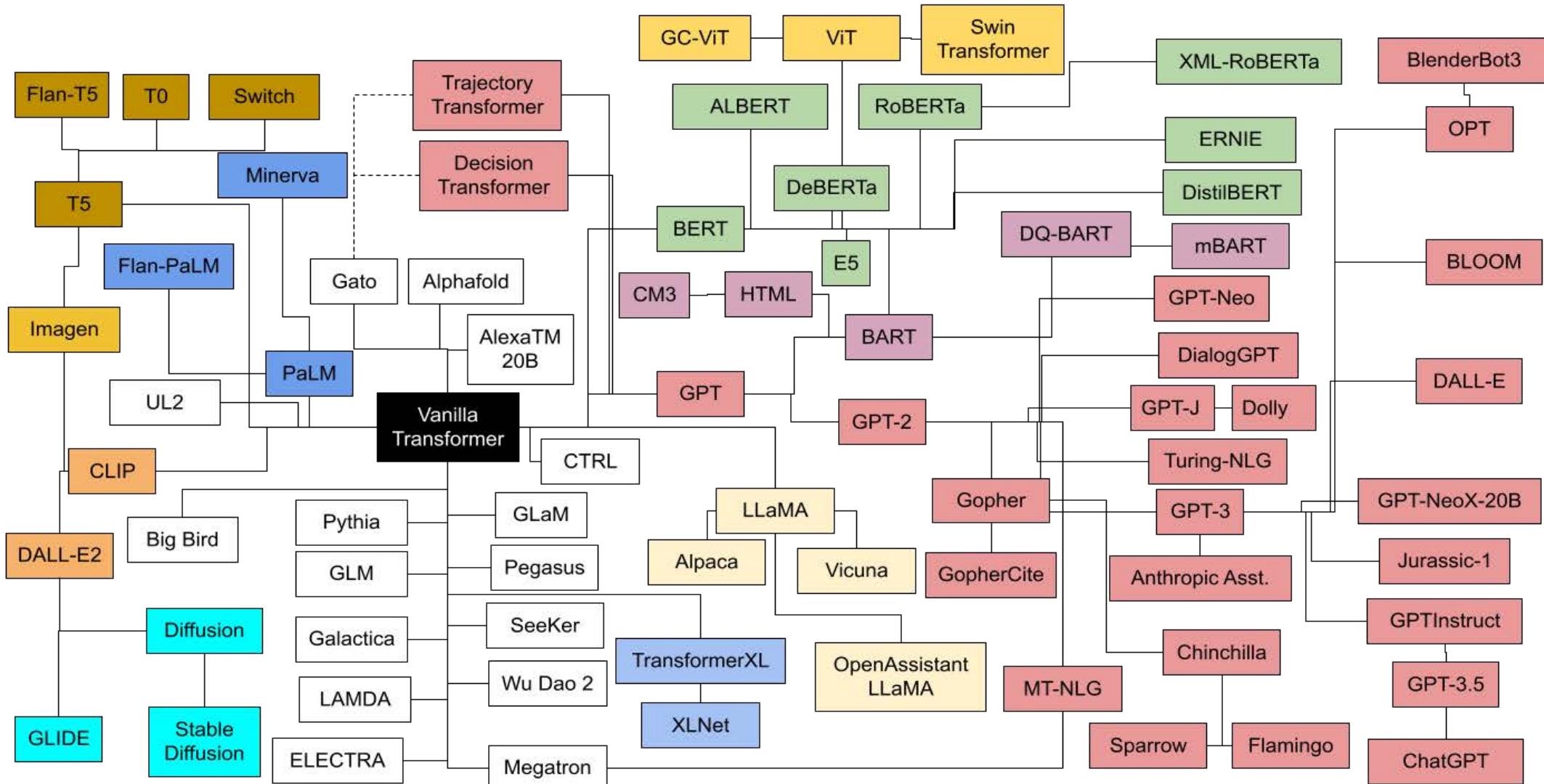


Acknowledgement from <http://jalammar.github.io/illustrated-bert/>

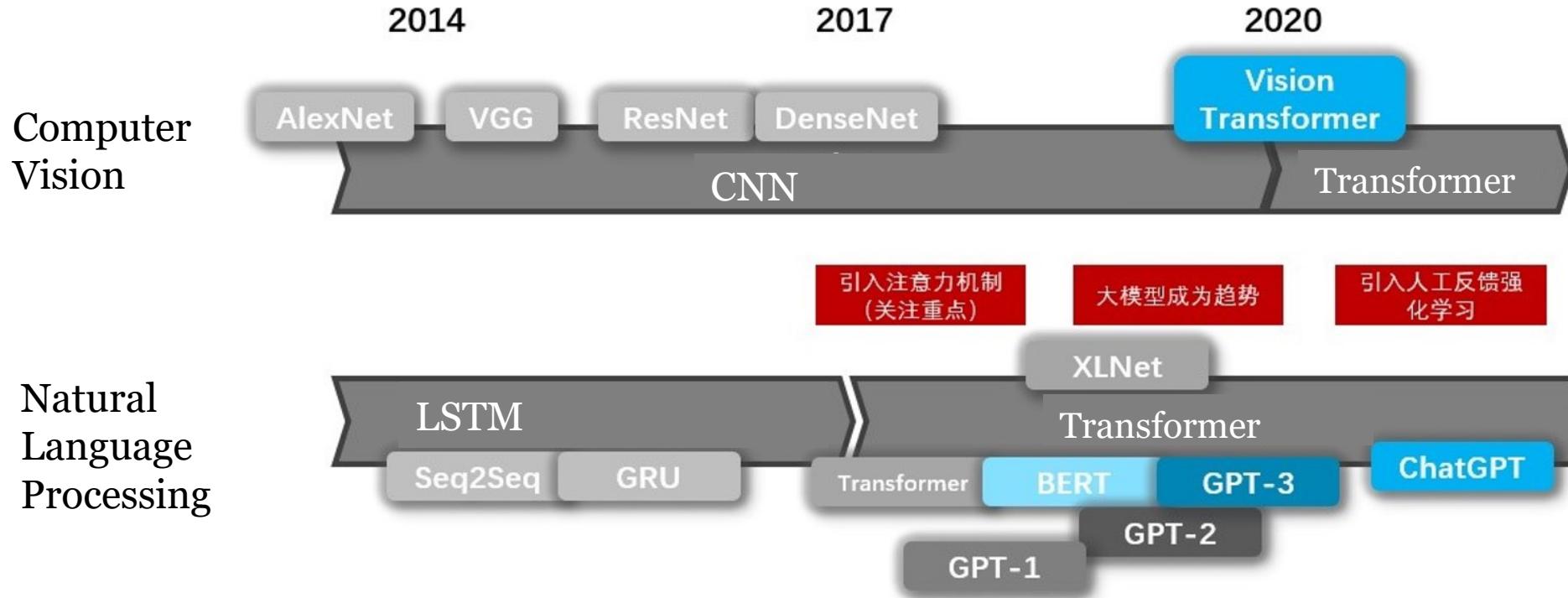
# - Content -

- Sequential Modeling: from Science
- The Transformer Model
- **The modern architecture**
- Some applications on scientific sequences
- Future direction

# The Development of Transformer

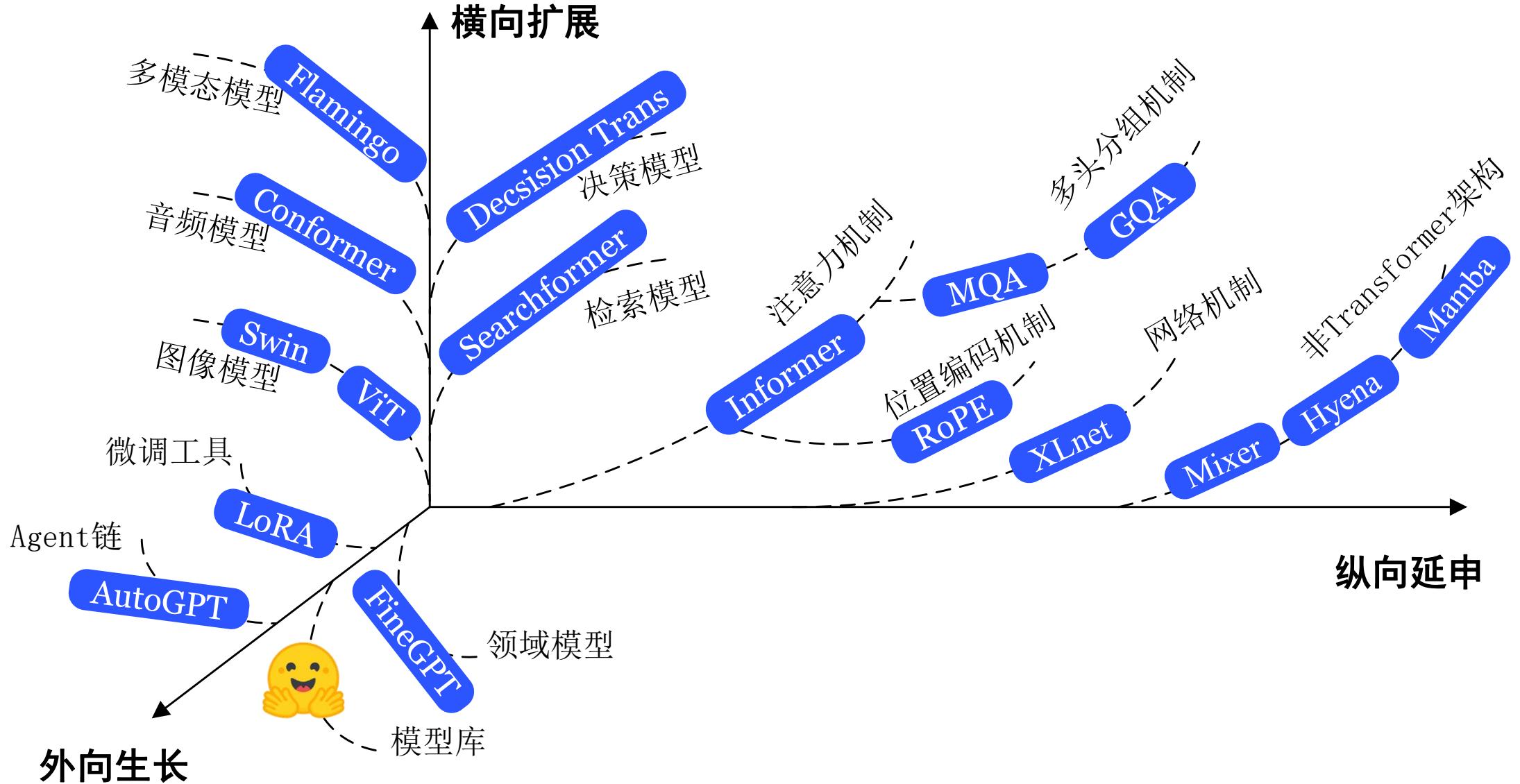


# The Development of Transformer



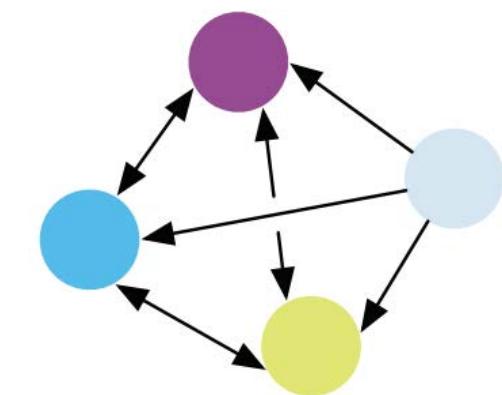
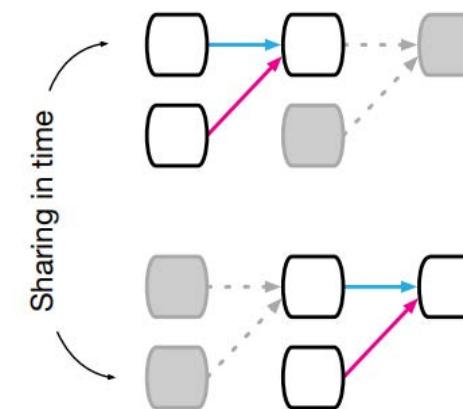
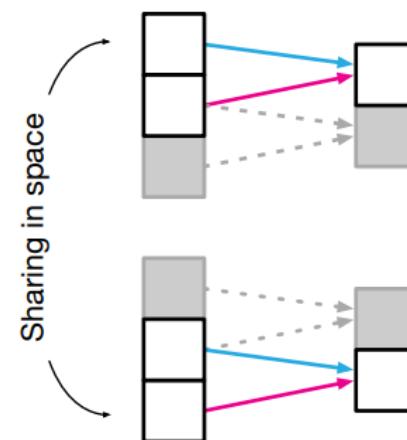
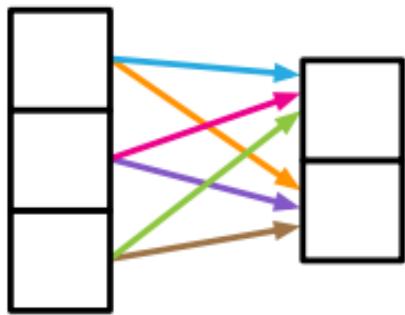
The Transformer has become the **main architecture** in the field of Computer Vision and Natural Language Processing.

# The Development of Transformer



# Inductive Bias and Invariance

Component	Entities	Relations	Rel. inductive bias	Invariance
Fully connected	Units	All-to-all	Weak	-
Convolutional	Grid elements	Local	Locality	Spatial translation
Recurrent	Timesteps	Sequential	Sequentiality	Time translation
Graph network	Nodes	Edges	Arbitrary	Node, edge permutations



# Inductive Bias and Invariance

**The attention architecture of the Transformer, compared to RNNs and CNNs, has a weaker inductive bias and relies on a large amount of training data to surpass the advantages of inductive bias.**

**RNN's temporal locality:** It can only represent the features of adjacent sequence units, while the Attention mechanism of the Transformer can flexibly depict the relevant features of long sequences.

**CNN's spatial locality and translation invariance:** It is suitable for representing local image features, but requires a very deep CNN architecture to expand its receptive field. Meanwhile, the Attention mechanism of ViTal can flexibly represent features of larger visual fields of pixel blocks.

**The Transformer is actually a dynamic MLP that can adjust weights through attention mechanism**

# Inductive Bias and Invariance

RNN



number off

**Limited information**

**Fast**

CNN



different group discussion

**Information Loss**

**Moderate**

Transformer



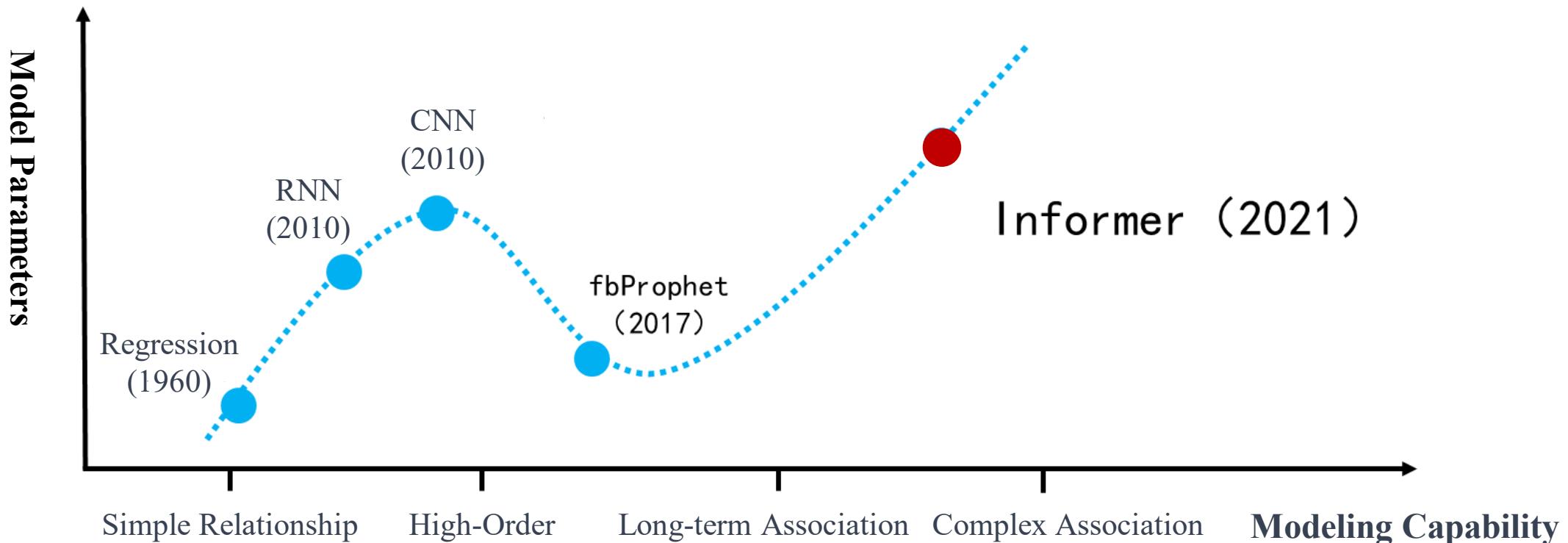
freely exchange

**Rich Information**

**Slow**

# Model Foundations for Sequential Architecture

Informer (AAAI 2021), a specialized neural network model for open sequence data, provides a standard solution for long sequence modeling.

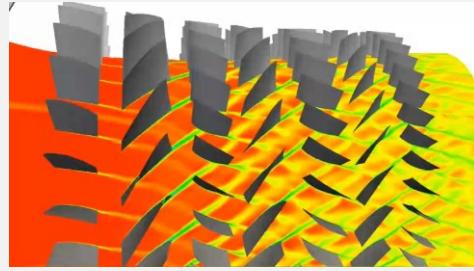


Since March 2021, tracking research on sequence processing solutions includes AutoFormer, EConformer, TSformer, FMMformer, Htransformer-1D, TCCT, TEGRU, and others.

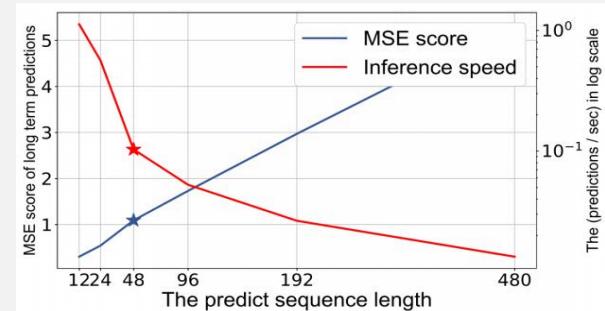
# Question: Representability - Modeling Long Scientific Sequences

## Problems

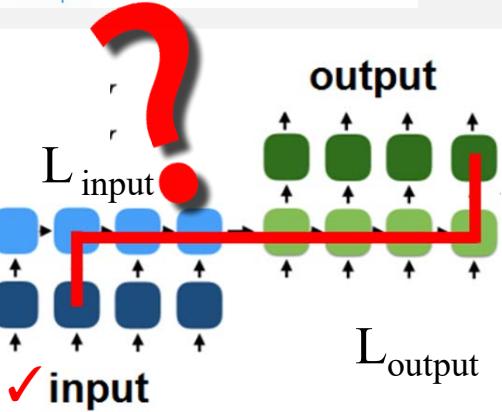
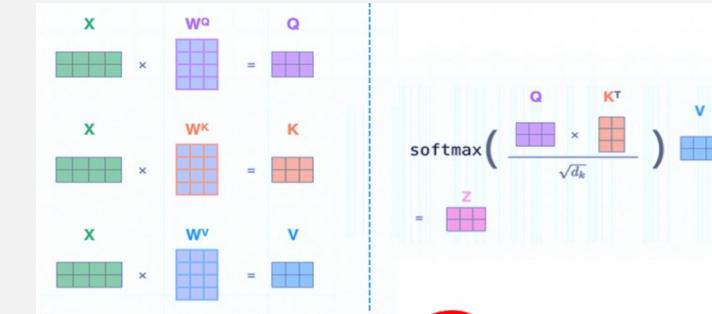
Poor accuracy in long scientific sequence modeling



Thousands of timesteps



Capturing long-term dependencies relies on large-capacity models



The shorter the path, the less traversal in parameter propagation within the model; the stronger the performance, the clearer the long-term dependencies are captured.

	ARIMA	Prophet	LSTM	Conv	Transformer
Long-Seq Input	✓	✓	✗	✓	✗
Long-Seq Output	✗	✗	✗	✗	✗
Complexity per Layer	O(L)	Not clear	O(L*d <sup>2</sup> )	O(k*L*d)	O(L <sup>2</sup> *d)
Longest Path			O(L)	O(log <sub>k</sub> L)	O(1)

# Challenges: Efficient Capture of Long Dependencies

1. **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  $O(L^2)$ .

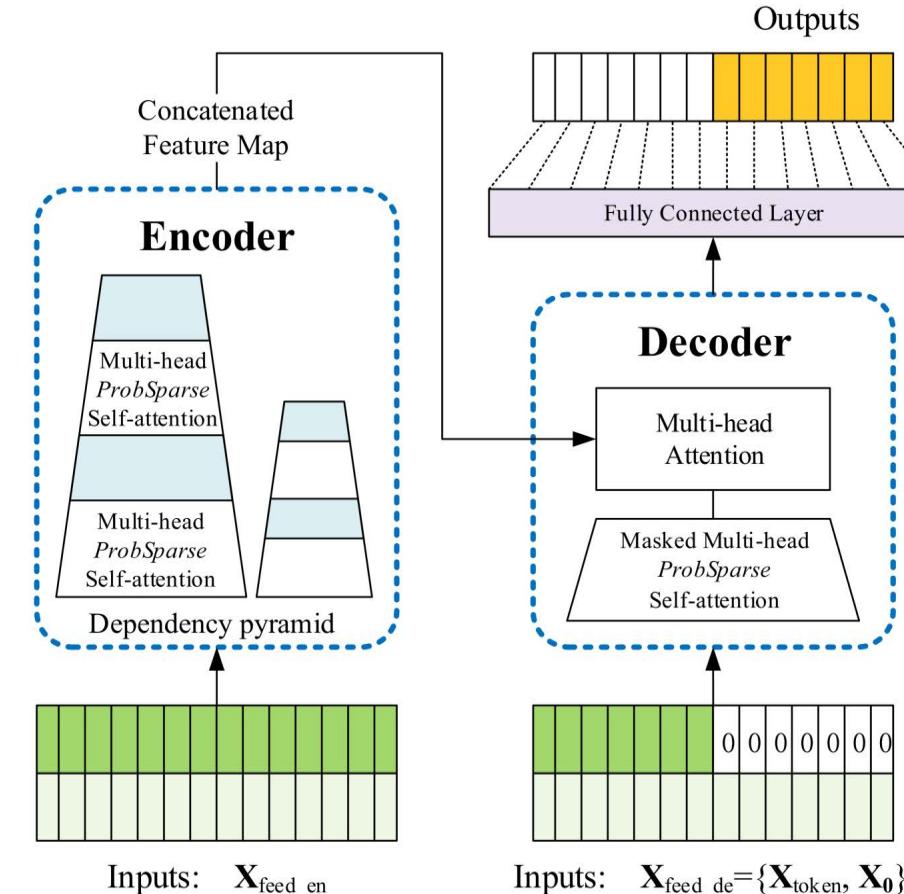
[Complexity/layer]

2. **The memory bottleneck in stacking layers.** The stack of  $J$  encoder/decoder layer makes total memory usage to be  $O(J \cdot L^2)$ , which limits the model scalability on receiving long sequence inputs.

[Long Input]

3. **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]



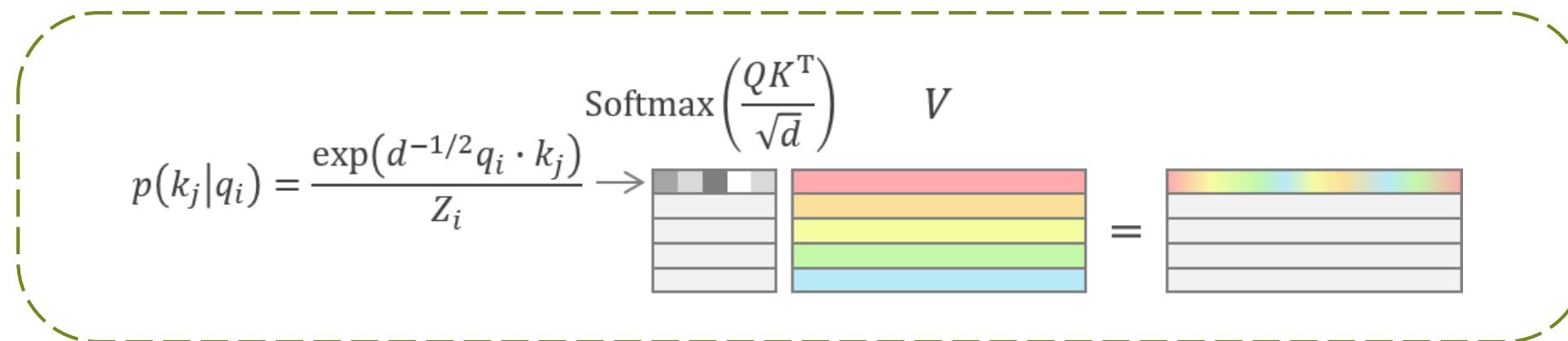
An overall graph of the Informer model.

[1] Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecasting, **AAAI 2021 (Best paper)**

# Challenge 1: Self-attention Mechanism

- Rewrite the original self-attention into the probability formulation

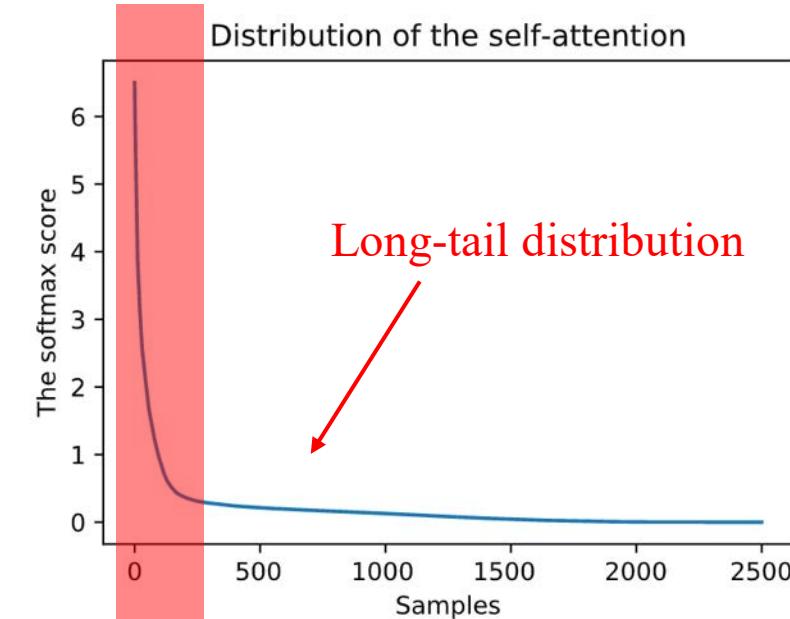
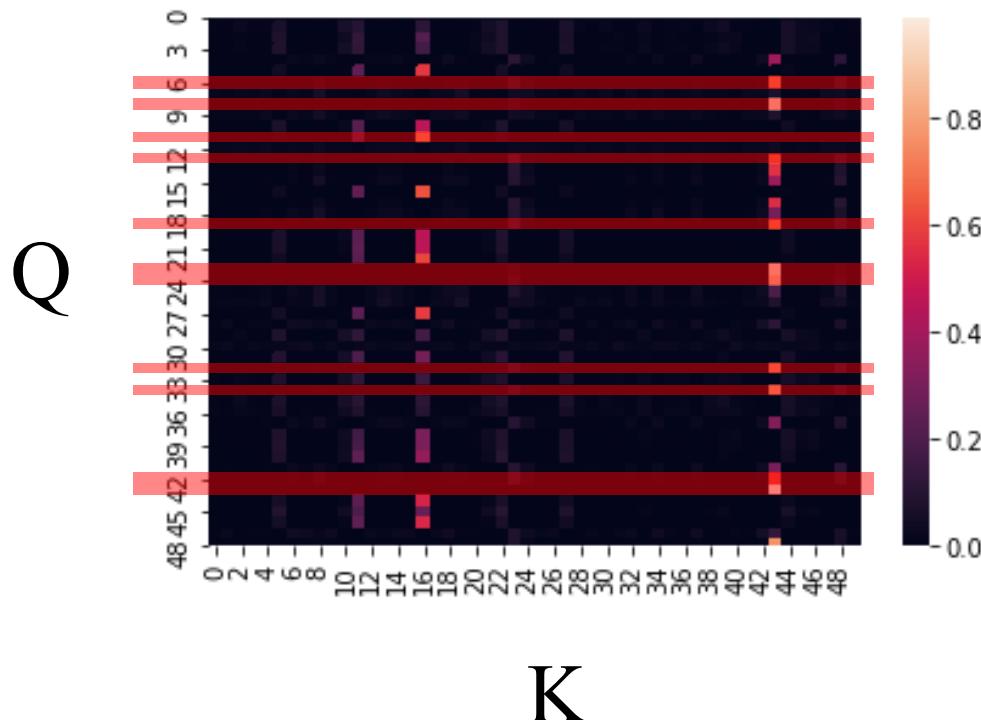
$$\mathcal{A}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V \quad \Rightarrow \quad \mathcal{A}(q_i, K, V) = \sum_j \frac{k(q_i, k_j)}{\sum_l k(q_i, k_l)} v_j = \mathbb{E}_{p(k_j|q_i)}[v_j]$$



It's a new way to investigate the dot-product attention.

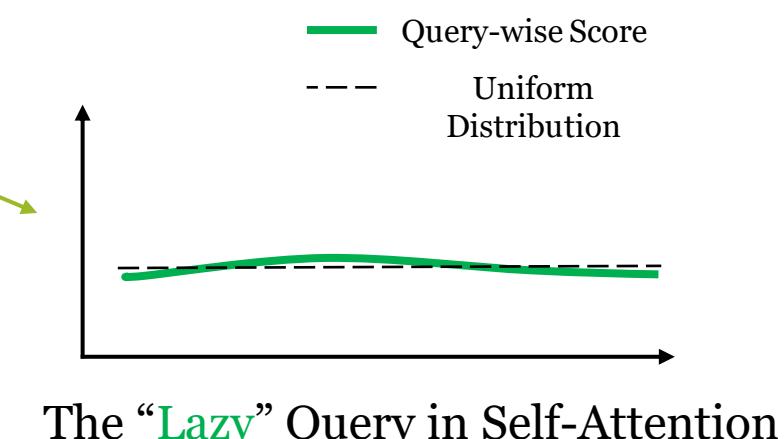
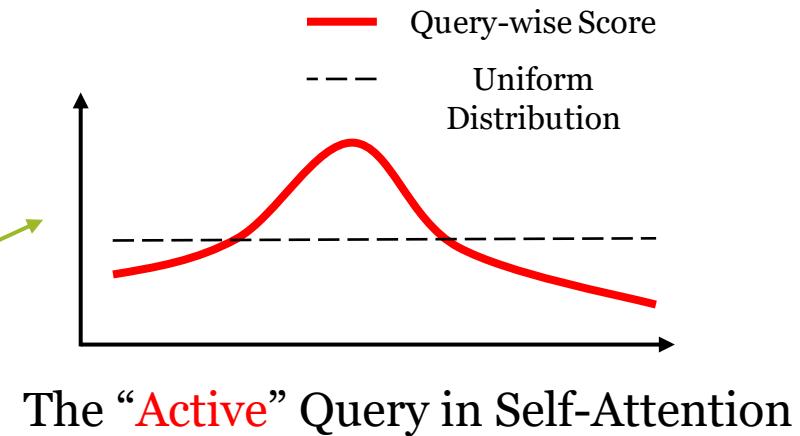
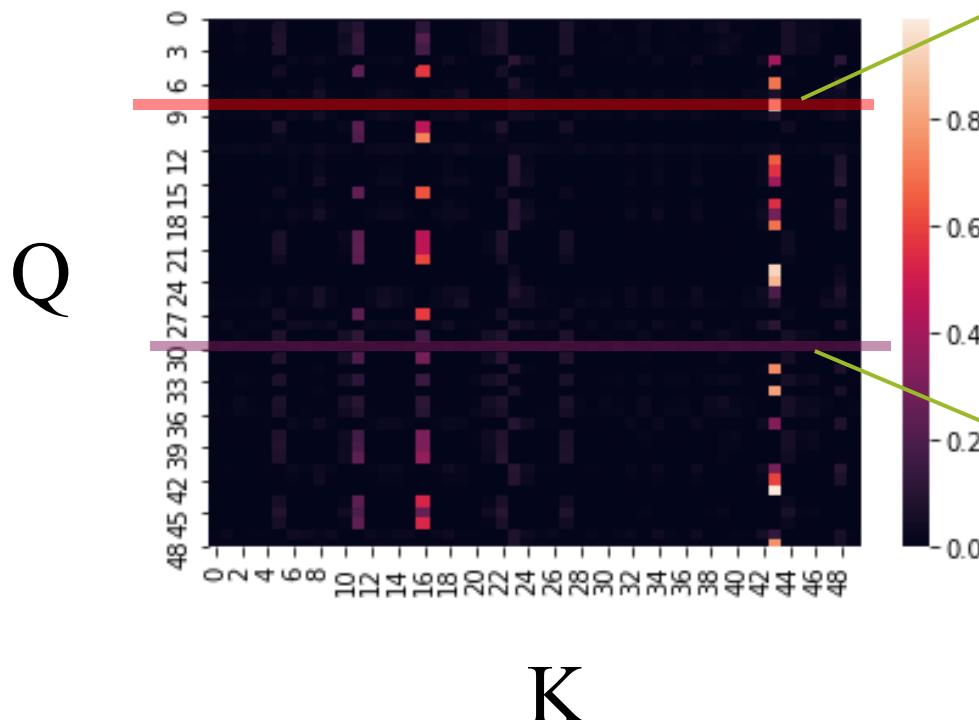
# Challenge 1: Self-attention Mechanism

$$\mathcal{A}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$



# Challenge 1: Self-attention Mechanism

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



# Challenge 1: Self-attention Mechanism

- Rewrite the original self-attention into the probability formulation

$$\mathcal{A}(Q, K, V) = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)V \quad \Rightarrow \quad \mathcal{A}(q_i, K, V) = \sum_j \frac{k(q_i, k_j)}{\sum_l k(q_i, k_l)} v_j = \mathbb{E}_{p(k_j|q_i)}[v_j]$$

- Establish the measurement

Attention probability  $p(k_j|q_i)$       ?      Uniform probability  $q(k_j|q_i) = \frac{1}{L_K}$

$$\begin{aligned} KL(q||p) &= \sum_{j=1}^{L_K} \frac{1}{L_K} \ln \frac{1/L_K}{k(q_i, k_j)/\sum_l k(q_i, k_l)} \\ &= \ln \sum_{l=1}^{L_K} e^{\frac{q_i k_l^\top}{\sqrt{d}}} - \frac{1}{L_K} \sum_{j=1}^{L_K} \frac{q_i k_j^\top}{\sqrt{d}} - \ln L_K \end{aligned}$$

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

$$M(q_i, K) = \ln \sum_{j=1}^{L_K} e^{\frac{q_i k_j^\top}{\sqrt{d}}} - \frac{1}{L_K} \sum_{j=1}^{L_K} \frac{q_i k_j^\top}{\sqrt{d}}$$

# Challenge 1: Self-attention Mechanism

- Rewrite the original self-attention into the probability formulation

$$\mathcal{A}(Q, K, V) = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)V \quad \Rightarrow \quad \mathcal{A}(q_i, K, V) = \sum_j \frac{k(q_i, k_j)}{\sum_l k(q_i, k_l)} v_j = \mathbb{E}_{p(k_j|q_i)}[v_j]$$

- Establish the measurement

$$M(q_i, K) = \ln \sum_{j=1}^{L_K} e^{\frac{q_i k_j^\top}{\sqrt{d}}} - \frac{1}{L_K} \sum_{j=1}^{L_K} \frac{q_i k_j^\top}{\sqrt{d}}$$

- 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*


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

where  $\bar{Q}$  is a sparse matrix of the same size of  $Q$  and it only contains the Top- $u$  queries under the sparsity measurement  $M(q, K)$ .

# Challenge 1: Self-attention Mechanism

- Rewrite the original self-attention into probability formulation
- Establish the measurement
- Define *ProbSparse* self-attention
  
- The measurement have to calculate each dot-product pairs, which still requires  $O(L^2)$  computation.  
And the log-sum-exp fun has the numerical stability issue.

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

Replace

$$\max_j \left\{ \frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d}} \right\}$$

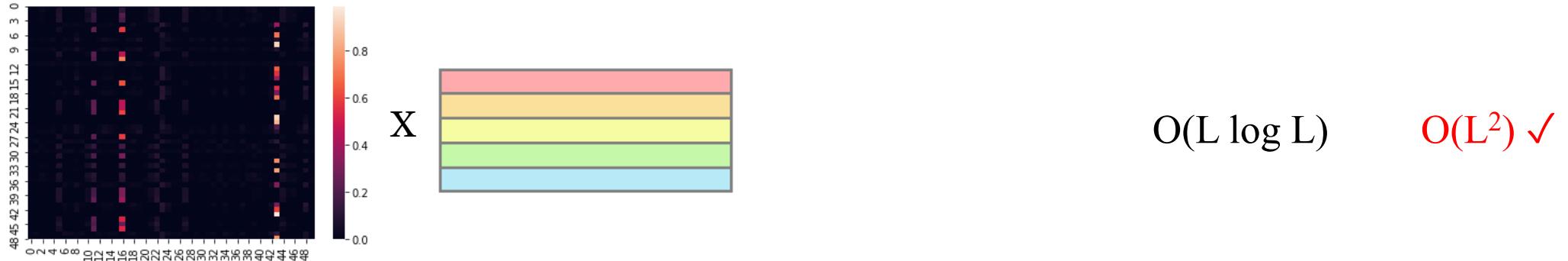
**LEMMA 1.** For each query  $\mathbf{q}_i \in \mathbb{R}^d$  and  $\mathbf{k}_j \in \mathbb{R}^d$  in the keys set  $\mathbf{K}$ , we have the following lower and upper bounds

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

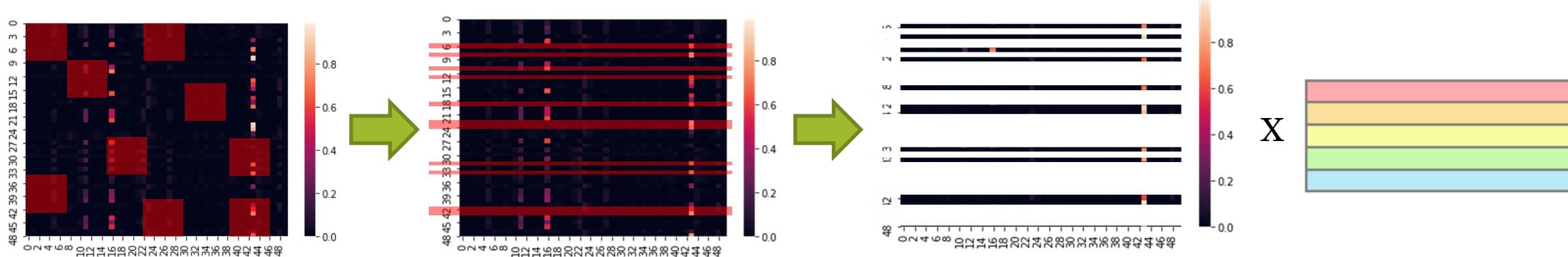
Especially when  $\mathbf{q}_i \in \mathbf{K}$ , it also holds.

# Challenge 1: Self-attention Mechanism

Transformer (2017):



Informer (2021):

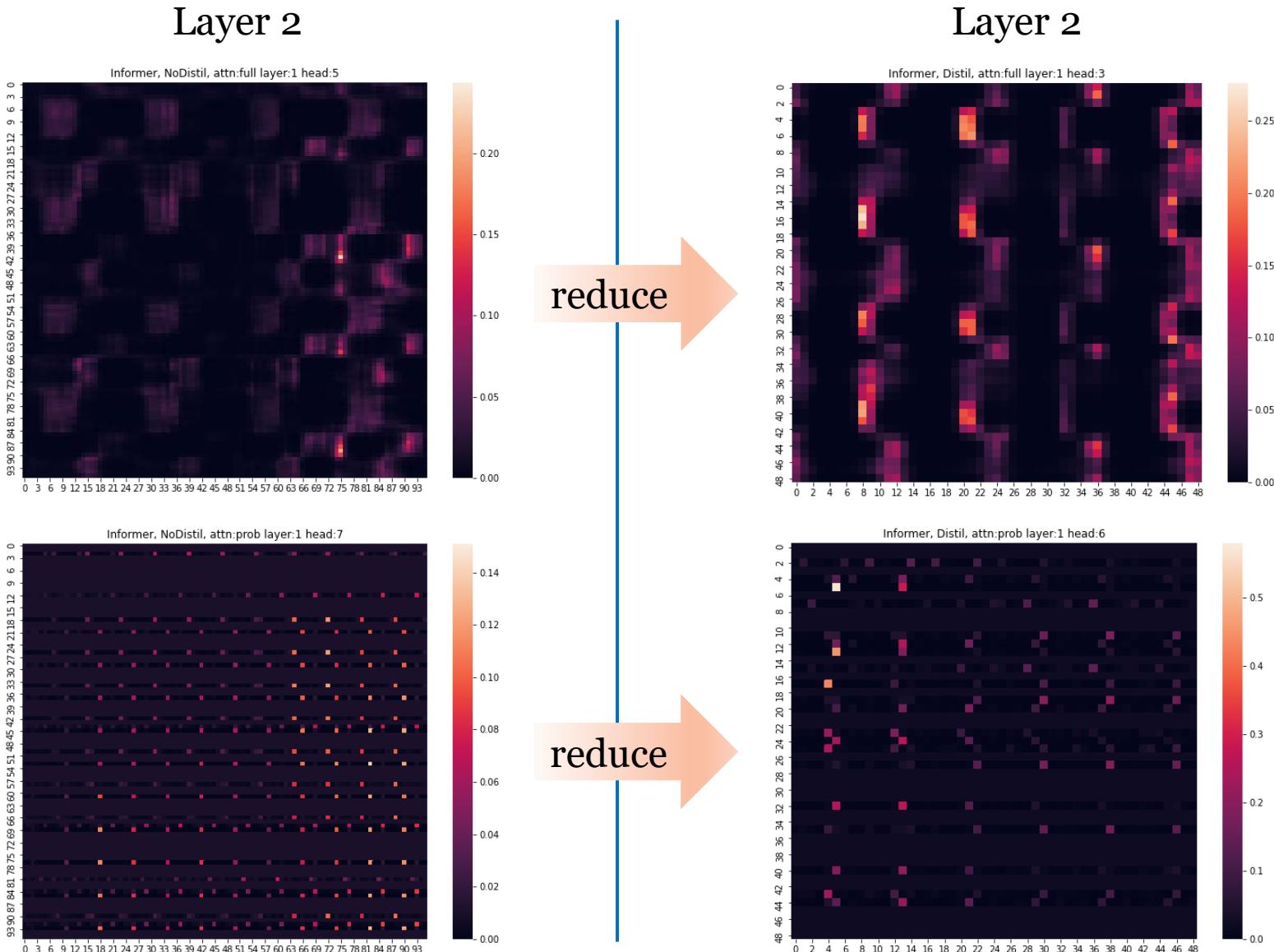


Sample ( $L \log L$ ) dot-product pairs

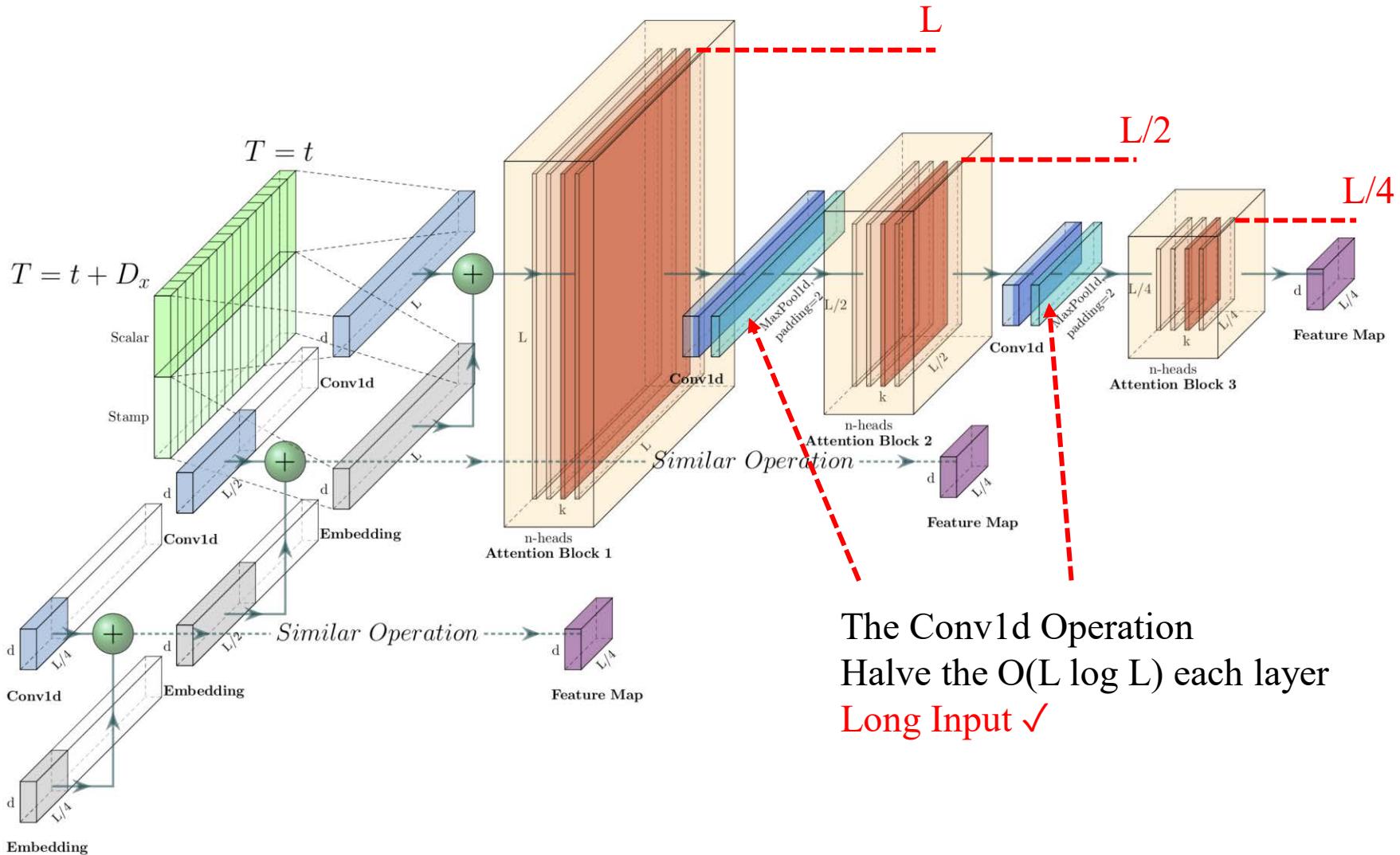
Use  $M$  select  $u$  queries

Calculate Probsparse Attention

# Challenge 2: Self-attention Distilling Operation

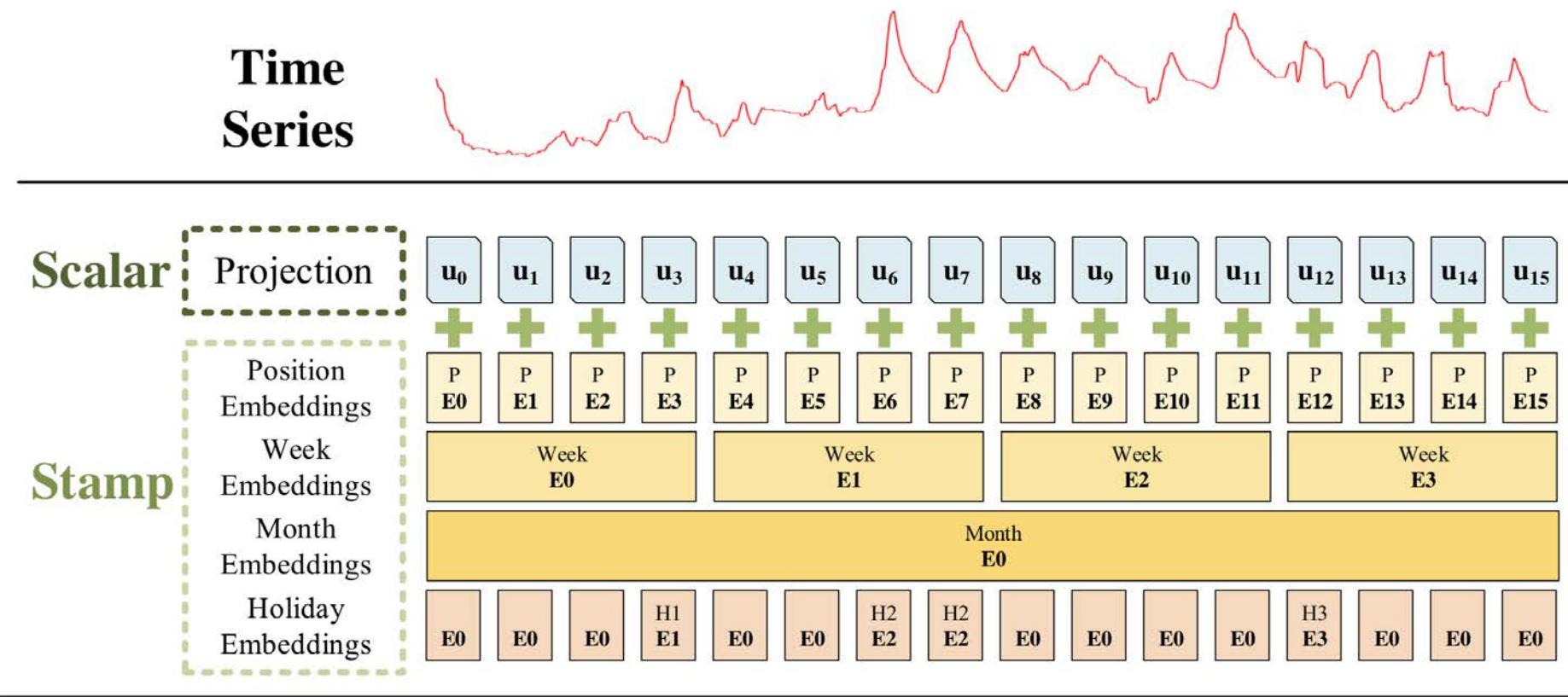


# Challenge 2: Self-attention Distilling Operation



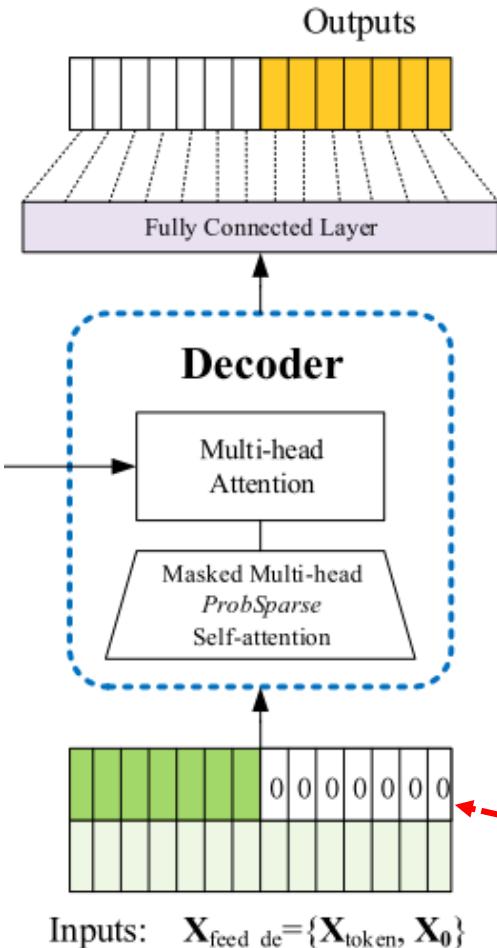
The Conv1d Operation  
Halve the  $O(L \log L)$  each layer  
Long Input ✓

# Challenge 3: Generative-style Decoder



$$\mathcal{X}_{\text{feed}[i]}^t = \alpha \mathbf{u}_i^t + \mathbf{PE}_{(D_x \times (t-1)+i, \cdot)} + \sum_p [\mathbf{SE}_{(D_x \times (t-1)+i)}]_p$$

# 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.

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

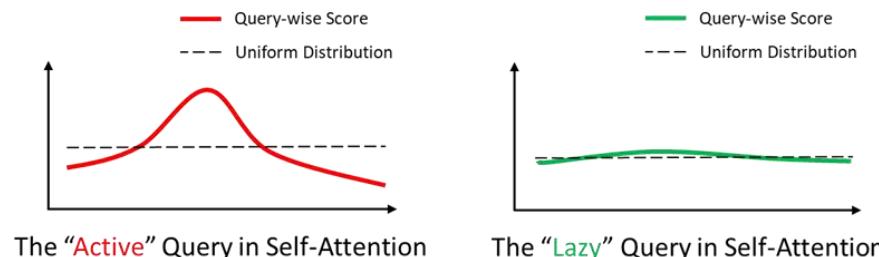
Long Output ✓

# Overview of Informer

## Challenge 1 [Complexity/layer]

**Problem:** The quadratic computation of self-attention.

**Solution:** *ProbSparse Attention*



## Challenge 2 [Long Input]

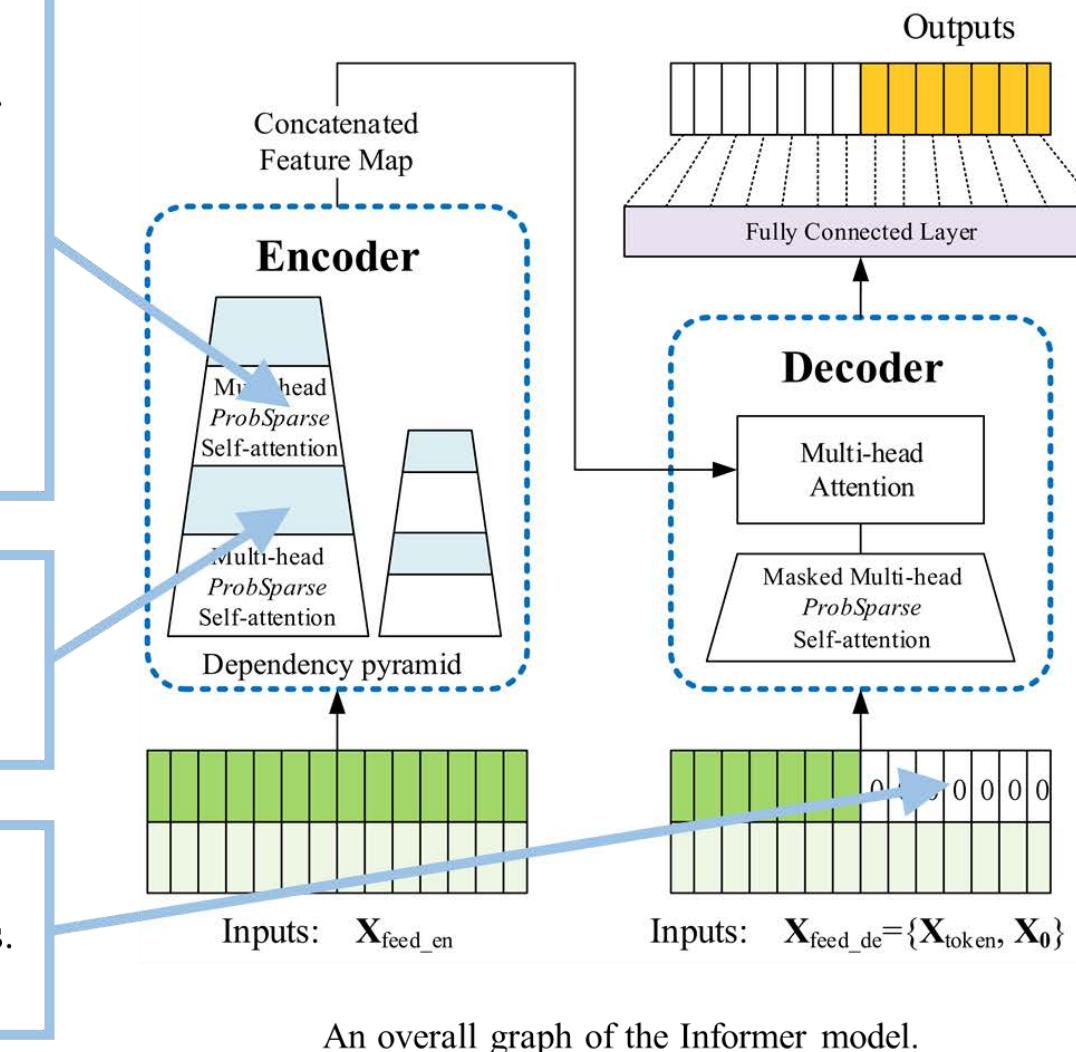
**Problem:** The memory bottleneck in stacking layers.

**Solution:** *Self-attention Distilling*

## Challenge 3 [Long Output]

**Problem:** The speed plunge in predicting long outputs.

**Solution:** *Generative-style Decoder*



# Experiment settings

## Datasets

- ETT (Electricity Transformer Temperature): 2 stations, 2 years, every 15 minutes. Train/Val/Test is 1<sup>st</sup> year / 4 months / 8 months.
- ECL (Electricity Consuming Load): 1 country, 2 years, every 1 hour. Train/Val/Test is 15 months / 3 months / 6 months.
- Weather: 1600 US locations, 4 years, every 1 hour. Train/Val/Test is 28 months / 10 months / 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.

# (1) Univariate Time-series Forecasting

Methods	Informer	Informer <sup>†</sup>	LogTrans	Reformer	LSTMa	DeepAR	ARIMA	Prophet
Metric	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE
ETTh <sub>1</sub>	24	0.098 0.247	<b>0.092 0.246</b>	0.103 0.259	0.222 0.389	0.114 0.272	0.107 0.280	0.108 0.284
	48	<b>0.158 0.319</b>	0.161 0.322	0.167 0.328	0.284 0.445	0.193 0.358	0.162 0.327	0.175 0.424
	168	<b>0.183 0.346</b>	0.187 0.355	0.207 0.375	1.522 1.191	0.236 0.392	0.239 0.422	0.396 0.504
	336	0.222 0.387	<b>0.215 0.369</b>	0.230 0.398	1.860 1.124	0.590 0.698	0.445 0.552	0.468 0.593
	720	0.269 0.435	<b>0.257 0.421</b>	0.273 0.463	2.112 1.436	0.683 0.768	0.658 0.707	0.659 0.766
								2.735 3.253
ETTh <sub>2</sub>	24	<b>0.093 0.240</b>	0.099 0.241	0.102 0.255	0.263 0.437	0.155 0.307	0.098 0.263	3.554 0.445
	48	<b>0.155 0.314</b>	0.159 0.317	0.169 0.348	0.458 0.545	0.190 0.348	0.163 0.341	3.190 0.474
	168	<b>0.232 0.389</b>	0.235 0.390	0.246 0.422	1.029 0.879	0.385 0.514	0.255 0.414	2.800 0.595
	336	0.263 0.417	<b>0.258 0.423</b>	0.267 0.437	1.668 1.228	0.558 0.606	0.604 0.607	2.753 0.738
	720	<b>0.277 0.431</b>	0.285 0.442	0.303 0.493	2.030 1.721	0.640 0.681	0.429 0.580	2.878 1.044
								3.355 4.664
ETTm <sub>1</sub>	24	<b>0.030 0.137</b>	0.034 0.160	0.065 0.202	0.095 0.228	0.121 0.233	0.091 0.243	0.090 0.206
	48	0.069 0.203	<b>0.066 0.194</b>	0.078 0.220	0.249 0.390	0.305 0.411	0.219 0.362	0.179 0.306
	96	0.194 0.372	<b>0.187 0.372</b>	0.199 0.386	0.920 0.767	0.287 0.420	0.364 0.496	0.272 0.399
	288	<b>0.401 0.554</b>	0.409 0.548	0.411 0.572	1.108 1.245	0.524 0.584	0.948 0.795	0.462 0.558
	672	<b>0.512 0.644</b>	0.519 0.665	0.598 0.702	1.793 1.528	1.064 0.873	2.437 1.352	0.639 0.697
								2.747 1.174
Weather	24	<b>0.117 0.251</b>	0.119 0.256	0.136 0.279	0.231 0.401	0.131 0.254	0.128 0.274	0.219 0.355
	48	<b>0.178 0.318</b>	0.185 0.316	0.206 0.356	0.328 0.423	0.190 0.334	0.203 0.353	0.273 0.409
	168	<b>0.266 0.398</b>	0.269 0.404	0.309 0.439	0.654 0.634	0.341 0.448	0.293 0.451	0.503 0.599
	336	<b>0.317 0.416</b>	0.310 0.422	0.359 0.484	1.792 1.093	0.456 0.554	0.585 0.644	0.728 0.730
	720	<b>0.359 0.466</b>	0.361 0.471	0.388 0.499	2.087 1.534	0.866 0.809	0.499 0.596	1.062 0.943
								3.859 1.144
ETL	48	0.239 0.359	<b>0.238 0.368</b>	0.280 0.429	0.971 0.884	0.493 0.539	<b>0.204 0.357</b>	0.879 0.764
	168	0.447 0.503	0.442 0.514	0.454 0.529	1.671 1.587	0.723 0.655	<b>0.315 0.436</b>	1.032 0.833
	336	0.489 0.528	<b>0.501 0.552</b>	0.514 0.563	3.528 2.196	1.212 0.898	<b>0.414 0.519</b>	1.136 0.876
	720	<b>0.540 0.671</b>	0.543 0.578	0.558 0.609	4.891 4.047	1.511 0.966	0.563 0.595	1.251 0.933
	960	<b>0.582 0.608</b>	0.594 0.638	0.624 0.645	7.019 5.105	1.545 1.006	0.657 0.683	1.370 0.982
								6.901 4.264
Count	32	12	0	0	0	6	0	0

Informer (32)

Transformer (12)

- The *Informer* greatly improves the inference performance.
- The *Informer* and its canonical degradation *Informer*<sup>+</sup> show comparable performance and *Informer* beats *Informer*<sup>+</sup> mostly in wining-counts, i.e. 32>12.
- The *Informer* model shows significantly better results than time-series models.
- Our proposed methods achieve better results than Reformer and LogTrans.

## (2) Multivariate Time-series Forecasting

Methods	Informer		Informer <sup>†</sup>		LogTrans		Reformer		LSTMa		LSTnet		
	Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	
ETTh <sub>1</sub>	24	<b>0.577</b>	<b>0.549</b>	0.620	0.57	0.686	0.604	0.991	0.754	0.650	0.624	1.293	0.901
	48	<b>0.685</b>	<b>0.625</b>	0.692	0.67	0.766	0.757	1.313	0.906	0.702	0.675	1.456	0.960
	168	<b>0.931</b>	<b>0.752</b>	0.947	0.79	1.002	0.846	1.824	1.138	1.212	0.867	1.997	1.214
	336	1.128	0.873	<b>1.094</b>	<b>0.81</b>	1.362	0.952	2.117	1.280	1.424	0.994	2.655	1.369
	720	<b>1.215</b>	<b>0.896</b>	1.241	0.91	1.397	1.291	2.415	1.520	1.960	1.322	2.143	1.380
ETTh <sub>2</sub>	24	<b>0.720</b>	<b>0.665</b>	0.753	0.72	0.828	0.750	1.531	1.613	1.143	0.813	2.742	1.457
	48	<b>1.457</b>	<b>1.001</b>	1.461	1.07	1.806	1.034	1.871	1.735	1.671	1.221	3.567	1.687
	168	3.489	<b>1.515</b>	3.485	1.61	4.070	1.681	4.660	1.846	4.117	1.674	<b>3.242</b>	2.513
	336	2.723	1.340	2.626	<b>1.28</b>	3.875	1.763	4.028	1.688	3.434	1.549	<b>2.544</b>	2.591
	720	<b>3.467</b>	<b>1.473</b>	3.548	1.49	3.913	1.552	5.381	2.015	3.963	1.788	4.625	3.709
ETTm <sub>1</sub>	24	0.323	<b>0.369</b>	<b>0.306</b>	0.37	0.419	0.412	0.724	0.607	0.621	0.629	1.968	1.170
	48	0.494	0.503	<b>0.465</b>	<b>0.47</b>	0.507	0.583	1.098	0.777	1.392	0.939	1.999	1.215
	96	<b>0.678</b>	0.614	0.681	<b>0.61</b>	0.768	0.792	1.433	0.945	1.339	0.913	2.762	1.542
	288	<b>1.056</b>	<b>0.786</b>	1.162	0.879	1.462	1.320	1.820	1.094	1.740	1.124	1.257	2.076
	672	<b>1.192</b>	<b>0.926</b>	1.231	1.10	1.669	1.461	2.187	1.232	2.736	1.555	1.917	2.941
Weather	24	<b>0.335</b>	<b>0.381</b>	0.349	0.39	0.435	0.477	0.655	0.583	0.546	0.570	0.615	0.545
	48	0.395	0.459	<b>0.386</b>	<b>0.43</b>	0.426	0.495	0.729	0.666	0.829	0.677	0.660	0.589
	168	<b>0.608</b>	<b>0.567</b>	0.613	0.58	0.727	0.671	1.318	0.855	1.038	0.835	0.748	0.647
	336	<b>0.712</b>	<b>0.620</b>	0.707	0.634	0.754	0.670	1.930	1.167	1.657	1.059	0.782	0.683
	720	<b>0.851</b>	<b>0.731</b>	0.834	0.74	0.885	0.773	2.726	1.575	1.536	1.109	0.851	0.757
ETL	48	0.344	<b>0.393</b>	<b>0.334</b>	0.399	0.355	0.418	1.404	0.999	0.486	0.572	0.369	0.445
	168	0.368	0.424	<b>0.353</b>	<b>0.42</b>	0.368	0.432	1.515	1.069	0.574	0.602	0.394	0.476
	336	0.381	<b>0.431</b>	0.381	0.439	<b>0.373</b>	0.439	1.601	1.104	0.886	0.795	0.419	0.477
	720	0.406	0.443	<b>0.391</b>	<b>0.43</b>	0.409	0.454	2.009	1.170	1.676	1.095	0.556	0.565
	960	<b>0.460</b>	<b>0.548</b>	0.492	0.550	0.477	0.589	2.141	1.387	1.591	1.128	0.605	0.599
Count	33		14		1		0		0		2		

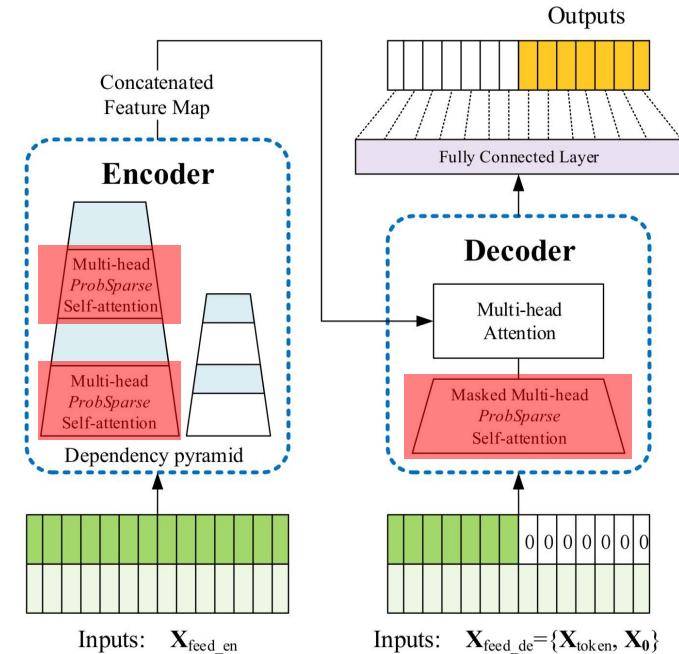
Informer (33)

Transformer (14)

- The *Informer* greatly improves the inference performance.
- The *Informer* and its canonical degradation *Informer*<sup>+</sup> show comparable performance and *Informer* beats *Informer*<sup>+</sup> mostly in winning-counts, i.e. 33>14.
- The previous conclusion in multivariate case holds.

### (3) Ablation study: The ProbSparse Attention

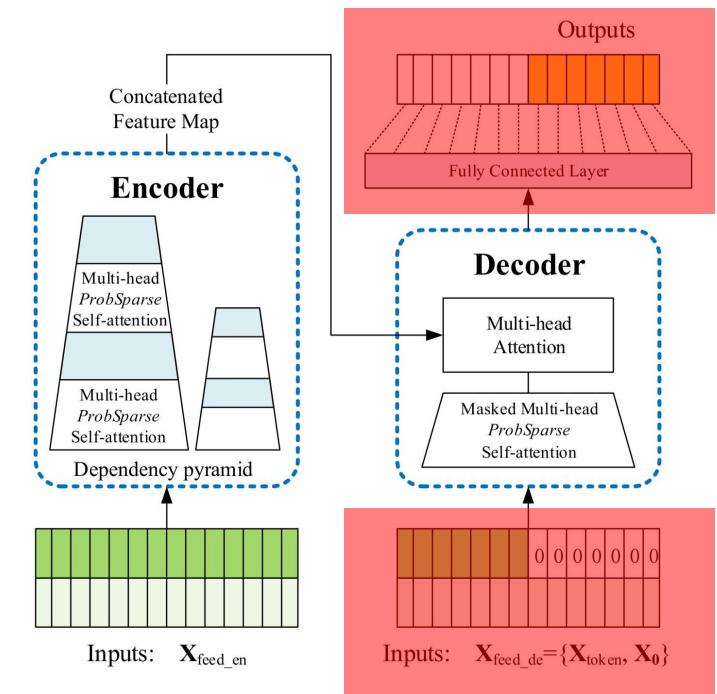
Prediction length		336		720	
Encoder's input		336	720	1440	720
Informer	MSE	0.249	0.225	0.216	0.271
	MAE	0.393	0.384	0.376	0.435
Informer <sup>†</sup>	MSE	0.241	0.214	-	0.259
	MAE	0.383	0.371	-	0.423
LogTrans	MSE	0.263	0.231	-	0.273
	MAE	0.418	0.398	-	0.463
Reformer	MSE	1.875	1.865	1.861	2.243
	MAE	1.144	1.129	1.125	1.536



- The *Informer* can handle longer inputs than vanilla Transformer (*Informer*<sup>†</sup>).
- The *Informer* achieves better results with longer inputs under same prediction windows.

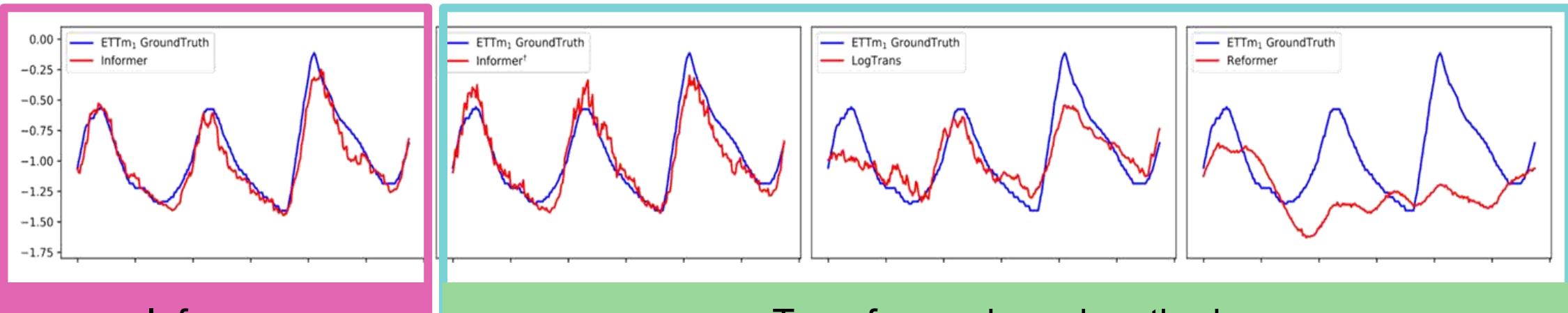
### (3) Ablation study: Generative Style Decoder

Prediction length	336				480					
Prediction offset	+0	+12	+24	+48	+72	+0	+48	+96	+144	+168
Informer <sup>‡</sup>	MSE	0.207	0.209	0.211	0.211	0.216	0.198	0.203	0.203	0.208
	MAE	0.385	0.387	0.391	0.393	0.397	0.390	0.392	0.393	0.401
Informer <sup>§</sup>	MSE	0.201	-	-	-	-	0.392	-	-	-
	MAE	0.393	-	-	-	-	0.484	-	-	-



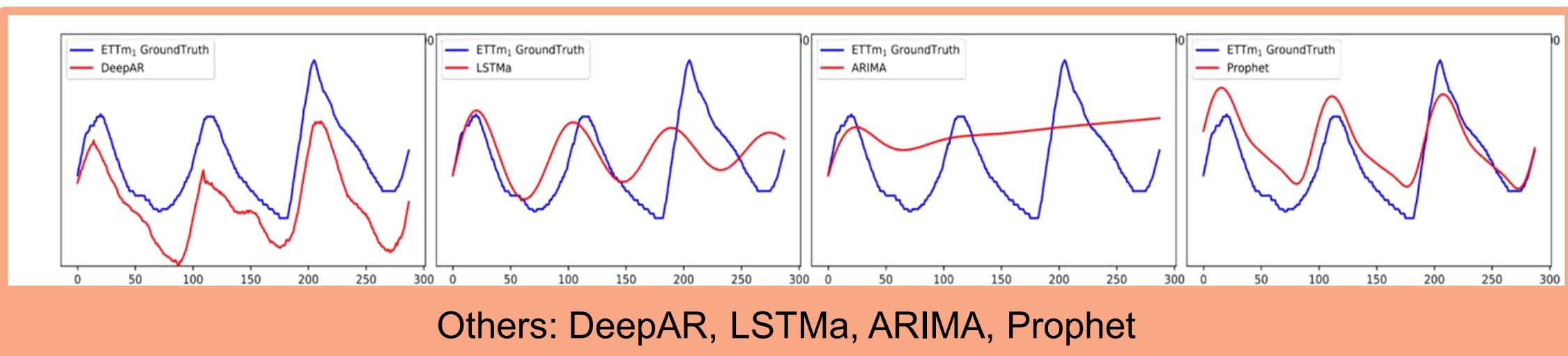
- Note that the *Informer*<sup>++</sup> uses the canonical self-attention and remove the distilling.
- We replace our proposed decoder with dynamic decoding, and it get worse results.
- Moreover, the generative style decoder can accurately predict under the shifting horizon without retraining the entire model.

# (4) Results Visualization

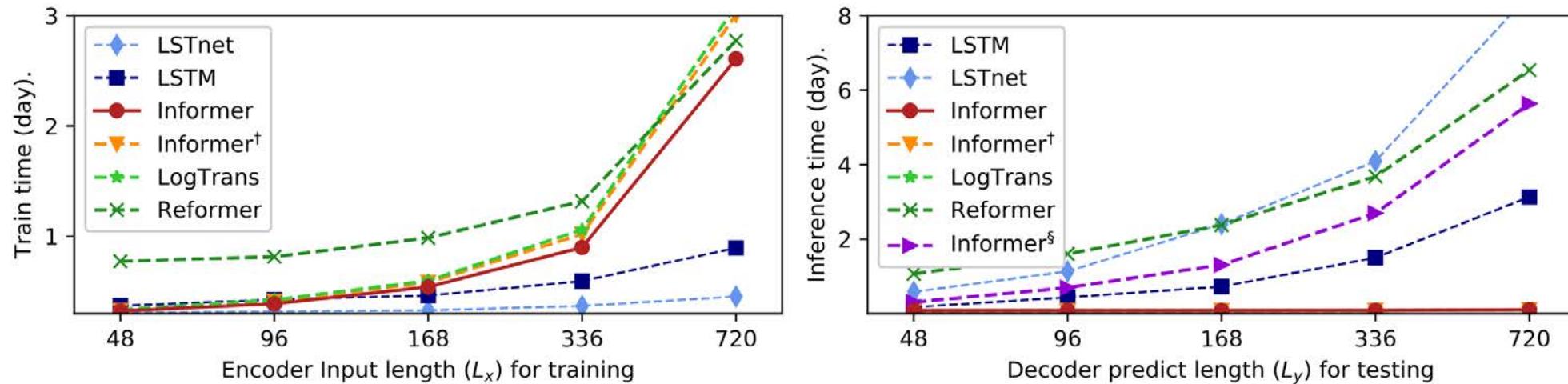


Informer

Transformer-based methods



# (5) Runtime



The total runtime of training/testing phase.

Table 6:  $L$ -related computation statics of each layer

Methods	Training		Testing Steps
	Time Complexity	Memory Usage	
Informer	$\mathcal{O}(L \log L)$	$\mathcal{O}(L \log L)$	1
Transformer	$\mathcal{O}(L^2)$	$\mathcal{O}(L^2)$	$L$
LogTrans	$\mathcal{O}(L \log L)$	$\mathcal{O}(L^2)$	$1^*$
Reformer	$\mathcal{O}(L \log L)$	$\mathcal{O}(L \log L)$	$L$
LSTM	$\mathcal{O}(L)$	$\mathcal{O}(L)$	$L$

<sup>1</sup> The LSTnet is hard to have a closed form.

# The Dataset

Code: <https://github.com/zhouhaoyi/Informer2020>

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.

### 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.

### ETT-small:

We donated two years of data, in which each data point is recorded every minute (marked by  $m$ ), and they were from two regions of a province of China, named ETT-small-m1 and ETT-small-m2, respectively. Each dataset contains 2 year \* 365 days \* 24 hours \* 60 mins = 1,051,200 data point. Besides, we also provide the hourly-level variants for fast development (marked by  $h$ ), i.e. ETT-small-h1 and ETT-small-h2. Each data point consists of 8 features, including the date of the point, the predictive value "oil temperature", and 6 different types of external power load features.

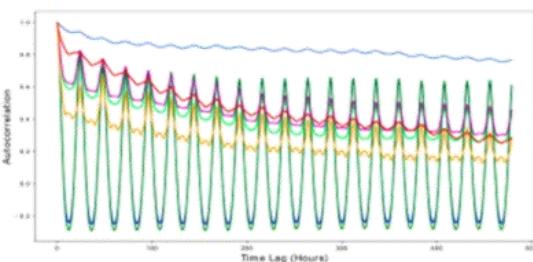
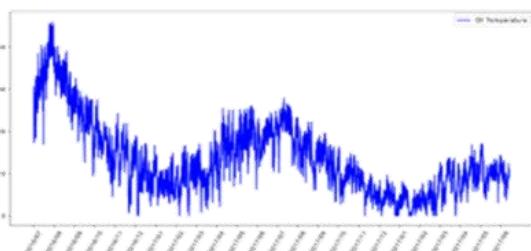


Figure 1.The overall view of "OT" in the ETT-small. Figure 2.The autocorrelation graph of all variables.

The screenshot shows a Jupyter Notebook interface with the title 'Informer.ipynb'. The left sidebar has a '目录' (Table of Contents) section with links to 'Informer Demo', 'Experiments: Train and Test', 'Prediction', 'Visualization', 'Custom Data', and '部分' (Part). The main area shows code execution cells. The first cell clones the repository: [ ] 1 !git clone https://github.com/zhouhaoyi/Informer2020 2 !git clone https://github.com/zhouhaoyi/ETDataset 3 !ls. The second cell imports modules: [ ] 1 import sys 2 if not 'Informer2020' in sys.path: 3 | | sys.path += ['Informer2020']. The third cell installs dependencies: [ ] 1 # !pip install -r ./Informer2020/requirements.txt. The fourth cell runs experiments: [ ] 1 from utils.tools import 2 from exp.exp\_informer import. The right sidebar shows sections for 'Informer Demo' and 'Download code and data'.



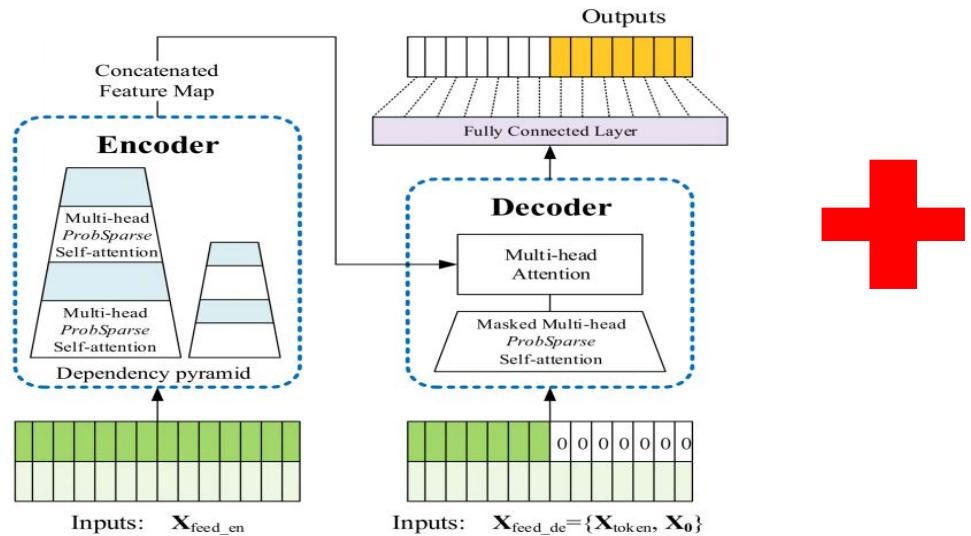
code&colab



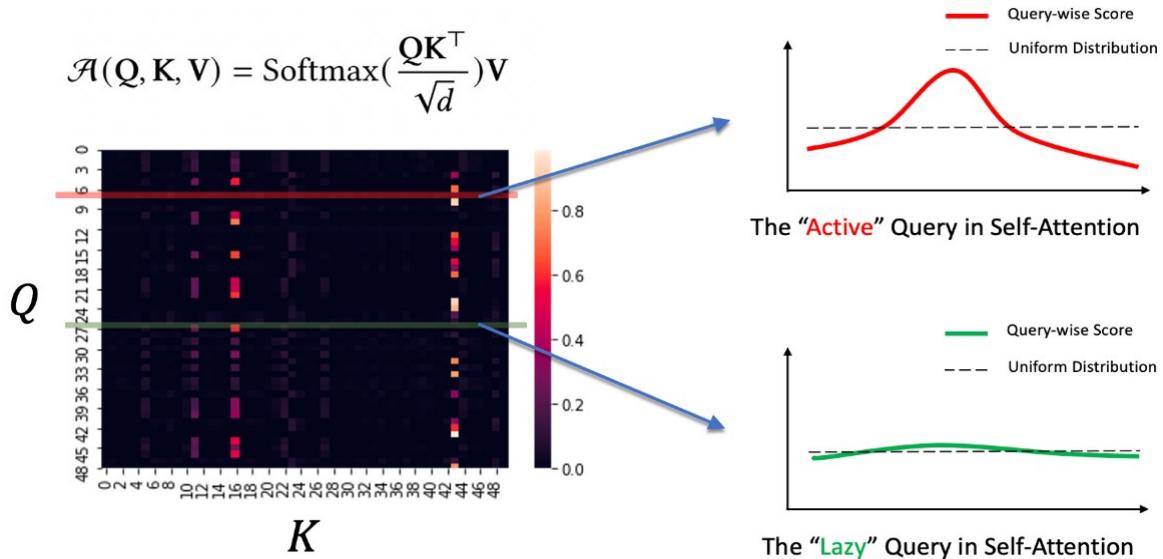
dataset

# Representable: Efficient Capture of Long Dependencies

Novel Architecture: Time Series Prediction Model - Informer

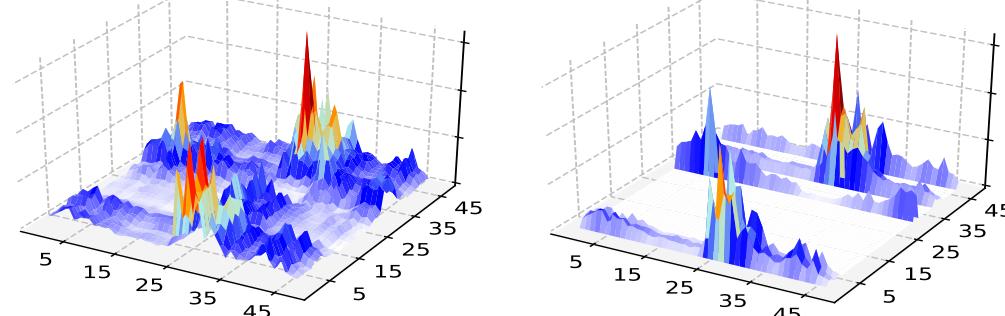


Reducing Complexity: ProbSparse Self-Attention



Probability Distribution-Guided Node Association Sparsity Measurement

Preserving Significant Probability Distribution



Sparse measure:  $\bar{M} = \max_j(\mathbf{q}_i \mathbf{k}_j^T) - \text{mean}_j(\mathbf{q}_i \mathbf{k}_j^T)$

Upper Bound Approximation:  $\bar{M} = \max_j(\mathbf{q}_i \mathbf{k}_j^T) - \text{mean}_j(\mathbf{q}_i \mathbf{k}_j^T)$

Logarithmic linear storage and computational complexity :  $\mathcal{O}(L^2) \rightarrow \mathcal{O}(L \log L)$

[1] Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecasting, **AAAI 2021 (Best paper)**

# Representable: Efficient Capture of Long Dependencies

## Challenge 1 [Diverse Sequential Representations]

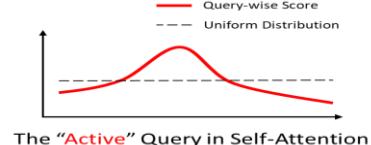
**Problem:** Fusion and Representation Challenges

**Solution:** Unified Input Representation

## Challenge 2 [Complexity per Layer]

**Problem:** Quadratic Complexity of Self Attention

**Solution:** Probabilistic Sparse Attention



## Challenge 3 [Long Sequence Input]

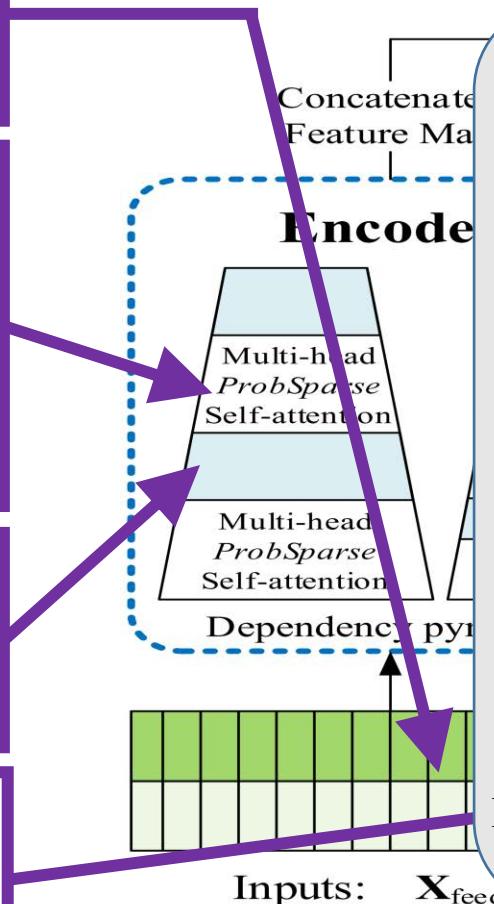
**Problem:** High Overhead for Multi-layer Stacking

**Solution:** Attention Distillation Mechanism

## Challenge 4 [Long Sequence Output]

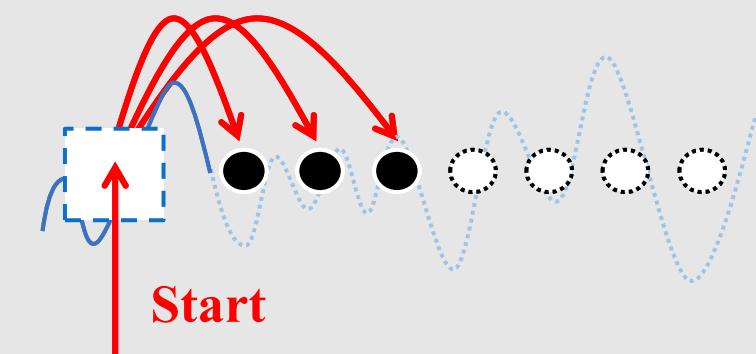
**Problem:** Accumulated Error in SA Decoding

**Solution:** Generative Decoding Mechanism



## Generative Shortest Path Inference

**Improvement:** Predicting Node Scale



$$\text{Coupling Token : } \mathbf{X}_{\text{token}}^t = \begin{array}{|c|} \hline \text{Diagonal Stripes} \\ \hline \end{array} + \tilde{\mathbf{X}}_i^t$$

$$\text{Time Series Representation : } \tilde{\mathbf{X}}_i^t = \begin{array}{|c|} \hline \text{Triangle} \\ \hline \end{array} + \text{PE}_i + \text{SE}_i$$

$$\text{Prediction Step Size : } \text{PE}_i = \begin{cases} \sin(pos/(2L_x)^{2j/d}) \\ \cos(pos/(2L_x)^{2j/d}) \end{cases}$$

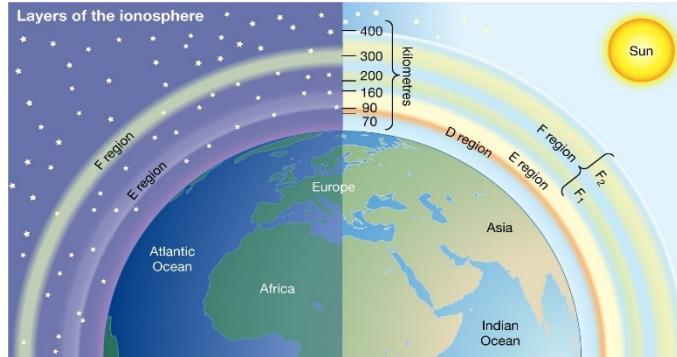
Overall Architecture of Informer

# - Content -

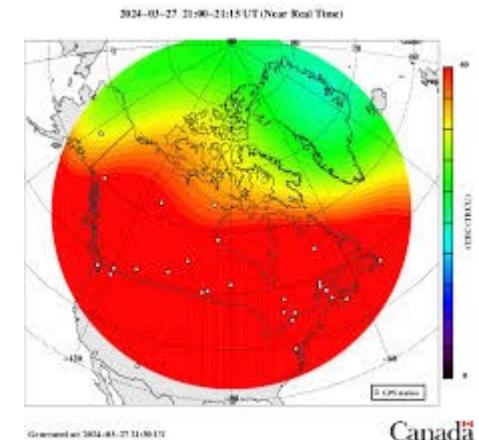
- Sequential Modeling: from Science
- The Transformer Model
- The modern architecture
- **Some applications on scientific sequences**
- Future direction

# Ionformer: Global TEC Prediction

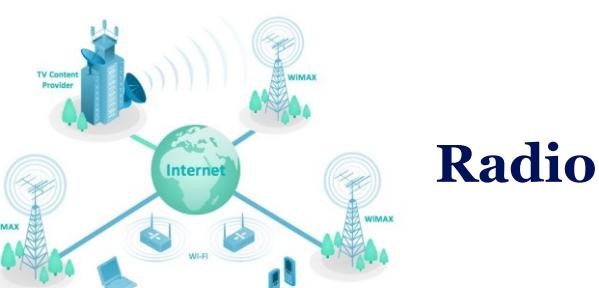
## Total Electron Content



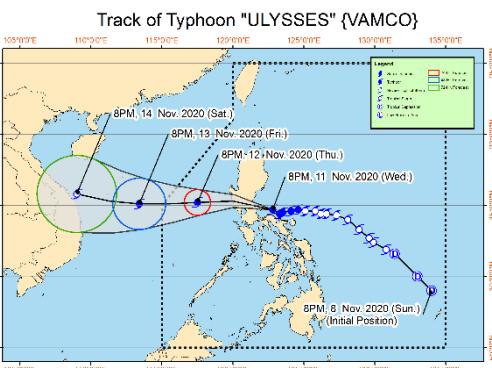
Total Electron Content (TEC) refers to the total **number of free electrons** along a path between two points in Earth's ionosphere, which is an important parameter for understanding the ionospheric effects on radio wave propagation.



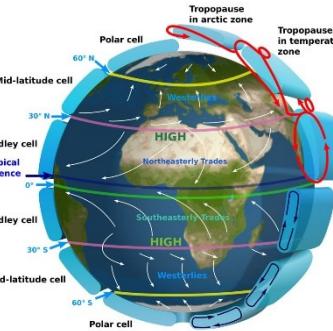
## Application of TEC prediction



## Radio



## Meteorology



## Earth Science

# Challenge in TEC prediction

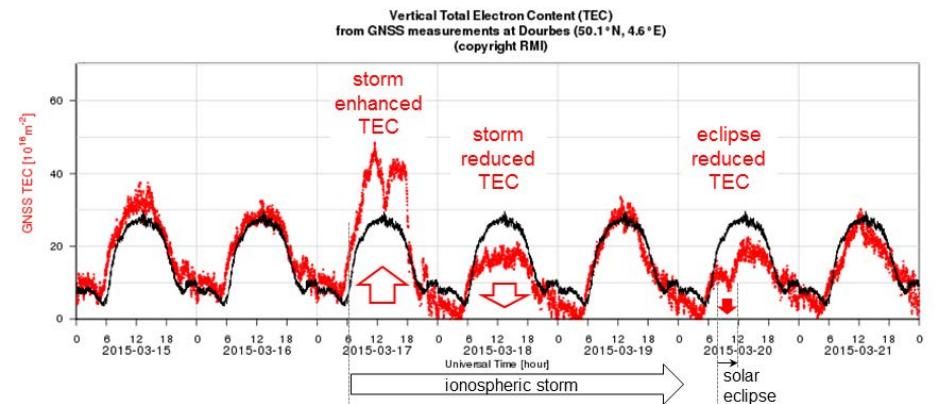
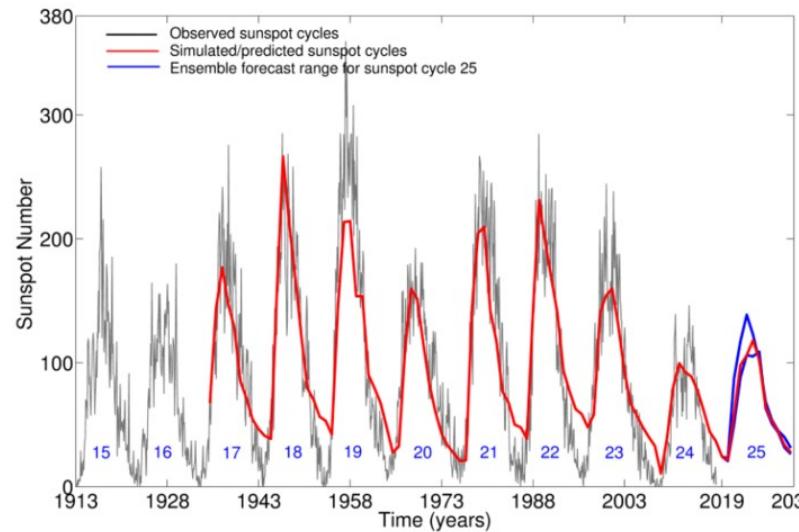
Intense ionospheric activity  
Complex influencing factors

↓ Bottleneck

Difficulty in giving accurate  
prediction under extreme  
climate conditions

↓ Reason

Lack of model architecture  
tailored to the unique  
properties of the ionosphere



# Ionformer: Sparse Attention

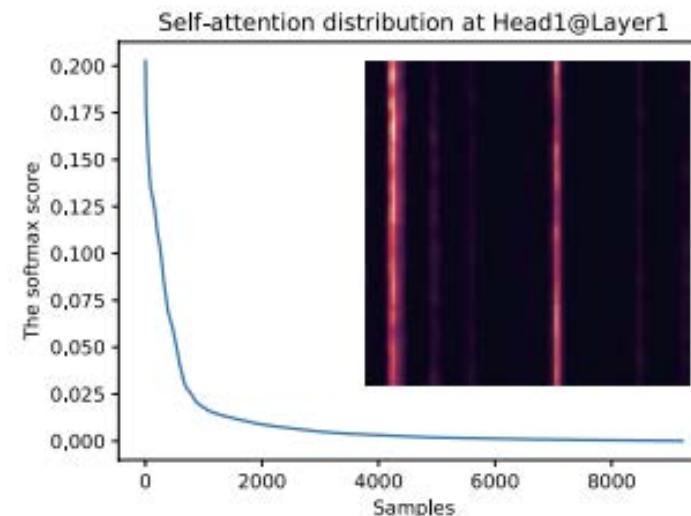
Reduce time complexity, increase robustness

## Challenge

- RNN and LSTM **lose information** in long sequence (TEC) prediction
- Traditional attention involves **quadratic time complexity** and **high memory complexity**

## Solution

- Use **sparse attention** mechanism to reduce time complexity
- Generate **different sparse Q-K pairs** for each head to avoid information loss



$$\overline{M}(q_i, K) = \max_j \left\{ \frac{q_i k_j^T}{\sqrt{d}} \right\} - \frac{1}{L_k} \sum_{j=1}^{L_K} \frac{q_i k_j^T}{\sqrt{d}}$$

**Sparse Measurement**

**Max value**

**Average attention**

# Ionformer: Patchwise Encoding

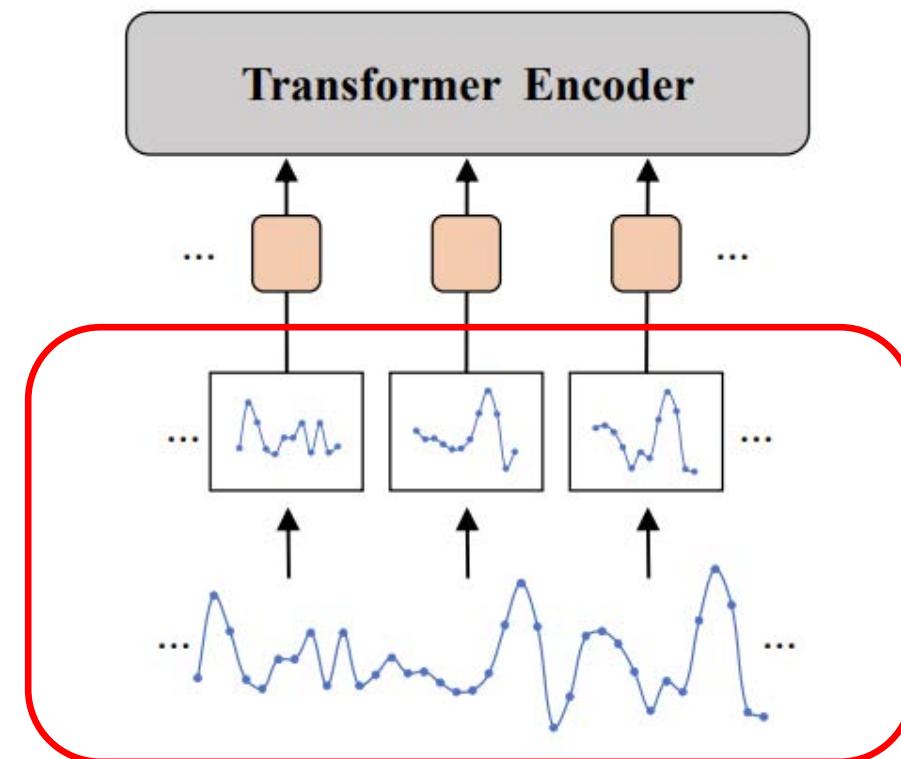
## Temporal local feature aggregation

### Challenge

- Extremely long TEC sequence limited information within single time
- Low-latitude regions are greatly affected by solar and geomagnetic activities

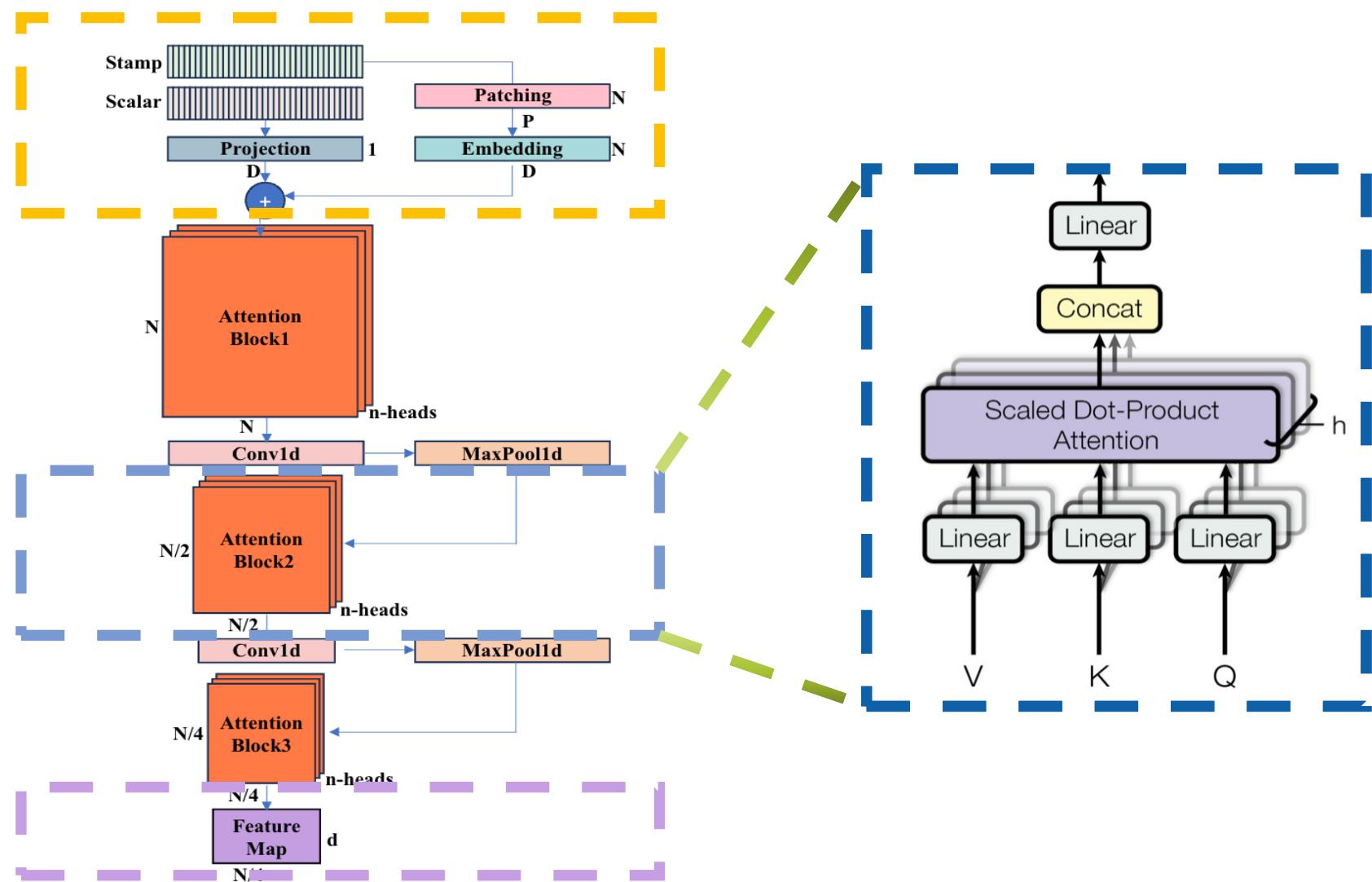
### Solution

- Model information over a period of time rather than a single time step
- Shift from Point-wise to Patch-wise, reducing computational complexity



# Structure of Ionformer

Time-series  
embedding

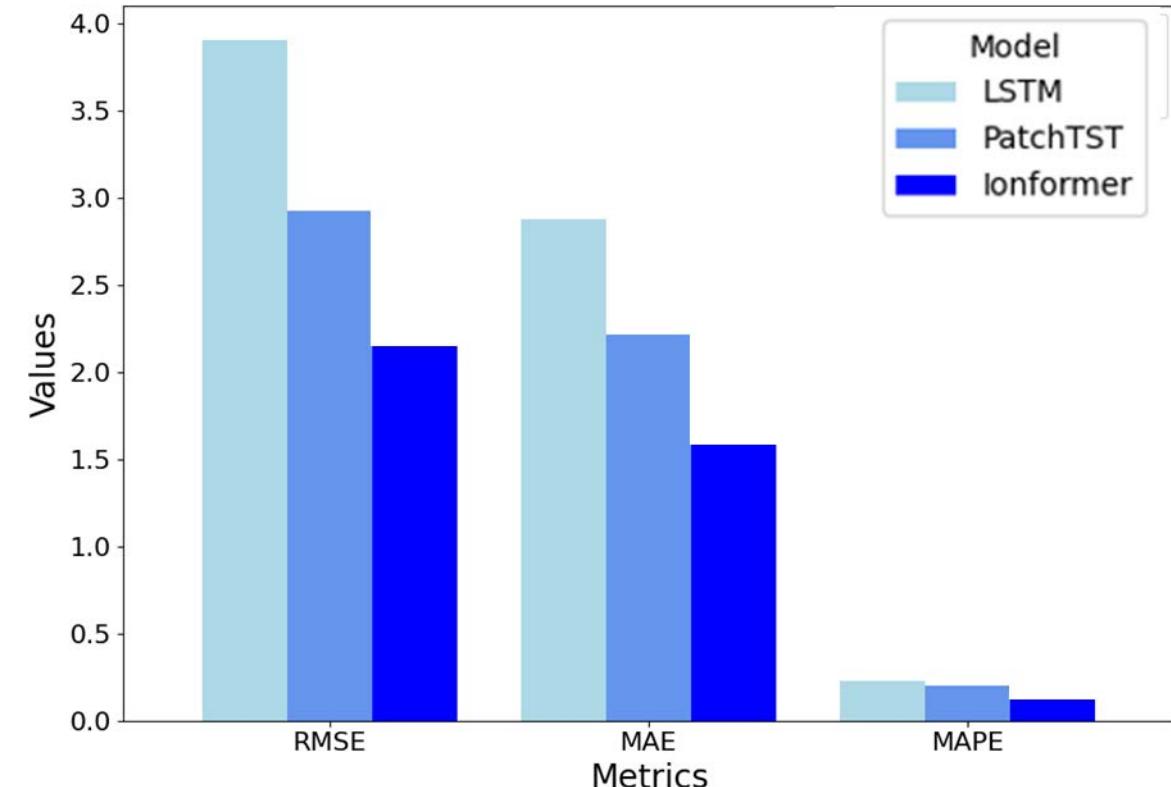


# Experimental Results

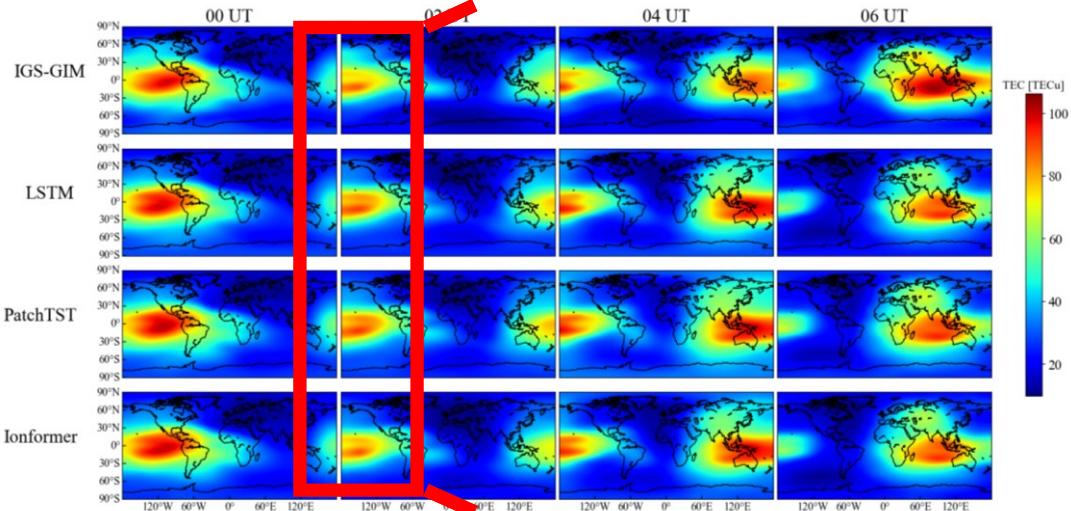
Year	2014			2017		
Metric	RMSE	MAE	MAPE	RMSE	MAE	MAPE
LSTM	4.725	3.245	0.153	1.933	1.376	0.136
Informer	4.582	3.096	0.114	1.787	1.269	0.123
PatchTST	4.322	2.910	0.135	1.739	1.221	0.119
Ionformer	4.209	2.794	0.122	1.638	1.117	0.107

**Ionformer achieves SOTA results on the TEC dataset**

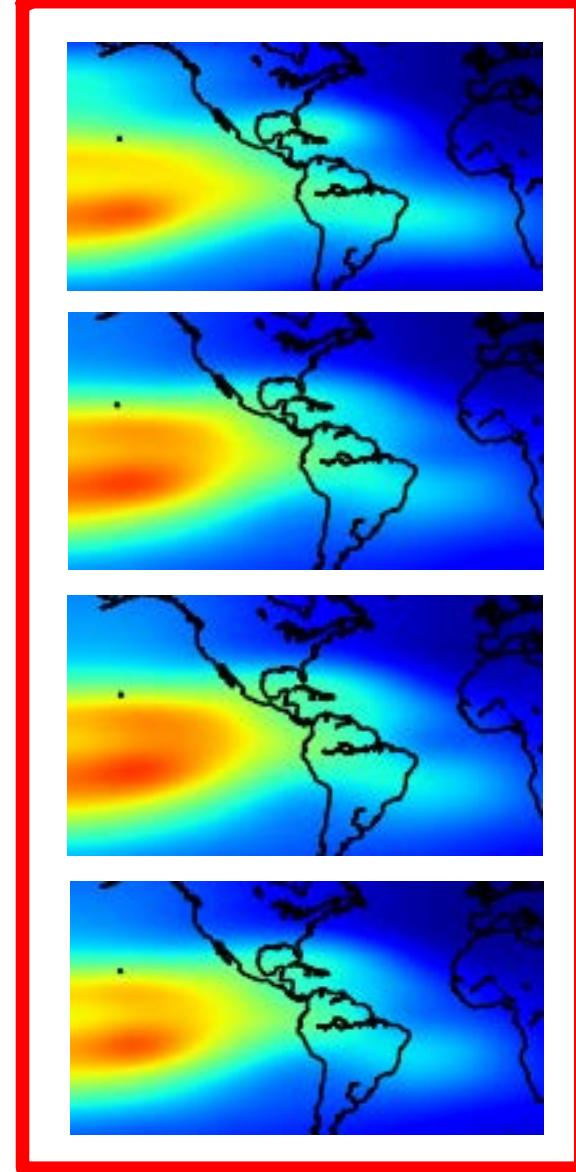
**20% of improvement compared to baselines**



# Results Visualization



**Ionformer can better capture the local pattern of TEC**



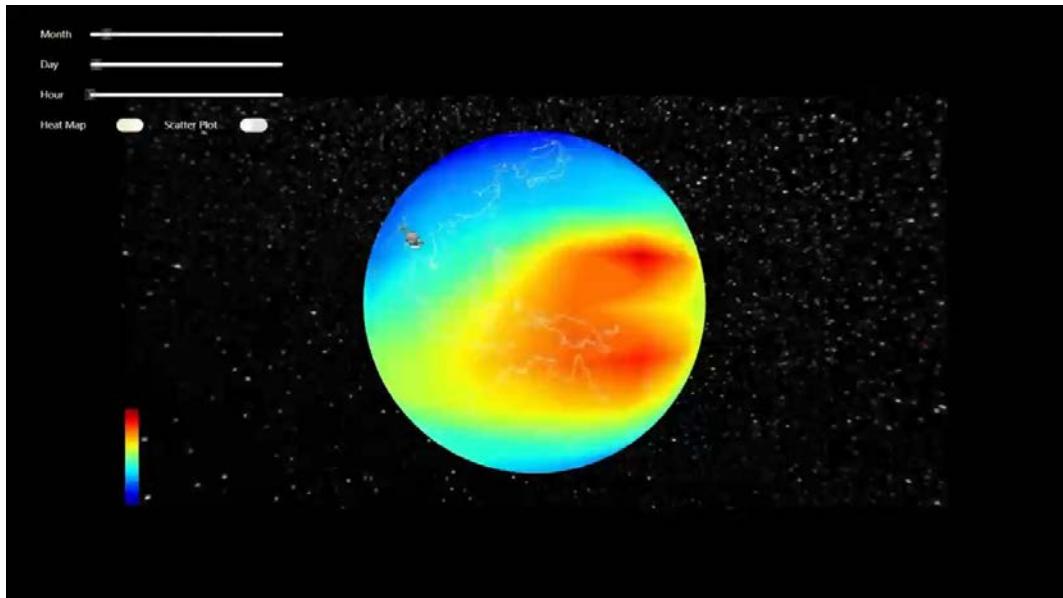
**Ground Truth**

**LSTM**

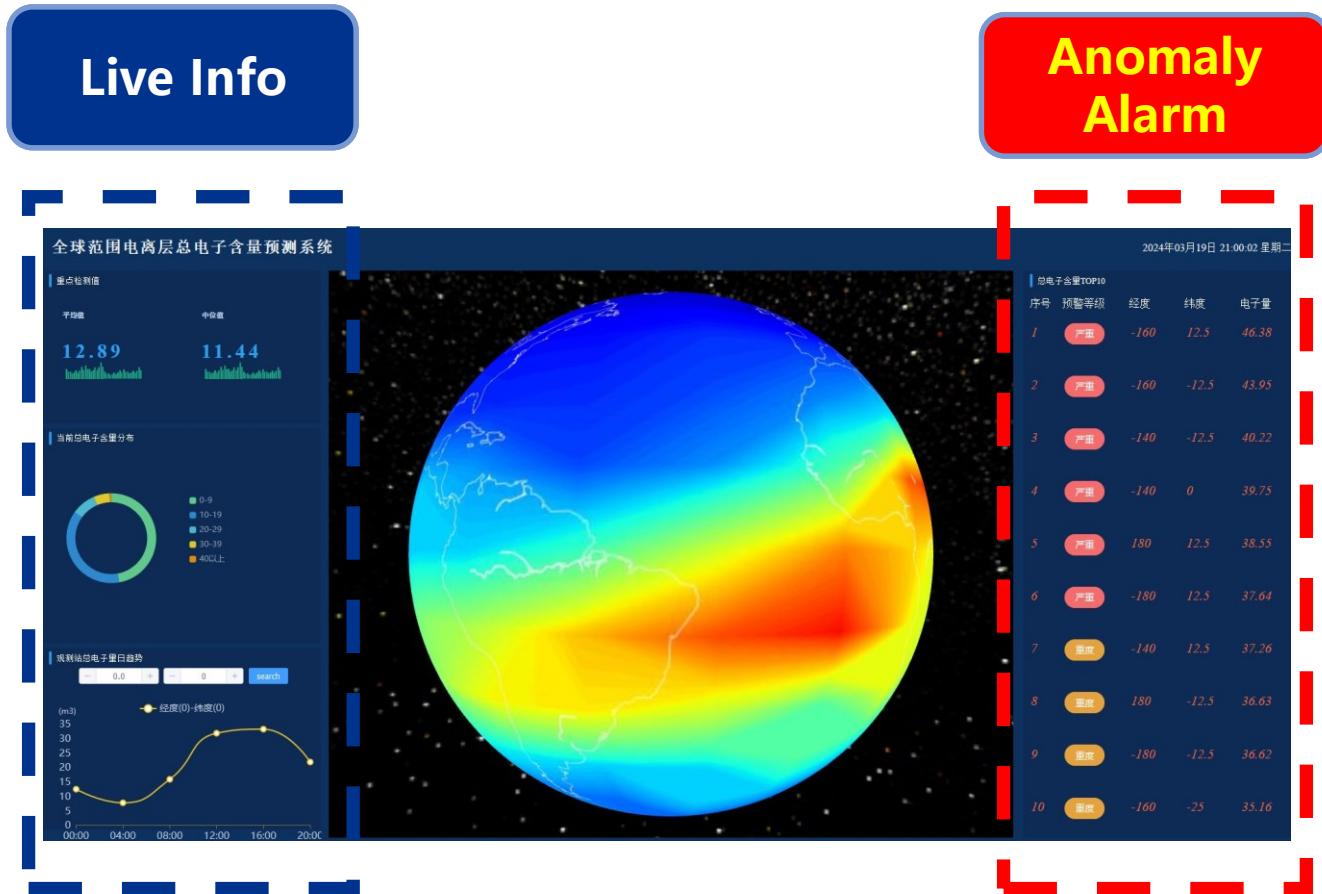
**PatchTST**

**Ionformer**

# TEC Forecasting System



Prediction Results



System Interface

A functional TEC prediction system based on Ionformer

学语文

学代码

学物理

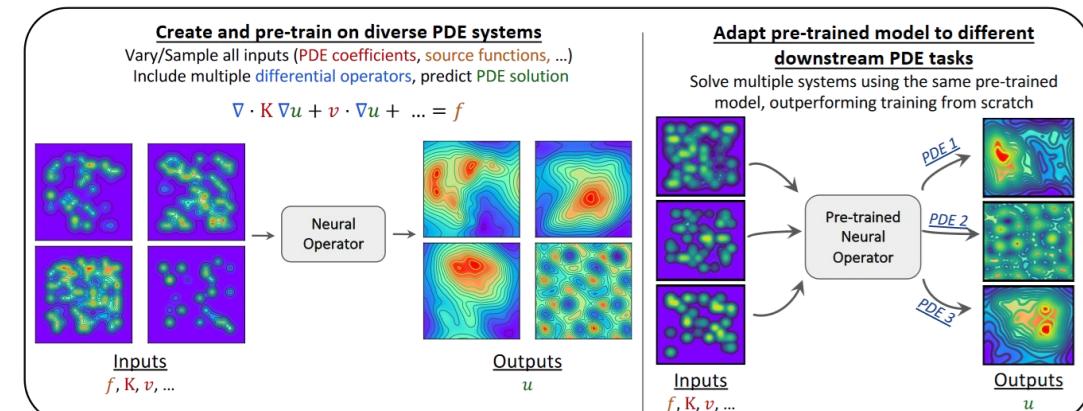
- “完形填空”式训练
- 缺乏数理逻辑能力
- 涌现能力：无

- “指令教育”式训练
- 部分任务获得逻辑推理能力
- 涌现能力：世界观、思维链、专业软件涌现能力

- “逻辑归纳”式训练
- 数据规则和物理模型的能力
- 涌现能力：科学知识生成能力、物理规律修正能力

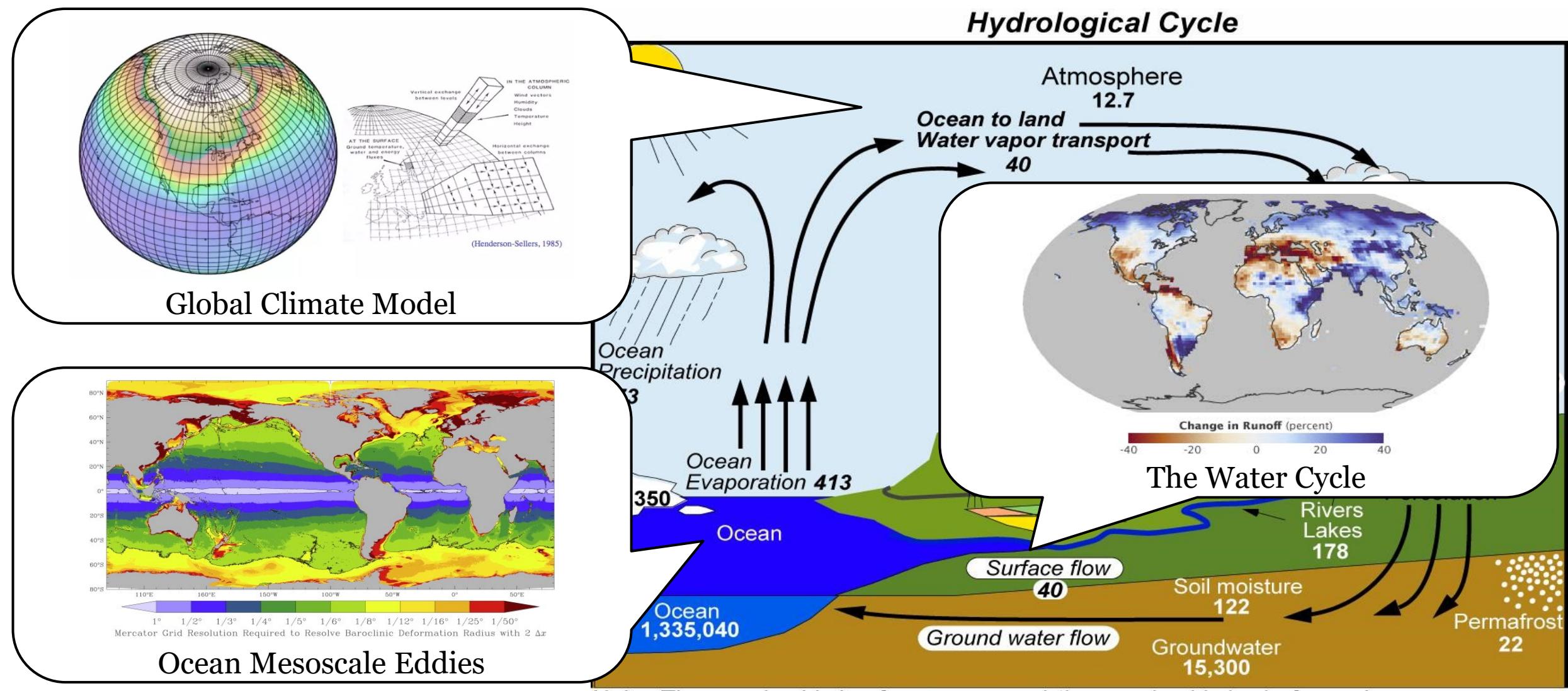
## 美国大规模启动科学智算大模型系列研究

- 【新组织】2023年1月2日，由美国能源部牵头，布朗大学、哈佛大学等成立专门研究组织SciML
- 【新理论】2023年5月27日，CMU率先提出可用于高通量、高维度物理数据的新型神经算子Hyena
- 【新模型】2023年10月9日，Meta公司支持，图灵奖得主Yann Lecun指导发布多物理场大模型MPP
- 【新能力】2023年11月10日，劳伦斯伯克利国家实验室证明神经算子可扩展到2B参数，**新方程发现、物理模型修正**等能力

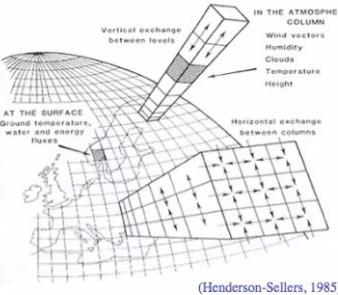
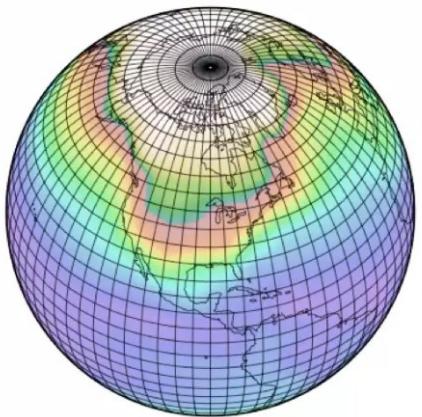


美国劳伦斯伯克利国家实验室新闻专栏报道

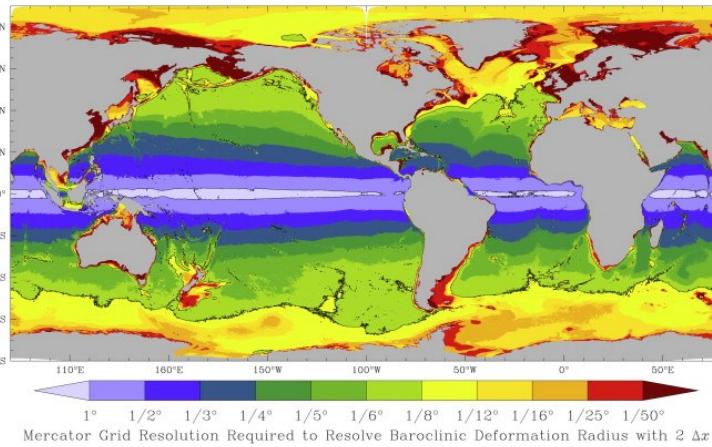
# Building Foundation Model for Scientific Computing



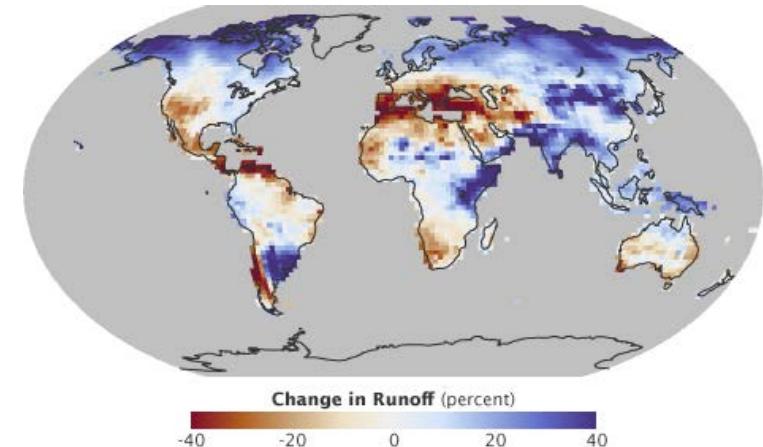
# Building Foundation Model for Scientific Computing



Global Climate Model



Ocean Mesoscale Eddies



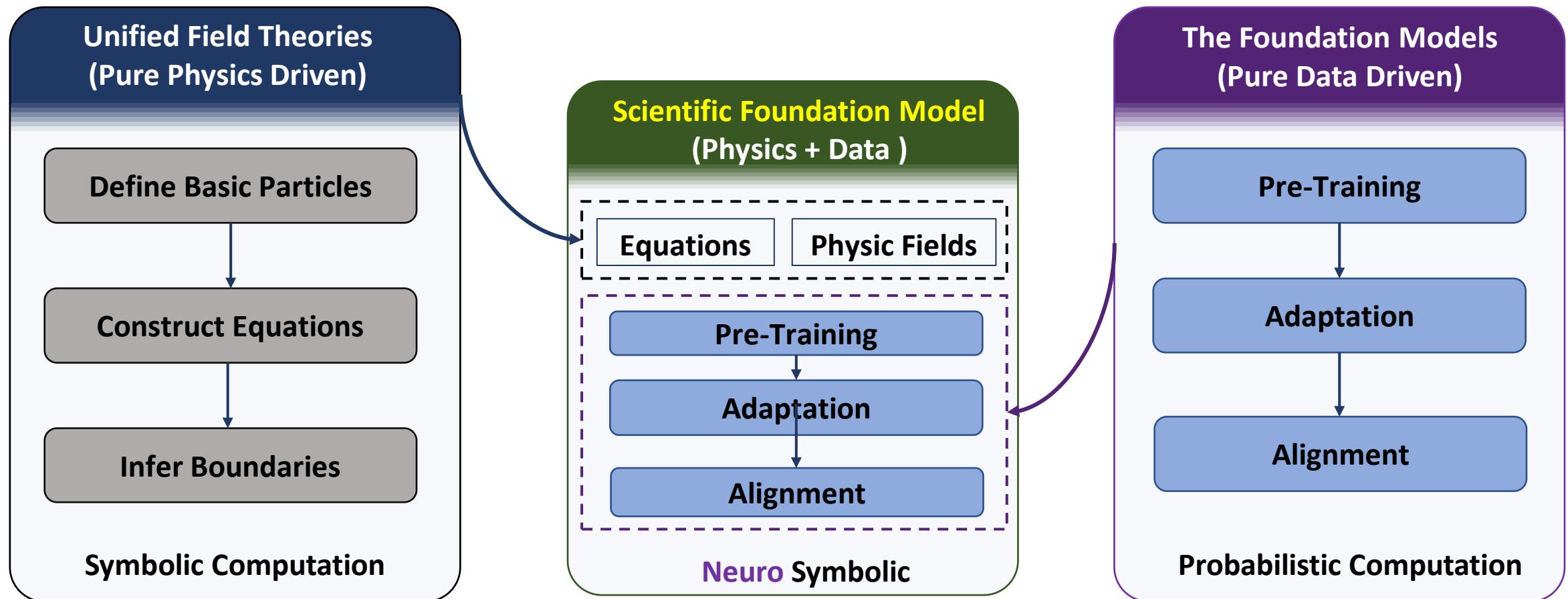
The Water Cycle



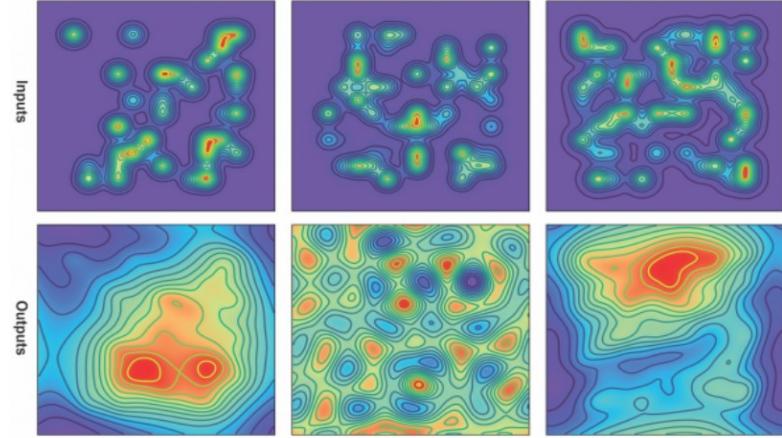
General Foundation Model?

# Building Foundation Model for Scientific Computing

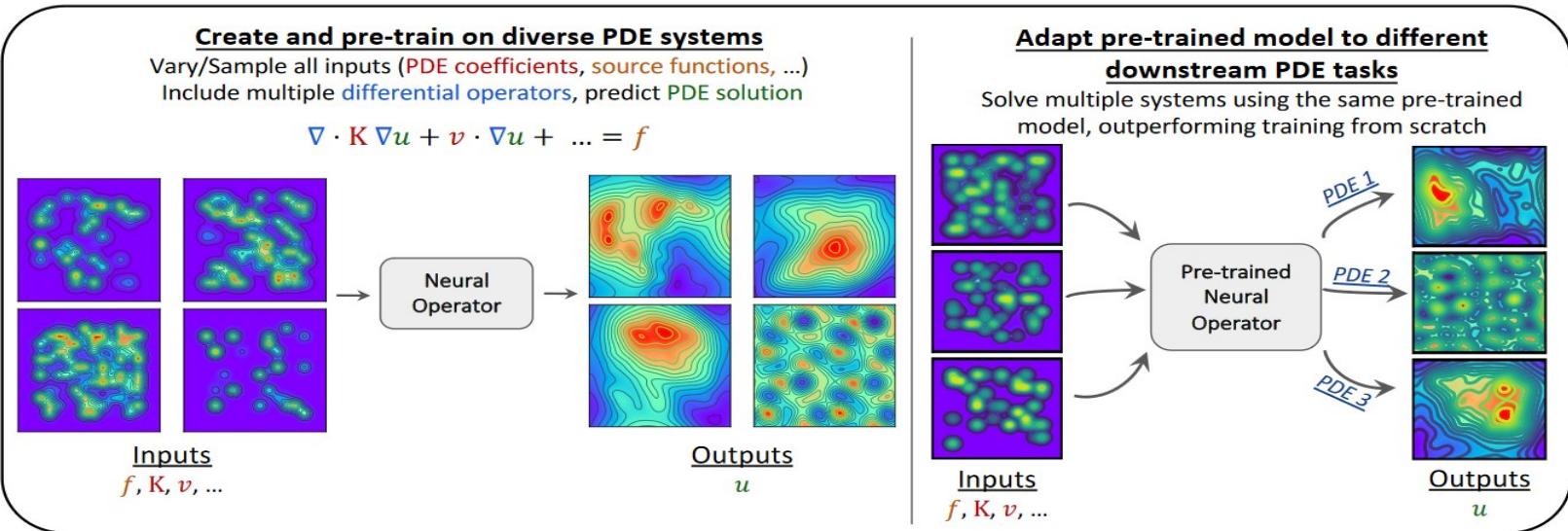
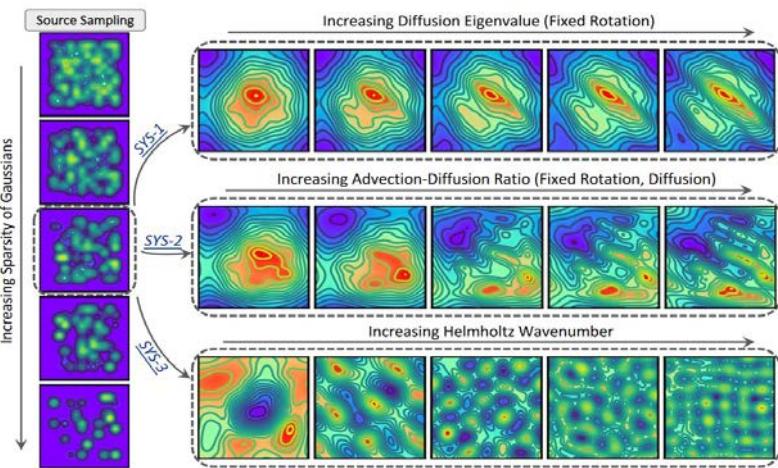
- Scientific Computing: The Grant Unified Field Theories
- Flexible Machine Learning: The Foundation Models (GPT-4, e.t.c)
- “Physics + data driven”: build flexible machine learning for Scientific Computing



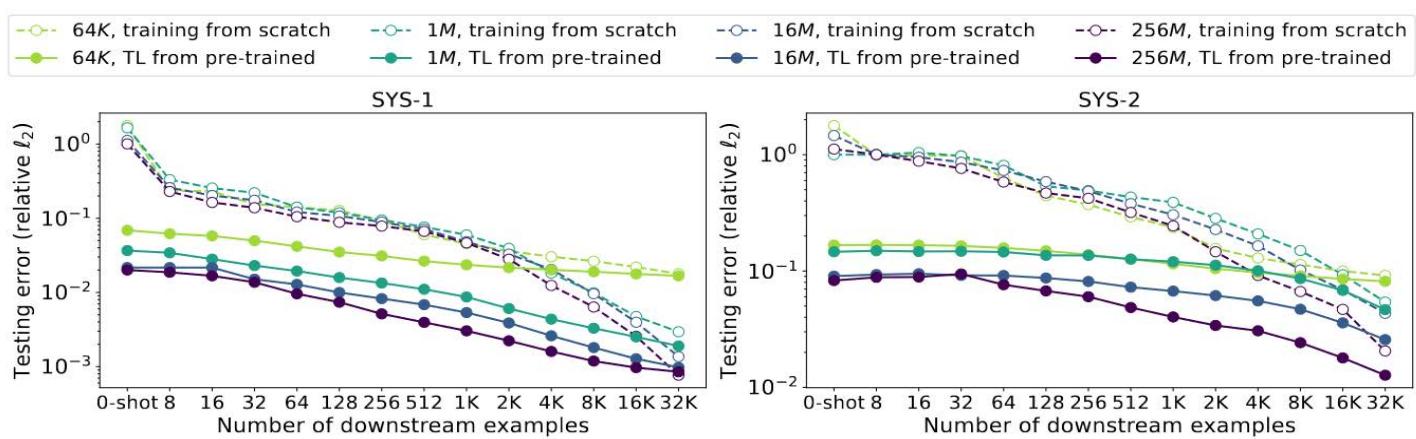
# Related work



Input Function  $\Rightarrow$  Output Function



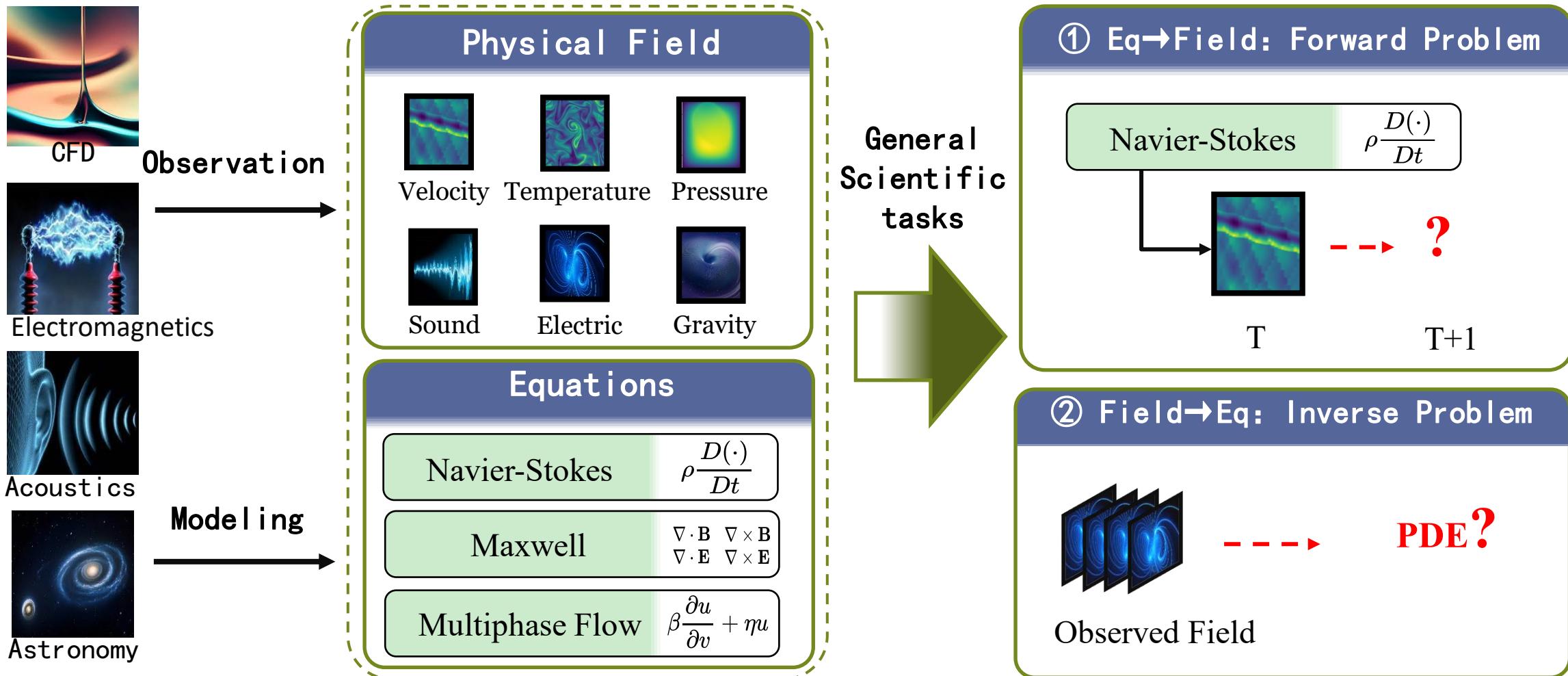
Pretrain + Finetune



Subramanian S, Harrington P, Keutzer K, et al. Towards Foundation Models for Scientific Machine Learning: Characterizing Scaling and Transfer Behavior[J]. arXiv preprint arXiv:2306.00258, 2023.

# Building Foundation Model for Scientific Computing

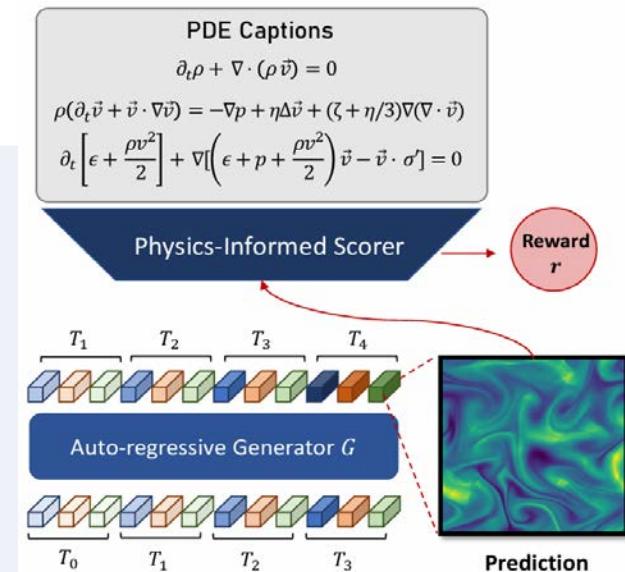
## Foundation model



# OmniArch: Physics-informed RL

## Physics-Informed Reinforcement Learning (PIRL)

- 目标：将PDE的文本描述与对应物理系统的观测对齐
- 方法：利用对比学习，构建PDE文本描述-物理场的“CLIP”模型
- 难点：PDE文本描述 (<0.1K) 数据量与物理场观测点 (>1M) 极度不匹配



## PDE Captions Augmentation (PDECap Aug)

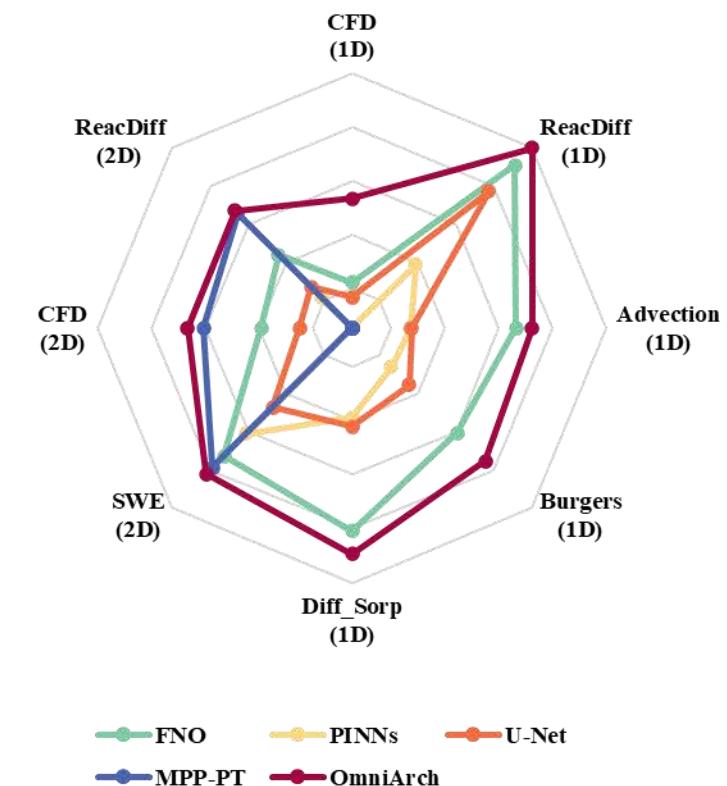
- 目标：在不改变方程语义的情况下，尽可能增加数据多样性
- 方法：规则 + 生成 + GPT-4 Check
- 效果：在PDE方程的基础上，合成约100K不同的PDE方程描述

步骤	方法名	具体操作
1	方程重写 (Equation Rewriting)	运用数学恒等式，对方程进行等价变换、
2	形式转换 (Form Transformation)	在微分形式和积分形式之间转换方程表达，应用格林函数等技术拓展方程的表示形式。
3	线性组合 (Linear Combination)	通过线性组合导出新的方程变体
4	符号替换 (Symbol Substitution)	系统地将变量符号替换为其他等价符号(如将x 替换为ξ),确保符号一致性和无歧义性。
5	专家物理检查 (Physical Checking)	由GPT-4驱动的智能专家评估增强后的方程，过滤掉那些不符合物理定律的非法实例。

# Results

**Achieves the SOTA performance:** (1) Previous foundation-like models like MPP-2D, PDEformer-1D often only supported a certain dimension. (2) By jointly training in 1D, 2D, and 3D, OmniArch surpasses over 40 PDE-specific models such as FNO, PINNs, U-Net with a single model.

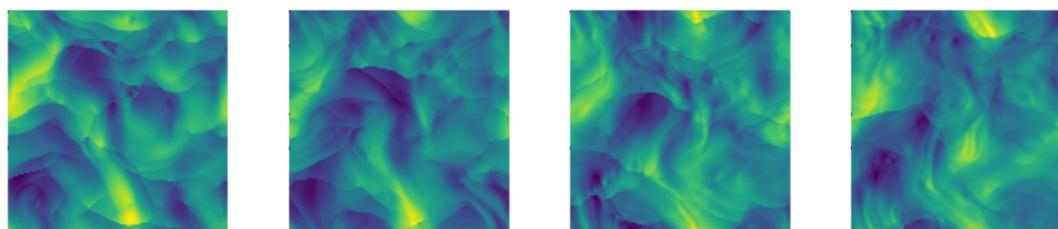
PDE	FNO	PINNs	U-Net	MPP	OmniArch	OmniArch + PIRL	相比SOTA提升	
1D	<b>CFD</b>	1.4100	-	2.6700	-	0.0981	0.0392	<b>97.2%</b>
	<b>ReacDiff</b>	0.0005	0.2140	0.0026	-	0.0004	0.0002	<b>64.0%</b>
	<b>Advection</b>	0.0091	0.8130	0.7760	-	0.0081	0.0045	<b>49.9%</b>
	<b>Burgers</b>	0.0174	0.9450	0.3200	-	0.0067	0.0032	<b>81.9%</b>
	<b>Diffusion-sorption</b>	0.0017	0.2200	0.1500	-	0.0019	0.0006	<b>62.2%</b>
2D	<b>CFD</b>	0.2060	-	1.0700	0.0178	0.0994	0.0153	<b>14.0%</b>
	<b>SWE</b>	0.0044	0.0170	0.0830	0.0022	0.0031	0.0015	<b>3.32%</b>
	<b>ReacDiff</b>	0.1200	1.6000	0.8400	0.0098	0.0993	0.0084	<b>14.0%</b>
	<b>NS-Incom</b>	0.2574	-	1.1200	-	0.1494	0.0827	<b>67.9%</b>
3D	<b>CFD</b>	0.3050	-	0.6150	-	0.6120	0.5072	-



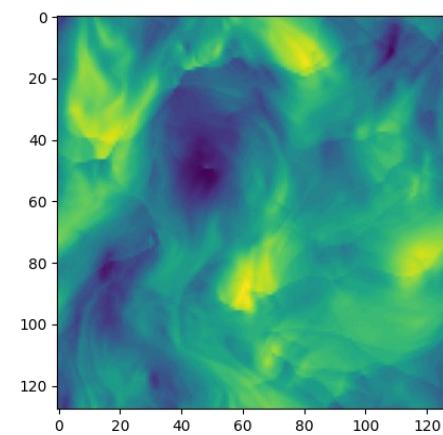
# Results

**In-Context Learning PDE:** (1) OmniArch learns physics operators directly from the given pre-processed physical field variations to dynamically handle out-of-distribution (OoD) physical systems without the need for retraining. (2) In physical systems such as Advection, CFD, SWE, we observe that as the length of the input physical field sequence context increases, the solution accuracy steadily improves.

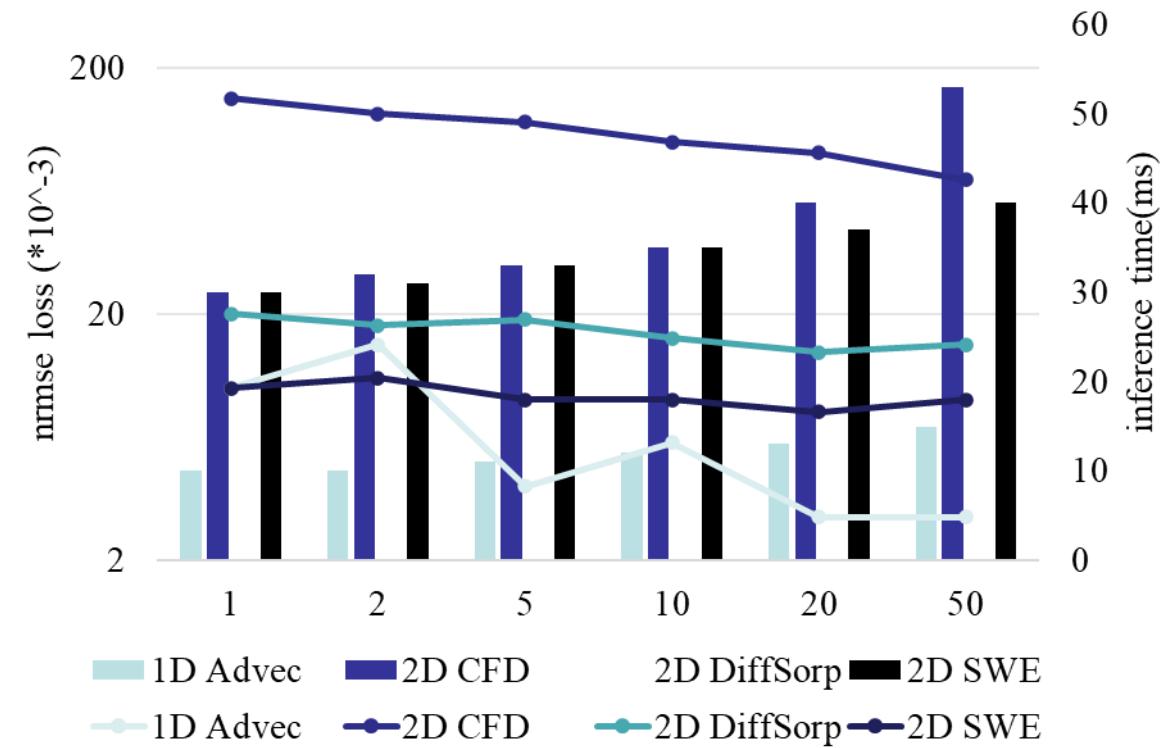
2DCFD input context (1-4s) :



2DCFD Predictions (5-20s) :

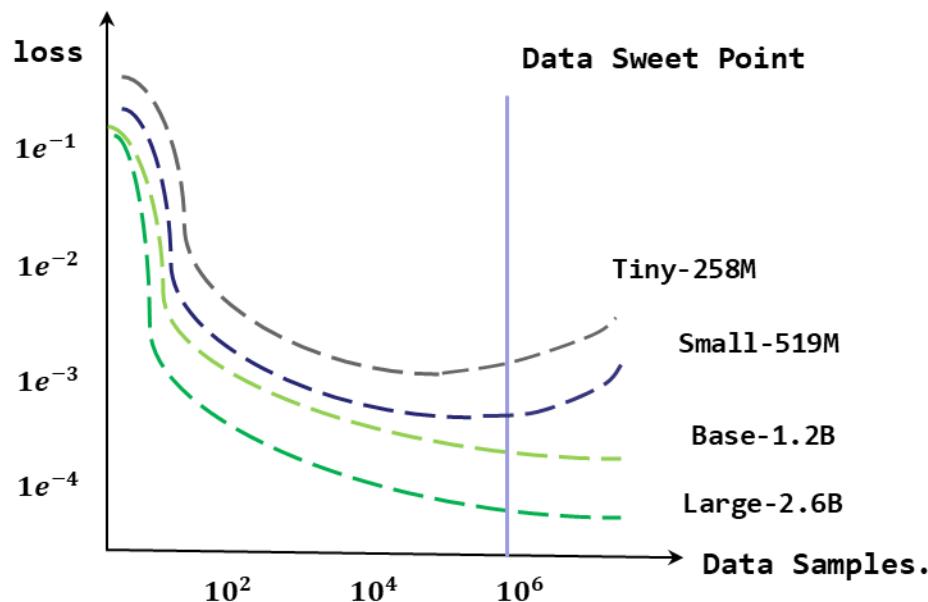


Time/loss ~ context length



# Results

**Balance of training efficiency and model size:** (1) The scaling-law: From 30M to 0.5B, the model's loss function grows proportionally to the size of the parameters. When the parameters exceed 1B, the model no longer overfits. (2) Data bottleneck: The growth in model performance depends on the increase in data points. Currently, the scale of open-source data is around  $10^6$ , limiting the maximum size of the model to about 2.6B.

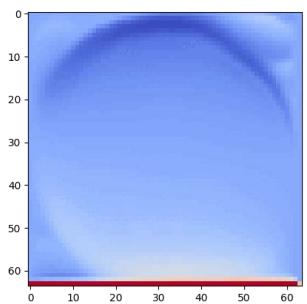


Model	Layers	Hidden	Intermediate	Heads	Params	GPU Hrs
Tiny	12	192	768	3	258M	267.6
Small	12	384	1536	6	519M	411.2
Base	12	768	3072	12	1.2B	608
Large	24	1024	4096	16	2.6B	776

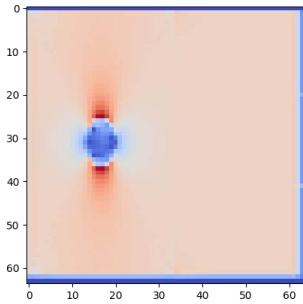
# Results

**Downstream tasks finetuning + few-shot testing:** (1) Fine-tuning is performed on four commonly used fluid equations: Cavity, Cylinder, Dam, and Tube. (2) Compared to models like MPP, the OmniArch model can achieve rapid improvement in accuracy with around 50 samples.

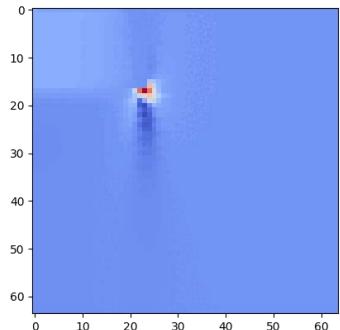
方腔流 (Cavity)



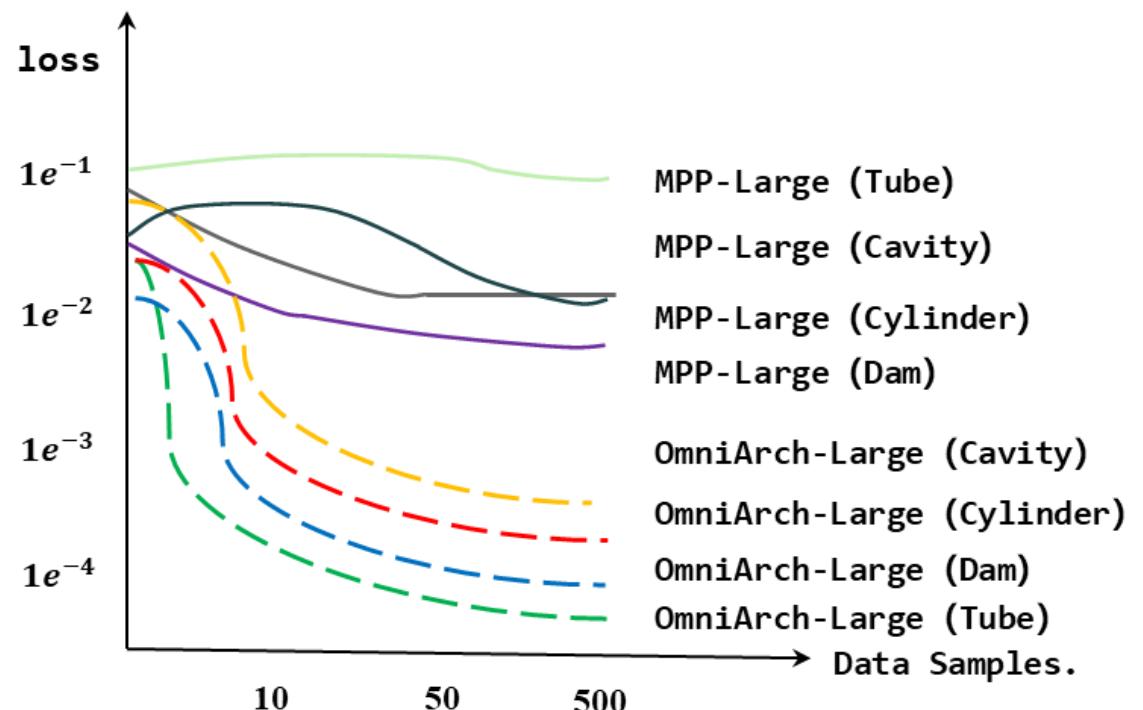
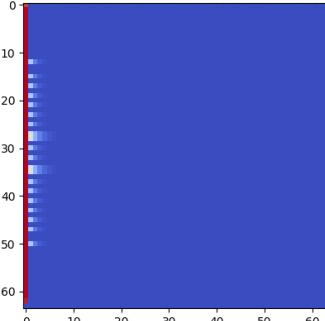
圆柱绕流 (Cylinder)



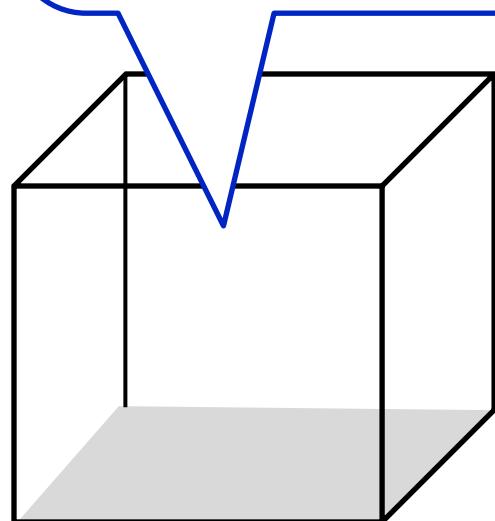
堤坝流 (Dam)



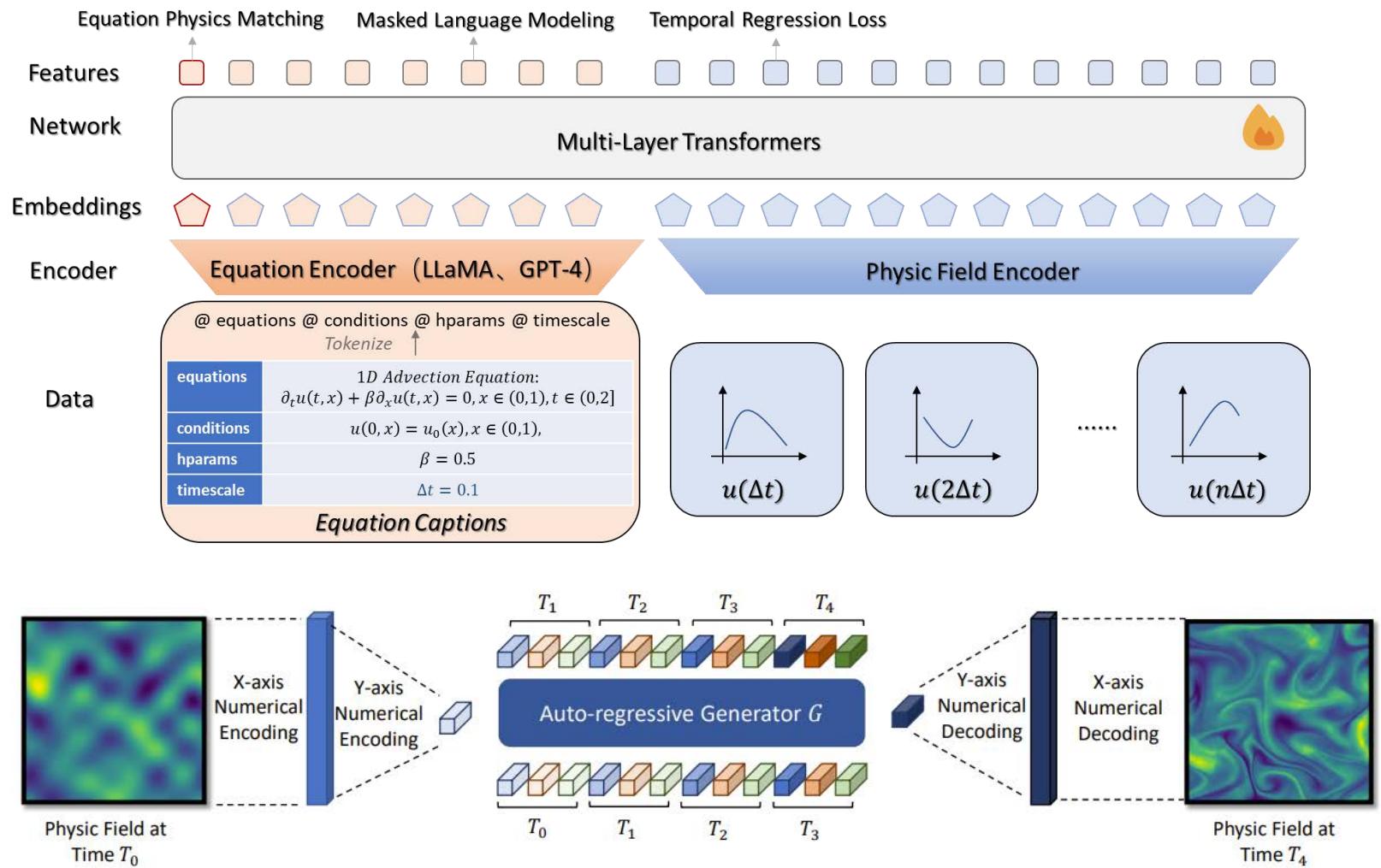
圆管流 (Tube)



# Aiming for Temporal Grounding



Temporal Knowledge



ProbSparse Attention

<https://arxiv.org/pdf/2402.16014.pdf>

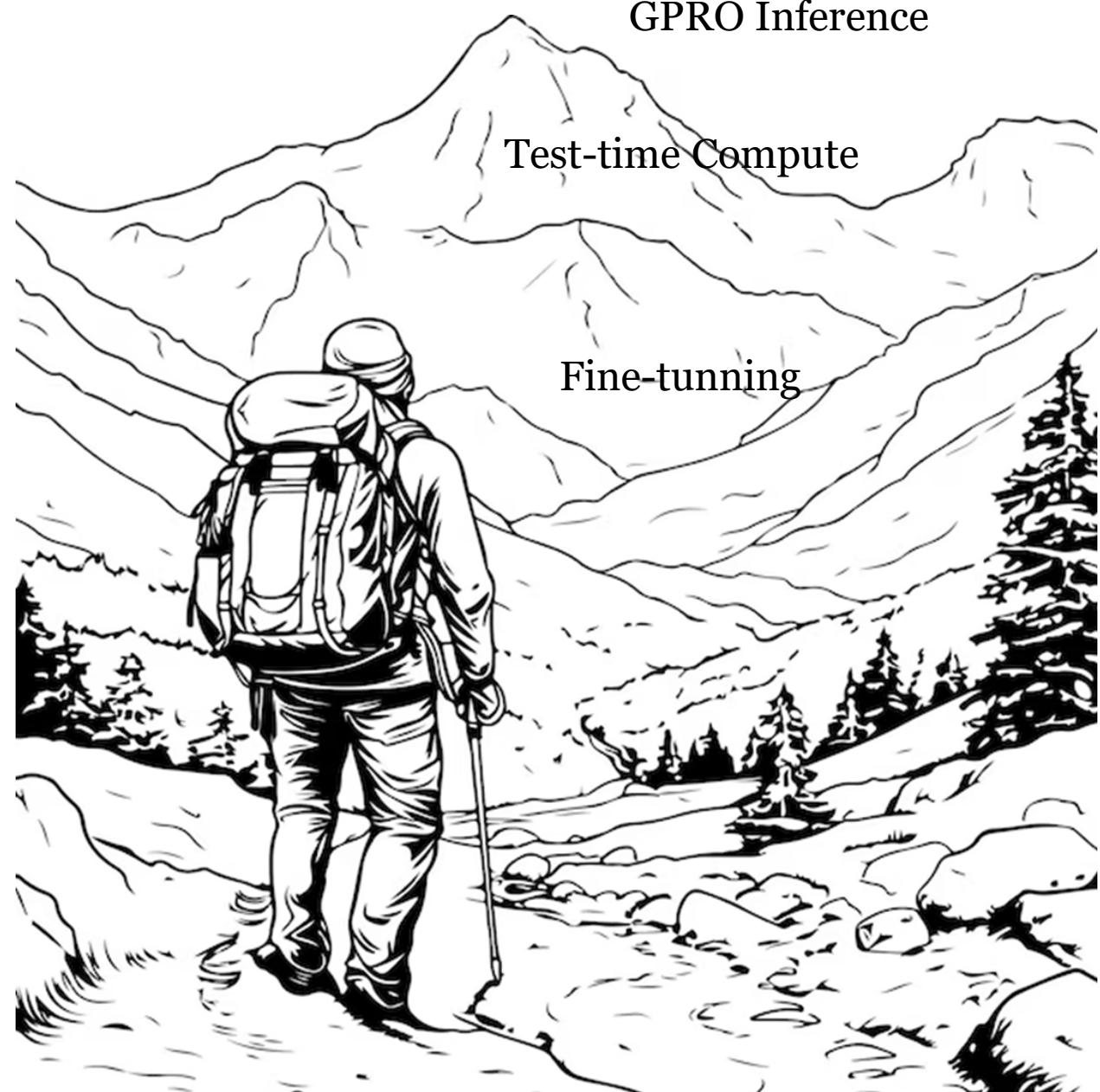
# - Content -

- Sequential Modeling: from Science
- The Transformer Model
- The modern architecture
- Some applications on scientific sequences
- **Future direction**

# Summary

**AI4S:** on the way!

**Time-series:** only start!



# Summary

70% Data

20% Computing

10% Architecture



**Unstable Trigon:** What is the next sequential architecture?

# Thanks!

Haoyi ZHOU, [www.zhouhaoyi.com](http://www.zhouhaoyi.com)



Beihang University