



Timer: Generative Pre-trained Transformers Are Large Time Series Models

Yong Liu^{*1} Haoran Zhang^{*1} Chenyu Li^{*1} Xiangdong Huang¹ Jianmin Wang¹ Mingsheng Long¹



Yong Liu



Haoran Zhang



Chenyu Li



Xiangdon Huang

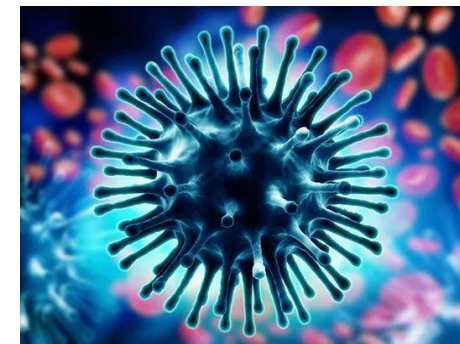


Jianmin Wang

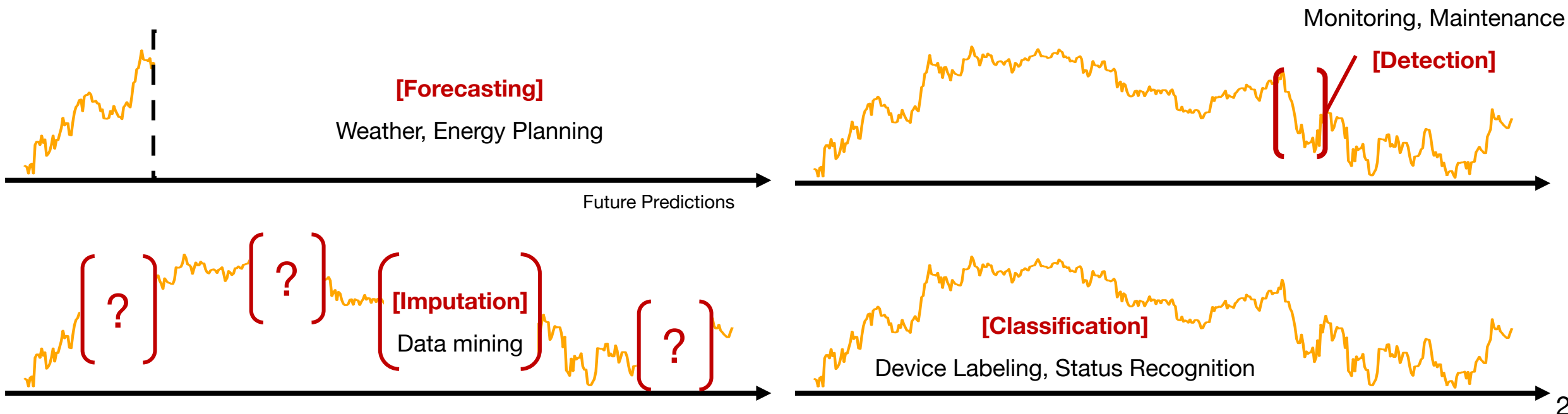


Mingsheng Long

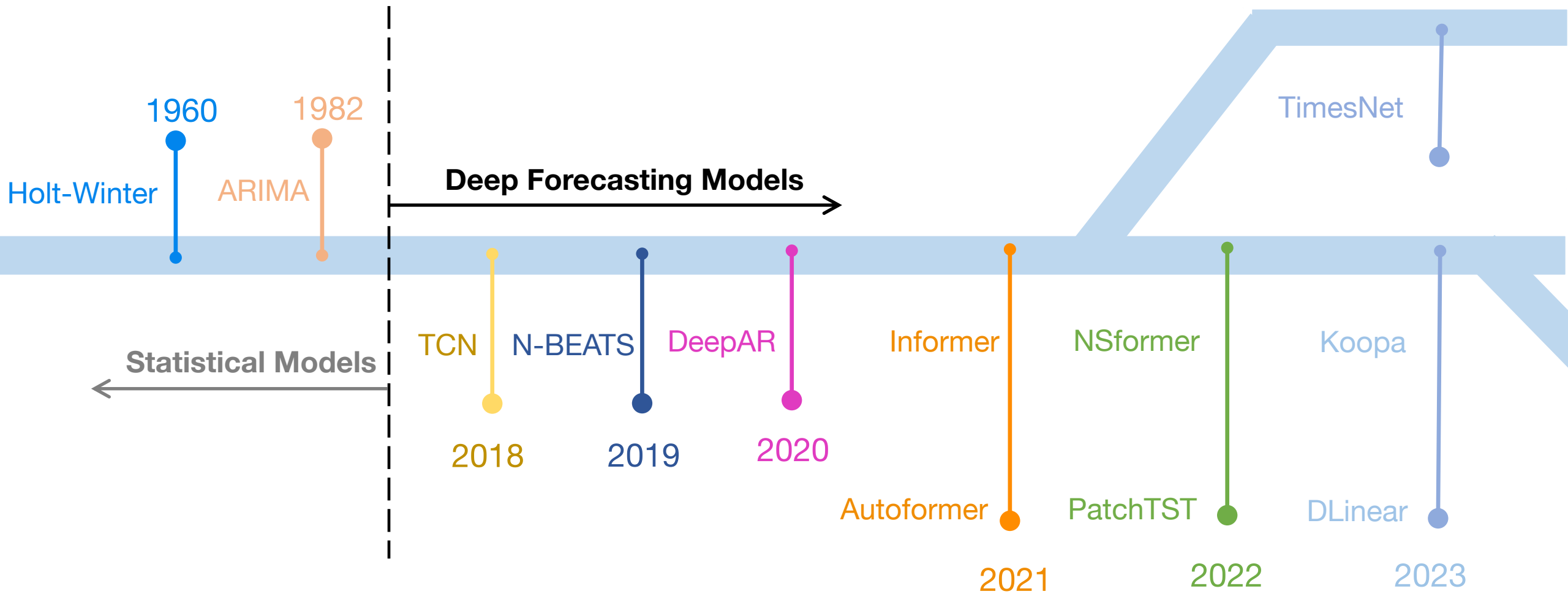
Time Series Applications



Time Series Analysis is Ubiquitous in Real World

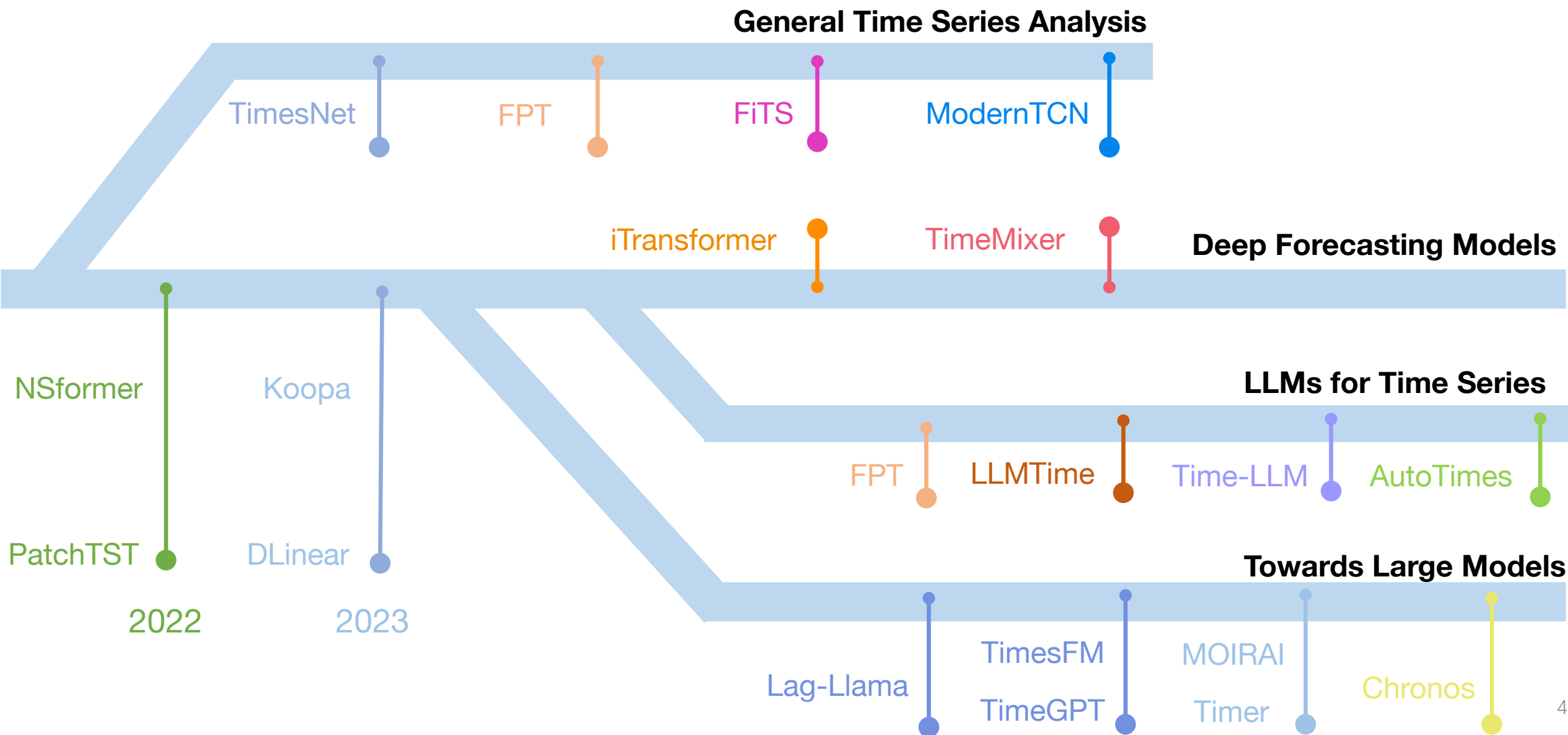


Deep Models for Time Series





Deep Models for Time Series



Motivations of Large Models

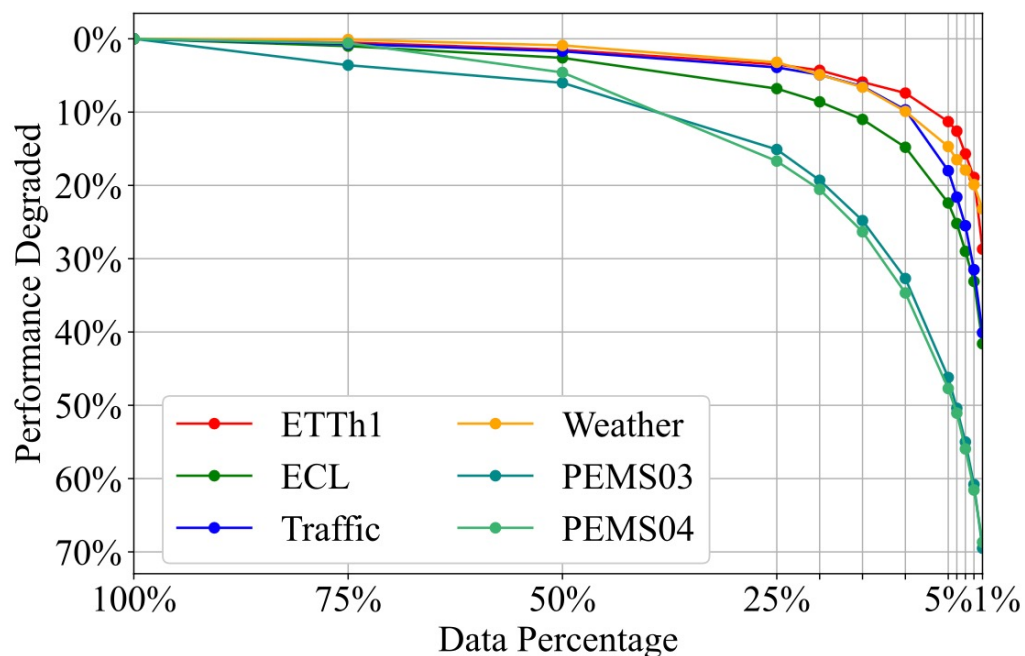


Status quo: Costly training small models in specific scenarios (tasks, datasets, settings)

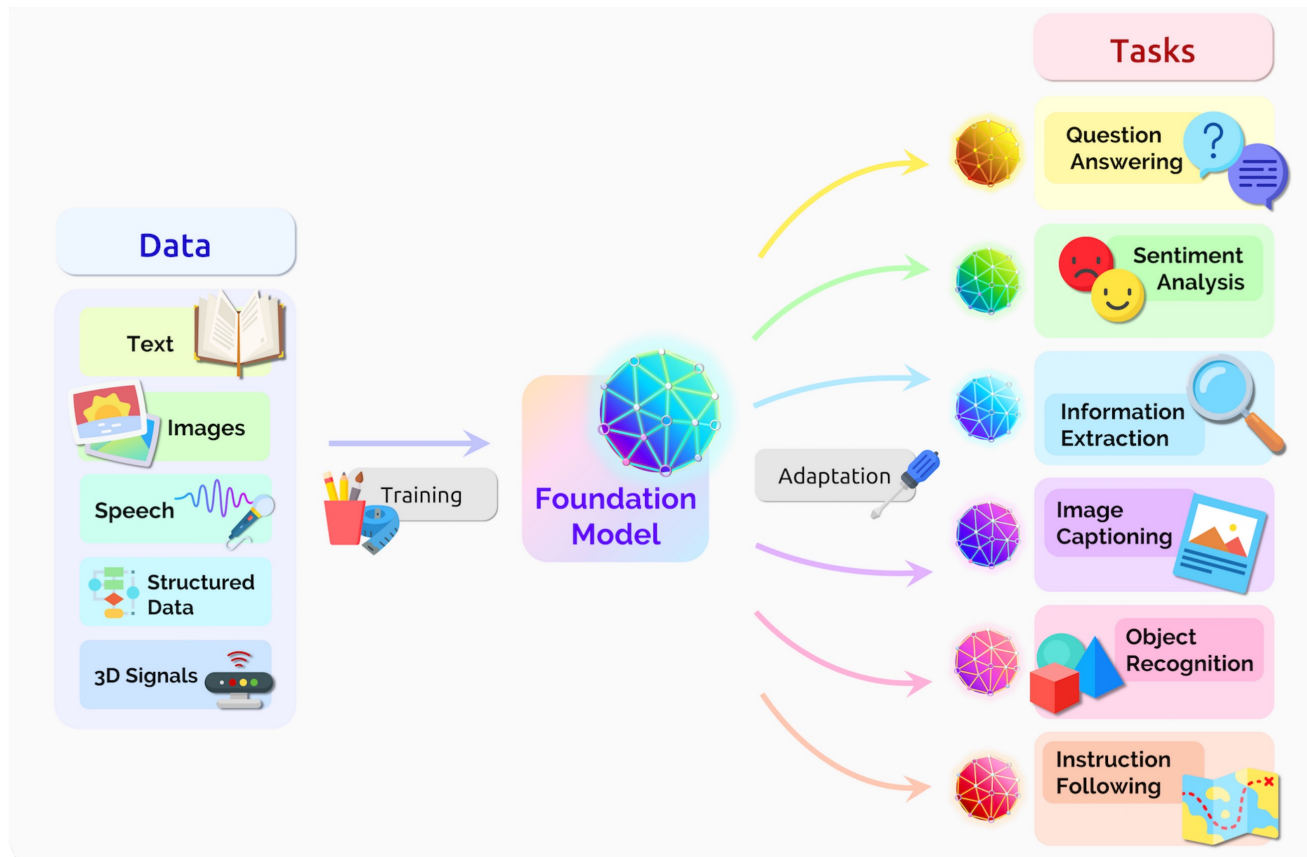


Data scarcity is common in real-world applications

- Real scenarios may lack training samples
- Performance of SOTA model degrades with limited data



Motivations of Large Models



😊 **[Data Universal]**

Learn from various datasets

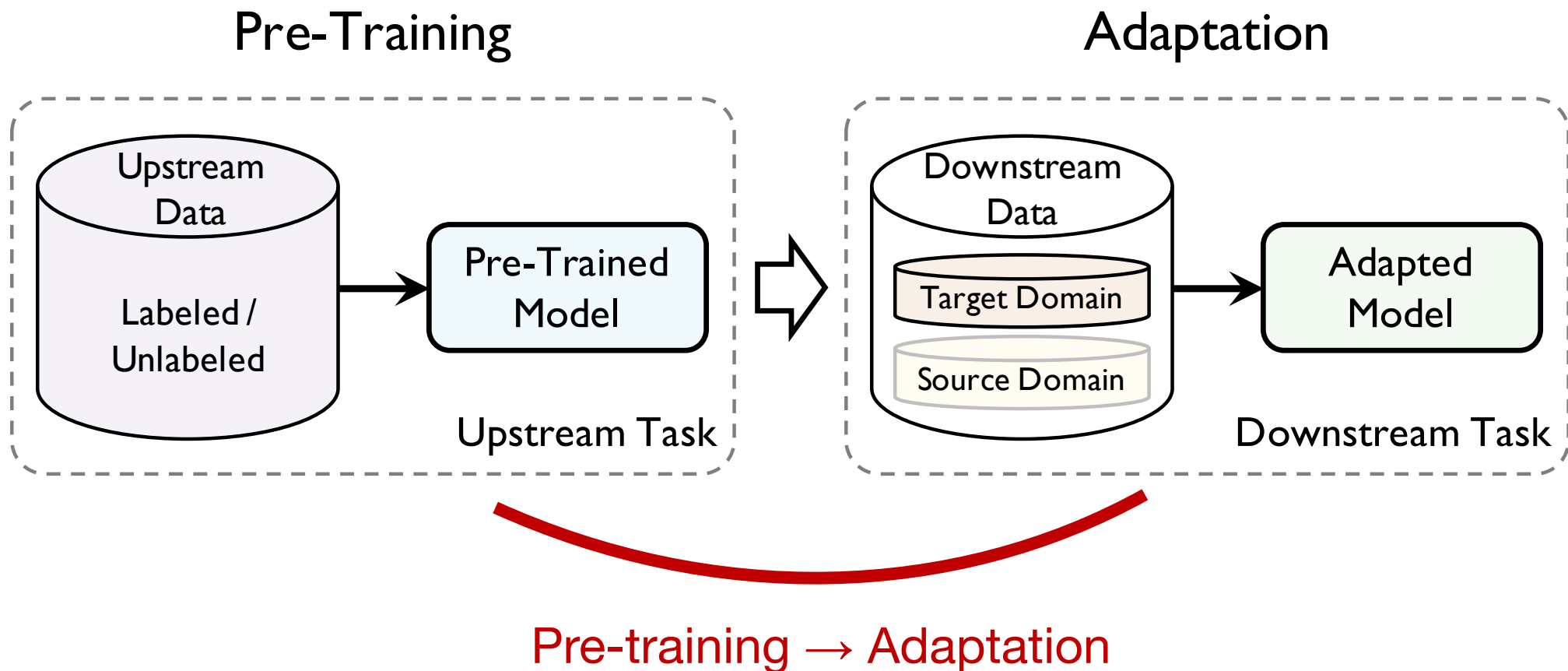
😊 **[Task Universal]**

Adapt to a wide range of
downstream tasks

Large Time Series Models

What is Large Model

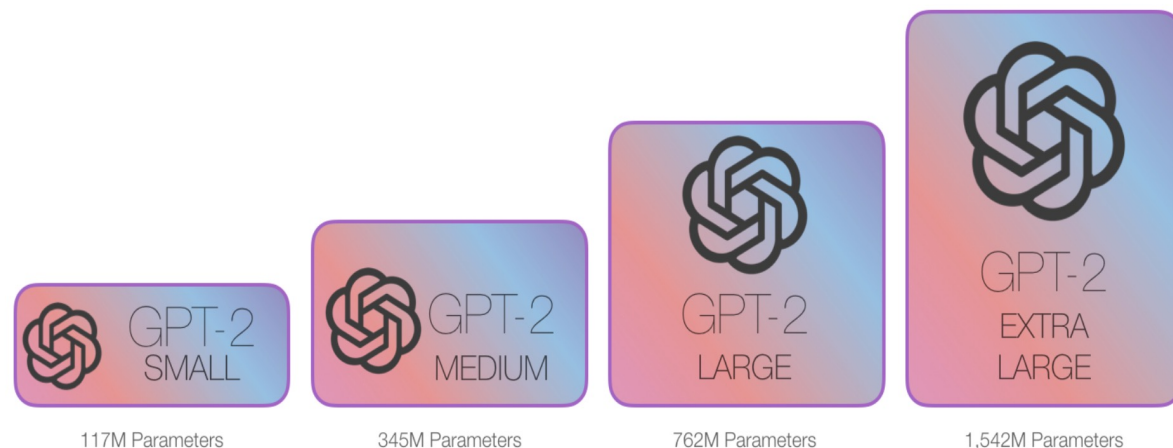
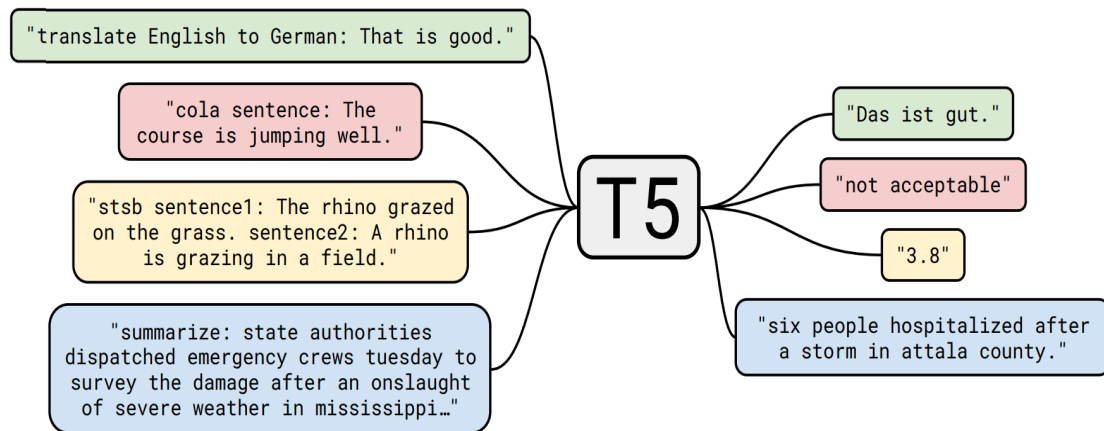
- **Generalizability:** One model fits different domains



Large Time Series Models

What is Large Model

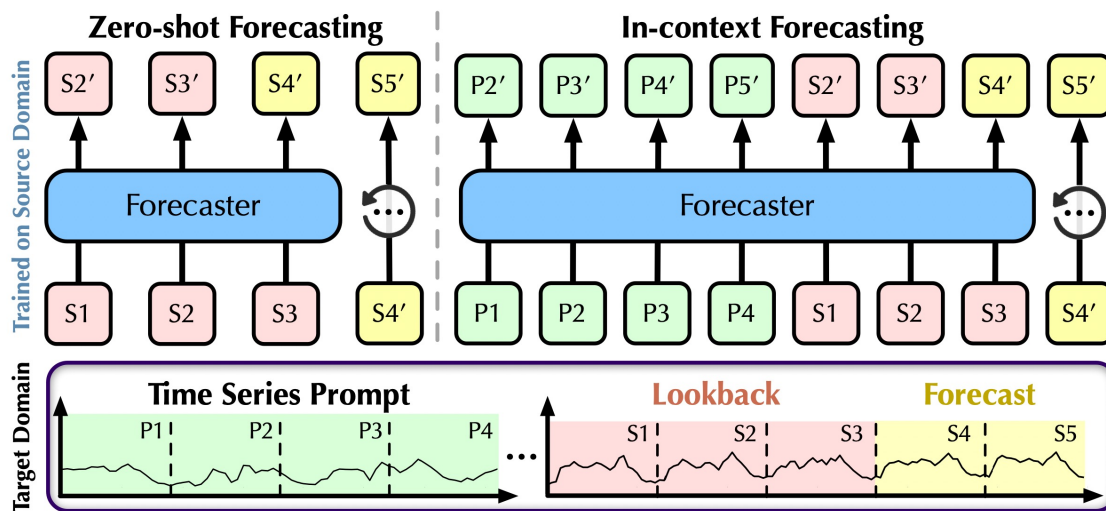
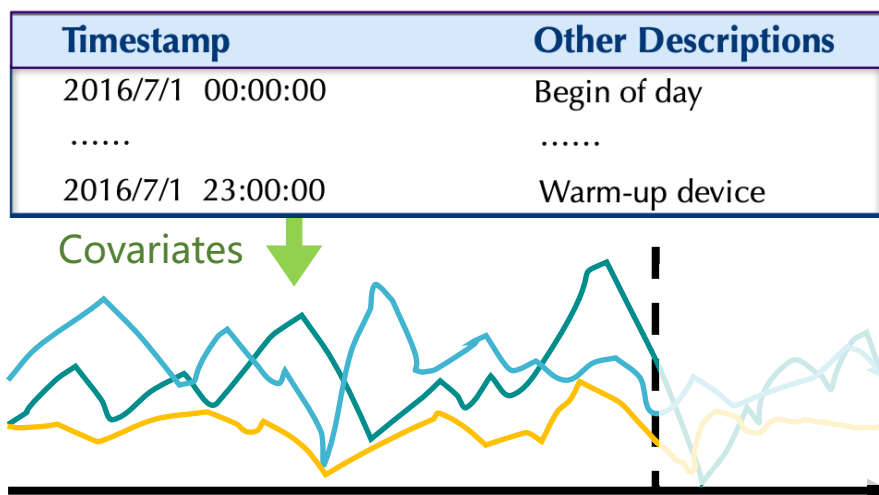
- **Generalizability:** One model fits different domains
- **Task Generality:** Versatility to cope with various scenarios/tasks
- **Scalability:** Performance improves with the scale of pre-training



Large Time Series Models

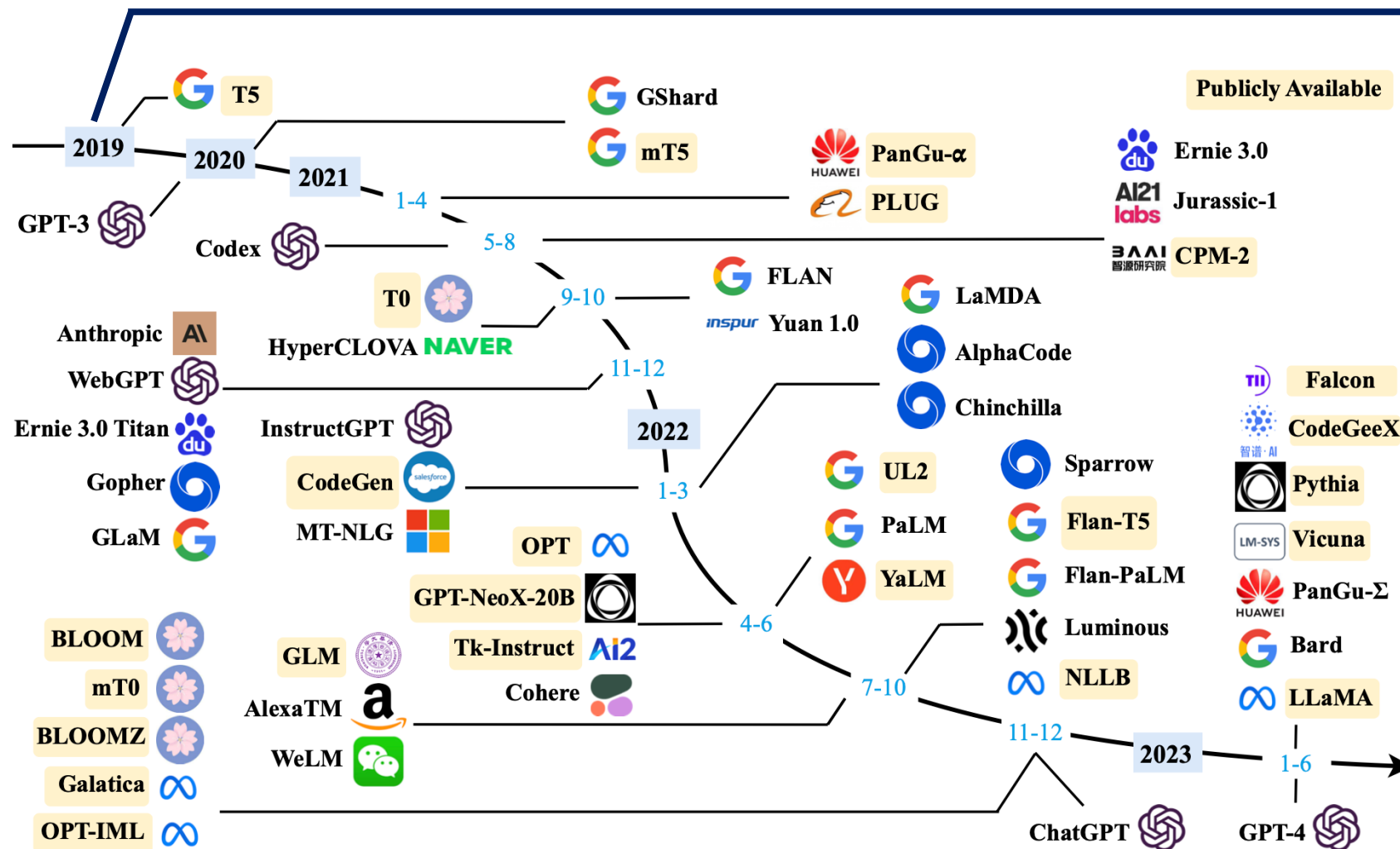
What is Large Model

- **Generalizability:** One model fits different domains
- **Task Generality:** Versatility to cope with various scenarios/tasks
- **Scalability:** Performance improves with the scale of pre-training
- **Emergence Abilities:** Multimodality, In-context Learning ...





Timeline of Large Language Models



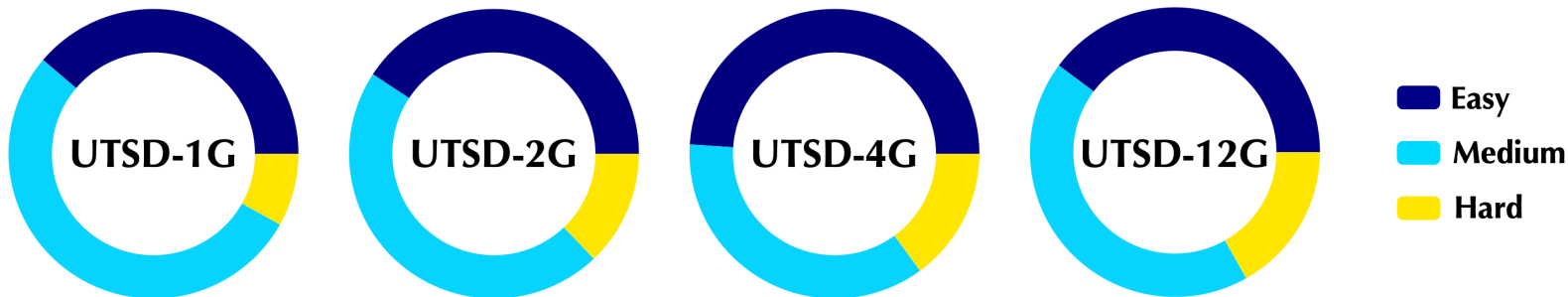
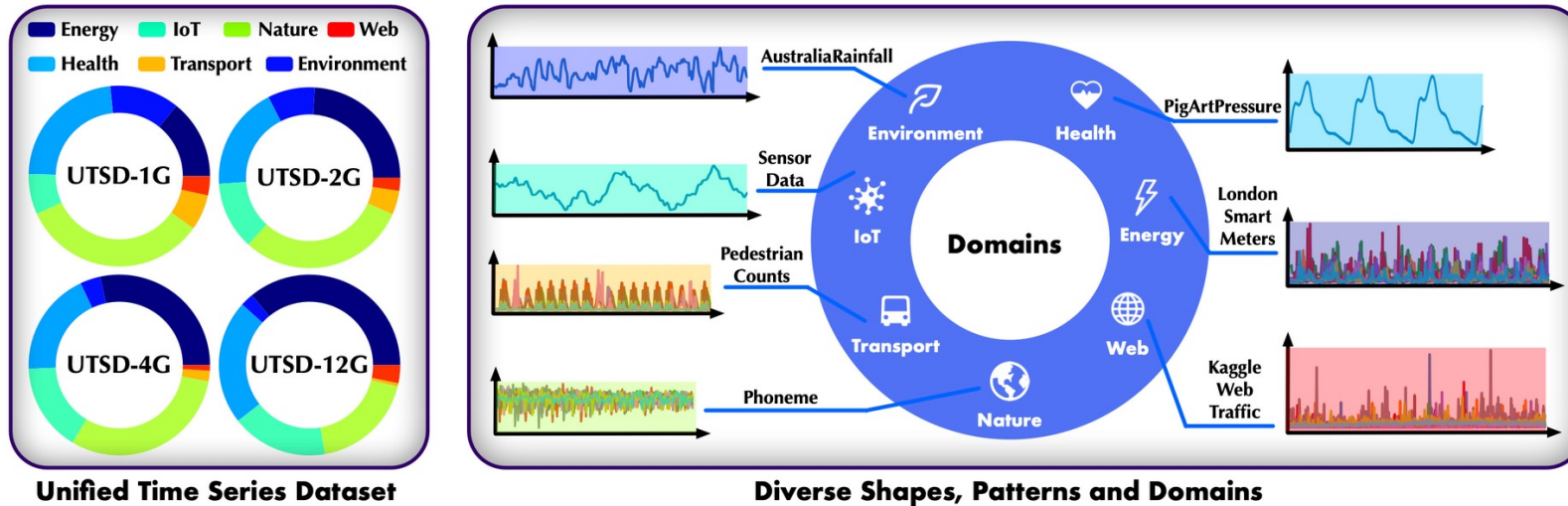
Large Time Series Models Are Still in Early Stages

Challenges

- Data Infrastructure
- Scalable Architecture
- Task Heterogeneity

Timer: Well-curated Datasets

❑ UTSD: Unified Time Series Dataset



Dataset: <https://huggingface.co/datasets/thuml/UTSD>

Data Quality

- Aggregation & Filter
- Preprocess & Evaluate
- Stacking up with a hierarchy



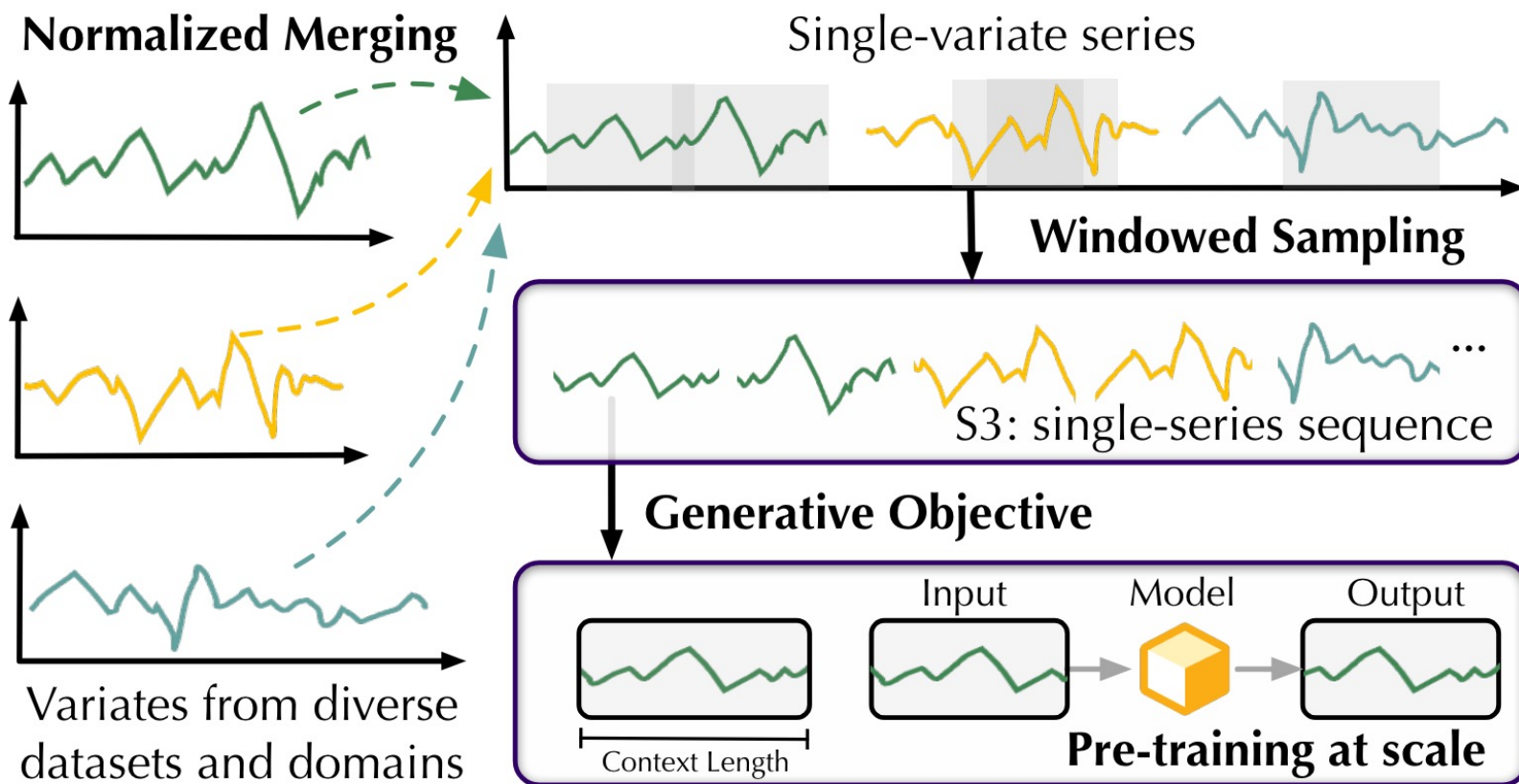
- 1 Billion Time Points
- 7 Typical Domain
- 4 Scalable Volumes
- Continuous Expansion...

Timer: Single-Series Sequence

□ Unified Format to Address **Data Heterogeneity**: single-series sentence (S3)

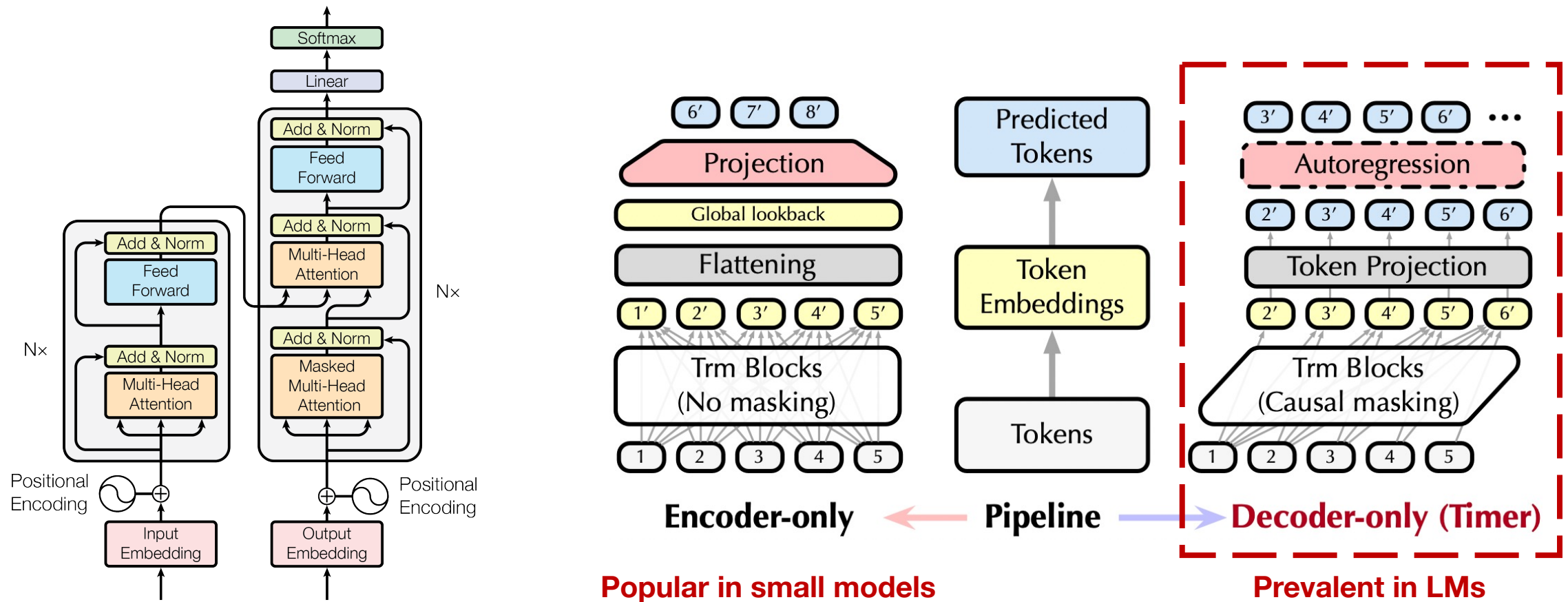
Distinct in Shape/Freq/Scale!

Dataset	Dim	Frequency
ETTh1, ETTh2	7	Hourly
ETTm1, ETTm2	7	15min
Exchange	8	Daily
Weather	21	10min
ECL	321	Hourly
Traffic	862	Hourly
Solar-Energy	137	10min
PEMS03	358	5min
PEMS04	307	5min
PEMS07	883	5min
PEMS08	170	5min



Timer: Explore Backbones for Large Model

□ Decoder-only Transformer with Autoregression



Timer: Generative Pre-training

□ Next Token Prediction (Both Training and Inference)

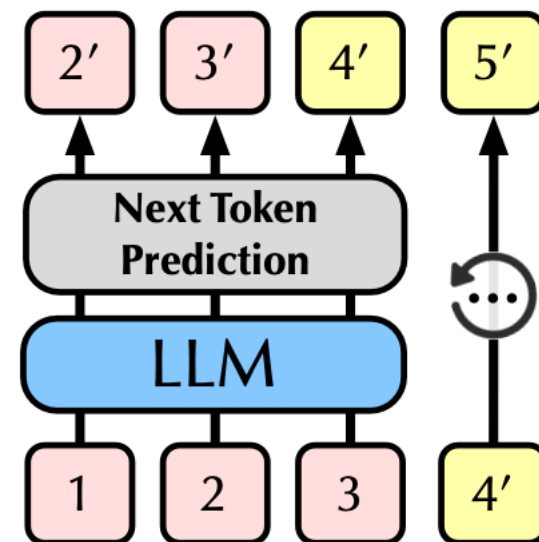
Tokenize : $\mathbf{s}_i = \{x_{(i-1)S+1}, \dots, x_{iS}\} \in \mathbb{R}^S$.

$\mathbf{h}_i^0 = \mathbf{W}_e \mathbf{s}_i + \mathbf{T} \mathbf{E}_i, i = 1, \dots, N,$

Forwarding : $\mathbf{H}^l = \text{TrmBlock}(\mathbf{H}^{l-1}), l = 1, \dots, L,$

$\{\hat{\mathbf{s}}_{i+1}\} = \mathbf{H}^L \mathbf{W}_d, i = 1, \dots, N,$

NTP : $\mathcal{L}_{\text{MSE}} = \frac{1}{NS} \sum \|\mathbf{s}_i - \hat{\mathbf{s}}_i\|_2^2, i = 2, \dots, N + 1.$

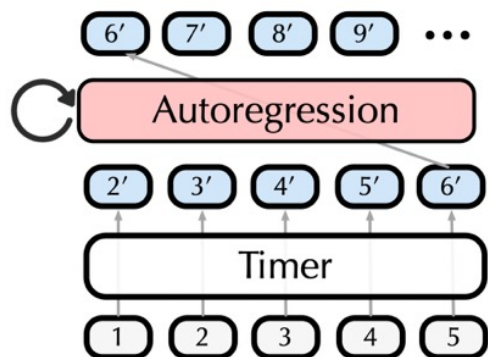
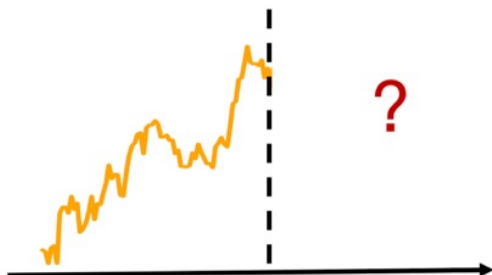


Token-wise supervision: generated token at each position is independently supervised

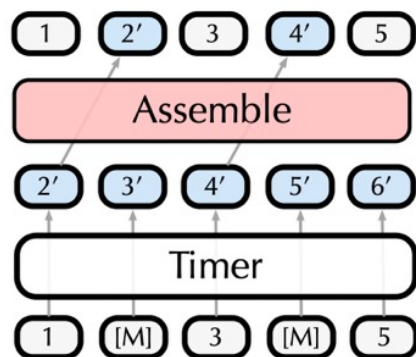
Timer: Unified Task Formulation

□ Unify Time Series Analysis into Generative Tasks

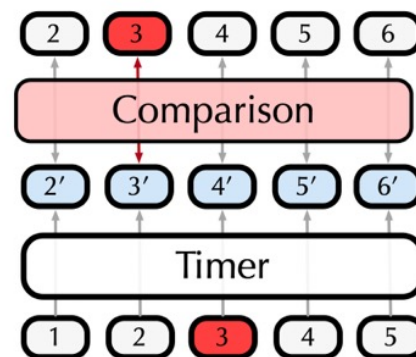
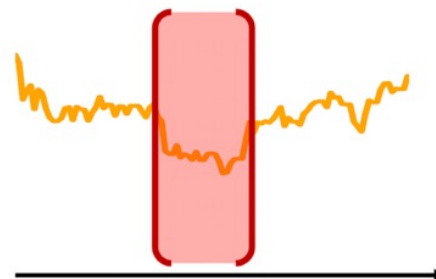
(1) Forecasting



(2) Imputation



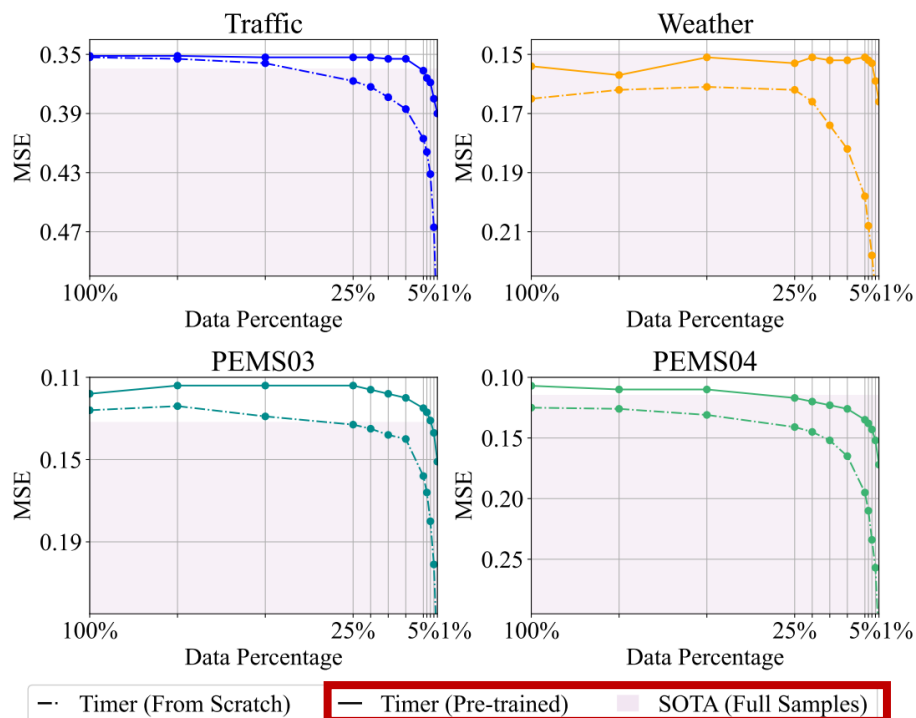
(3) Detection



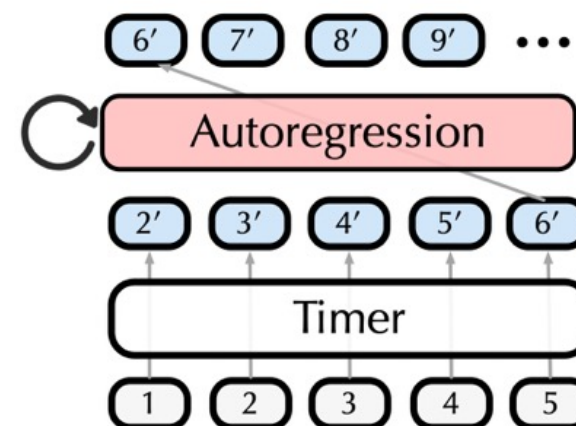
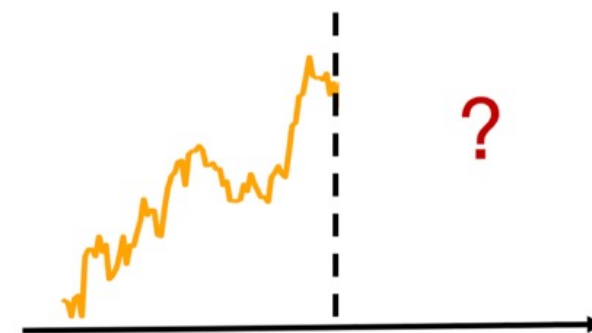
Task Generality

Time Series Forecasting

- Naturally predict the next token
- Timer trained with **1~5% samples** outperforms SOTA with **100% samples**



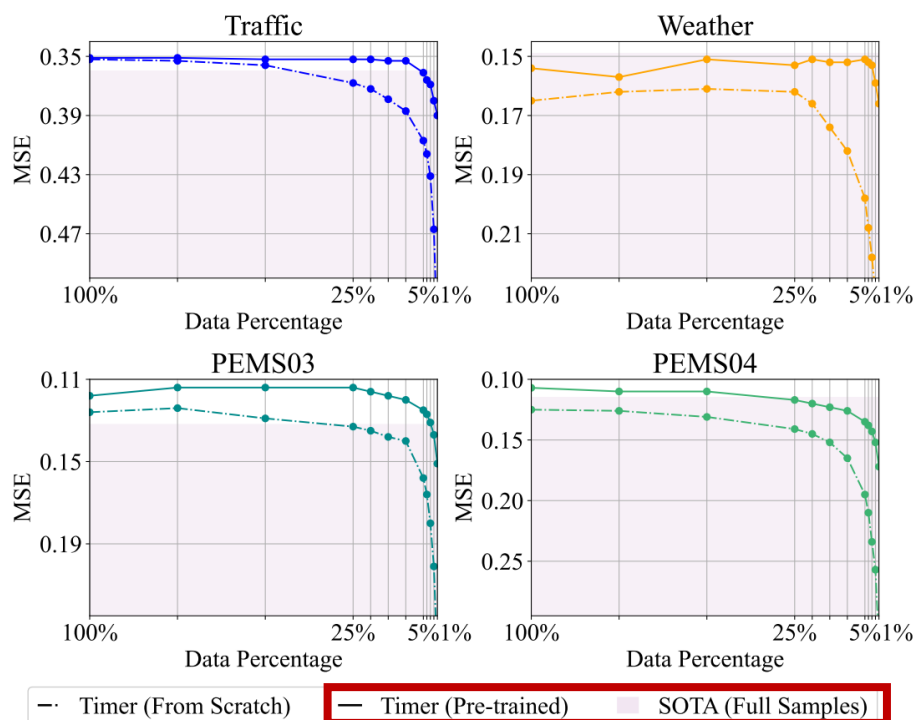
(1) Forecasting



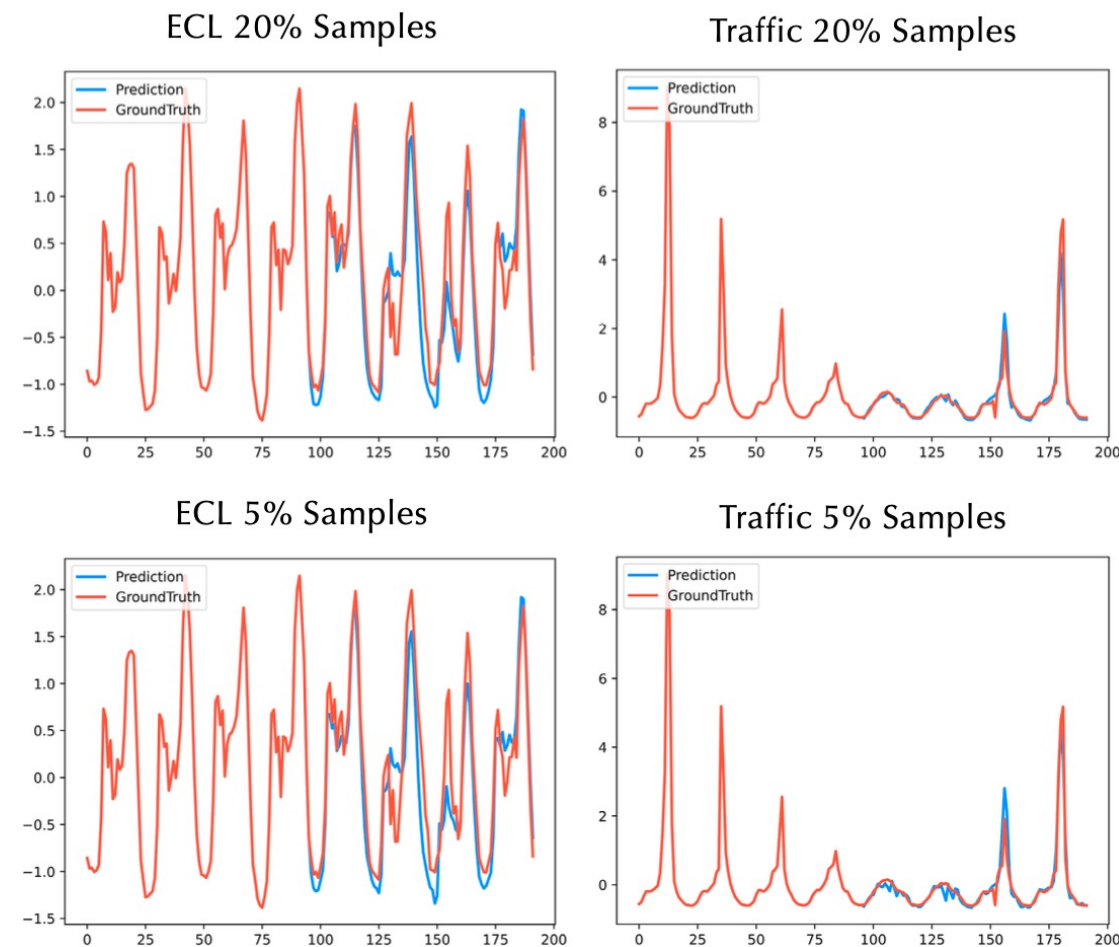
Task Generality

Time Series Forecasting

- Naturally predict the next token
- Timer trained with **1~5% samples** outperforms SOTA with **100% samples**



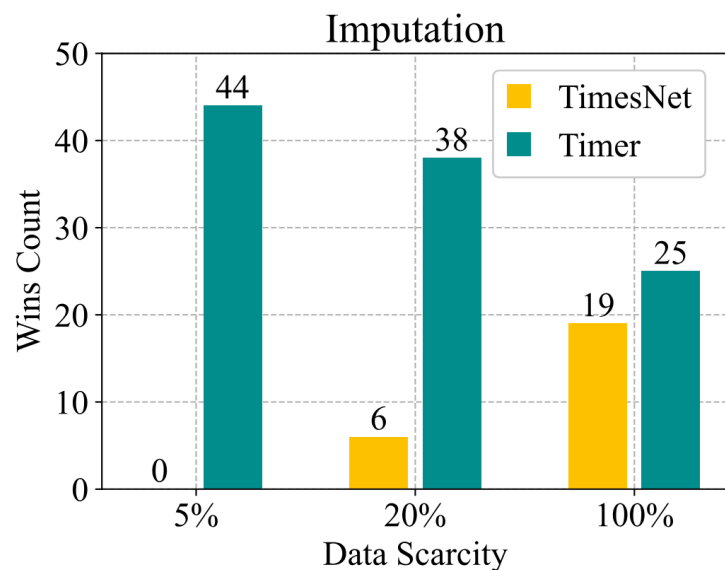
Showcases



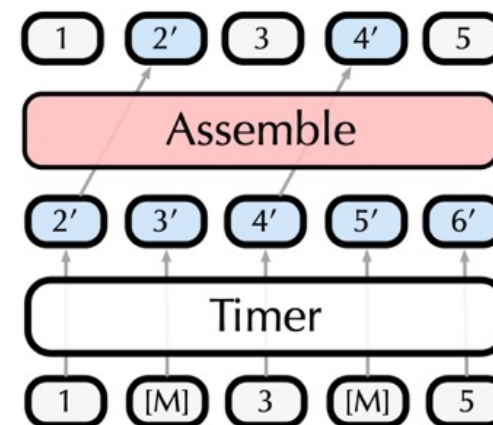
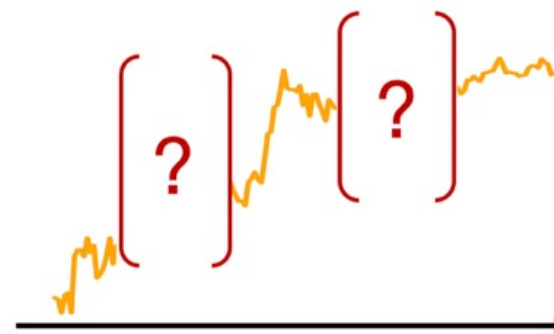
Task Generality

Time Series Imputation

- Imputation is performed by generating masked tokens with the previous context
- Surpass previous **SOTA TimesNet** in average masked cases and data scarcities.



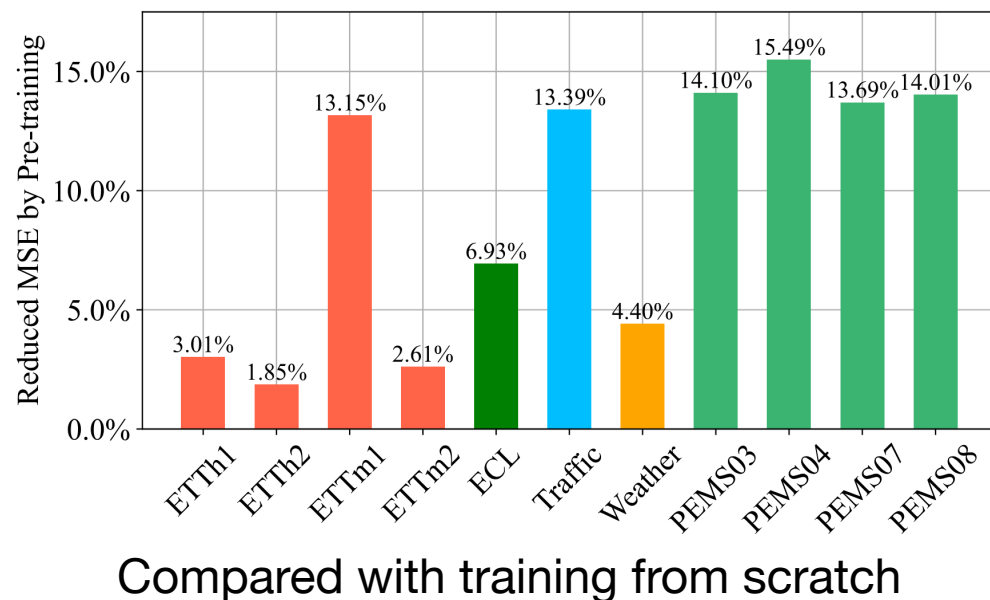
(2) Imputation



Task Generality

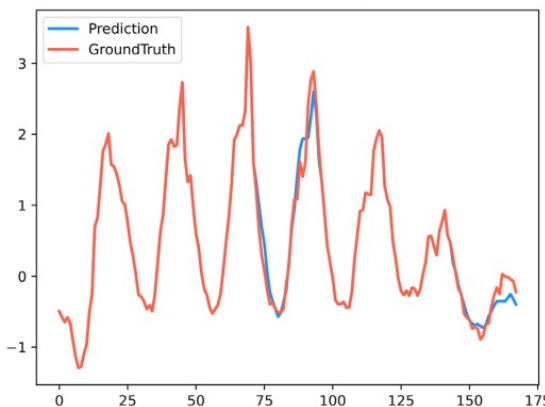
Time Series Imputation

- Imputation is performed by generating masked tokens with the previous context
- Stable improvement** exhibited in imputation by large-scale pre-training

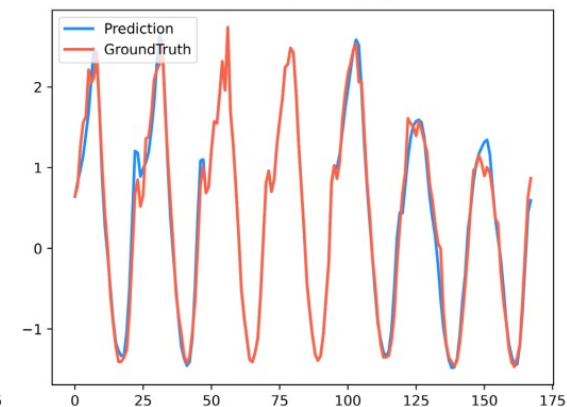


Showcases

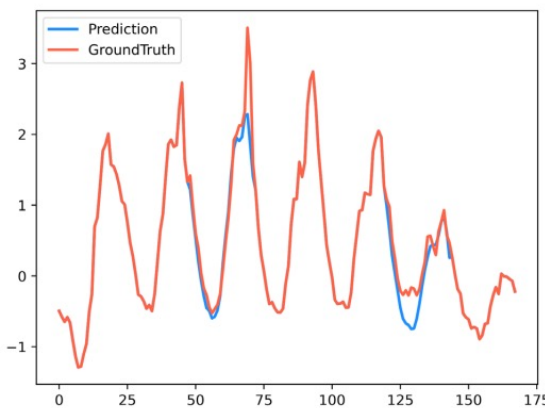
ECL 20% Samples



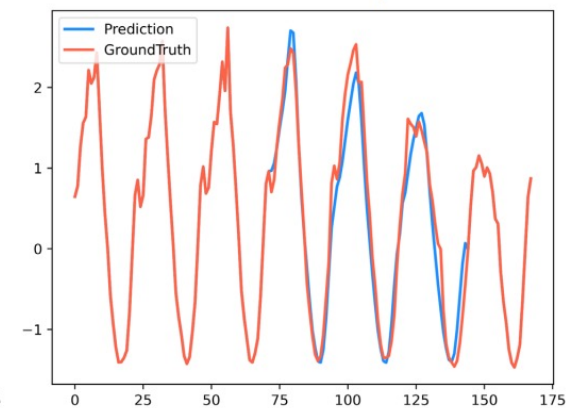
Traffic 20% Samples



ECL 5% Samples



Traffic 5% Samples

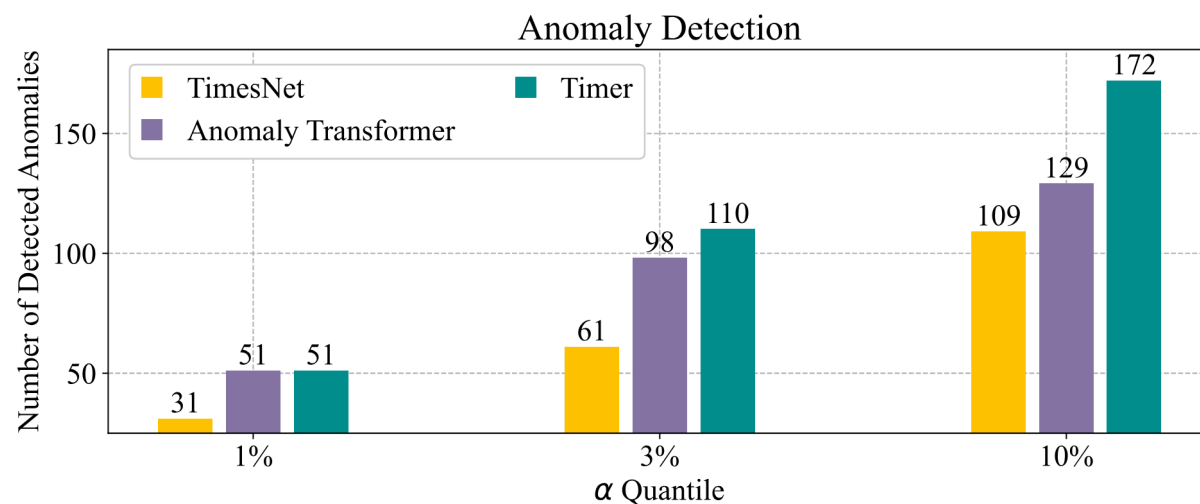
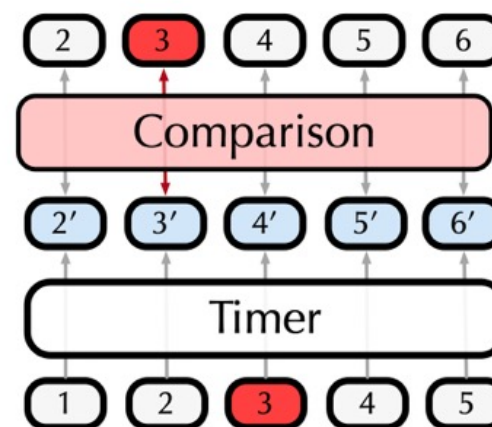
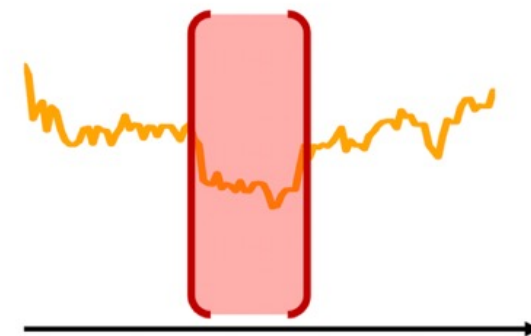


Task Generality

Anomaly Detection

- Conducted in a **predictive** approach by generating normal time series
- Quantile** the abnormal confidence in MSE
- Surpass **task-specific SOTA models** in the challenging UCR Anomaly Archive

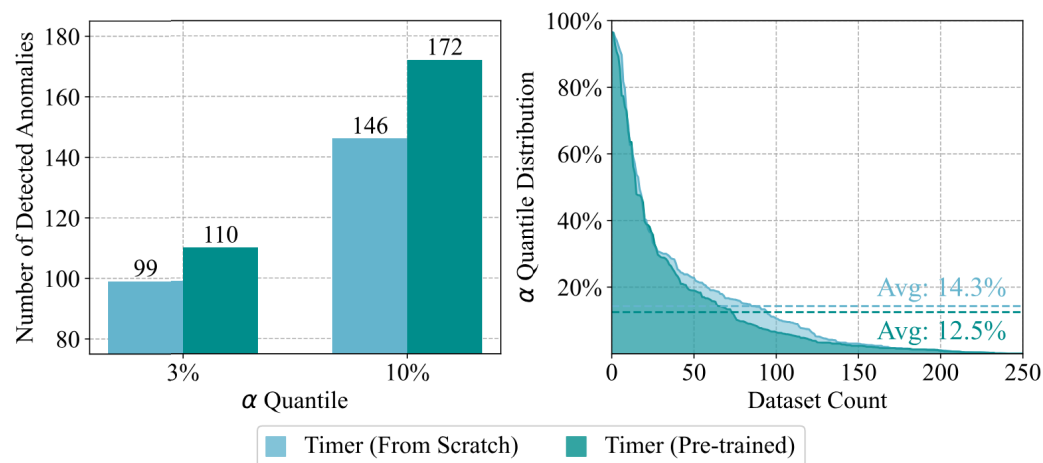
(3) Detection



Task Generality

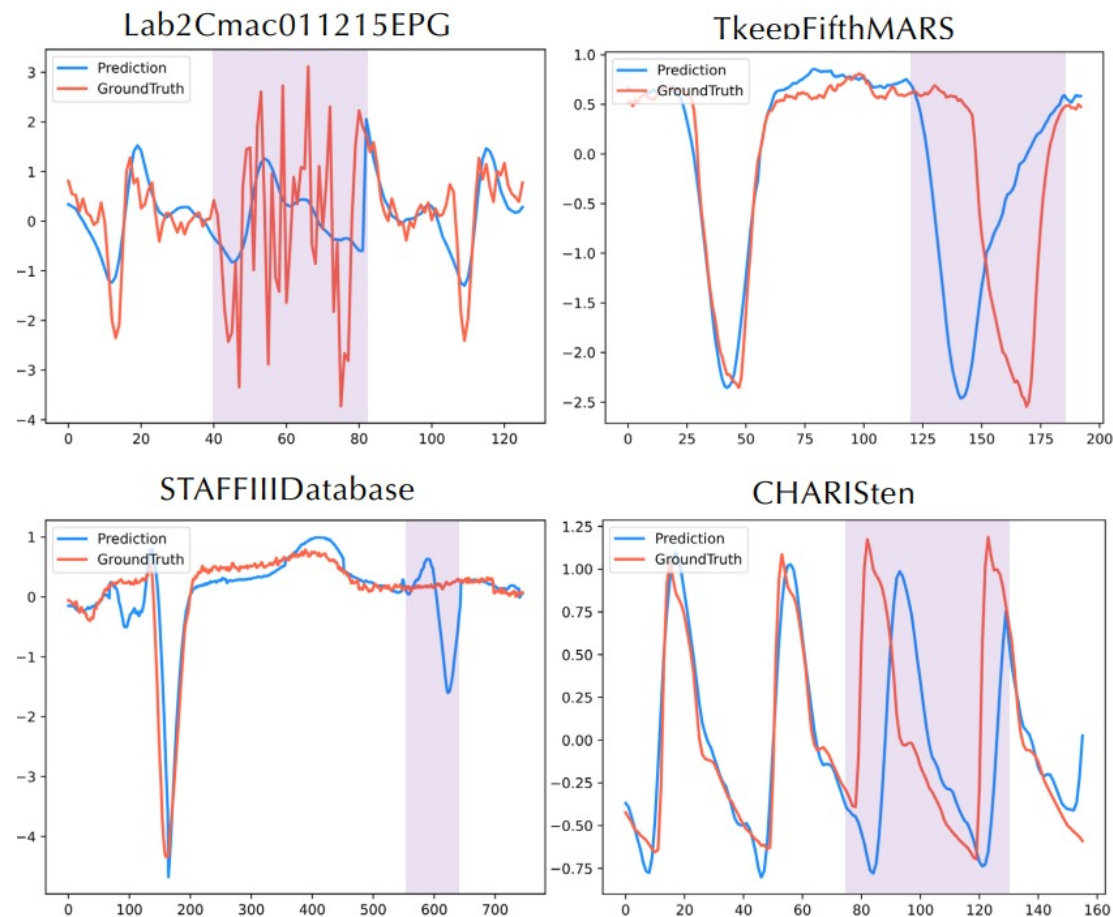
Anomaly Detection

- Conducted in a **predictive** approach by generating normal time series
- Quantile** the abnormal confidence in MSE
- Stable improvement** exhibited in anomaly detection by large-scale pre-training



Smaller α indicates better performance

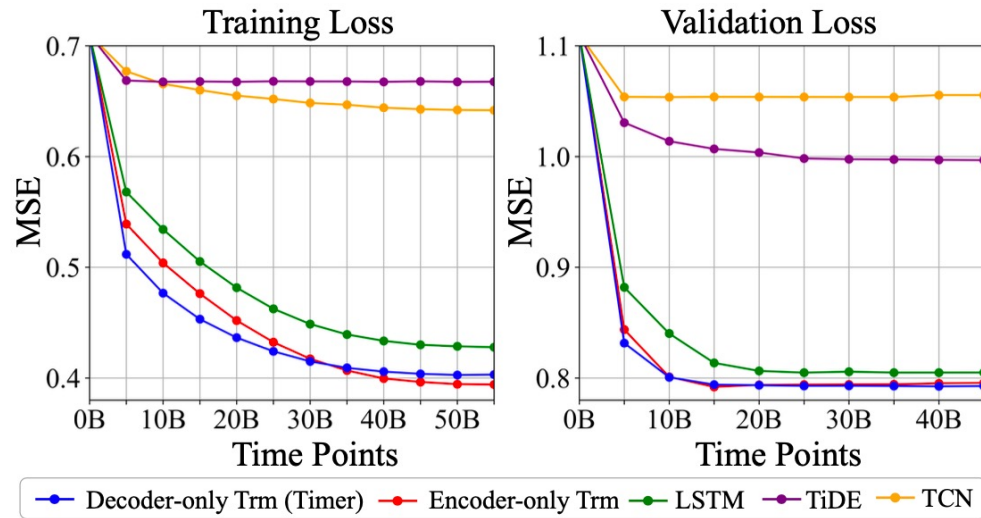
Showcases



Scalability



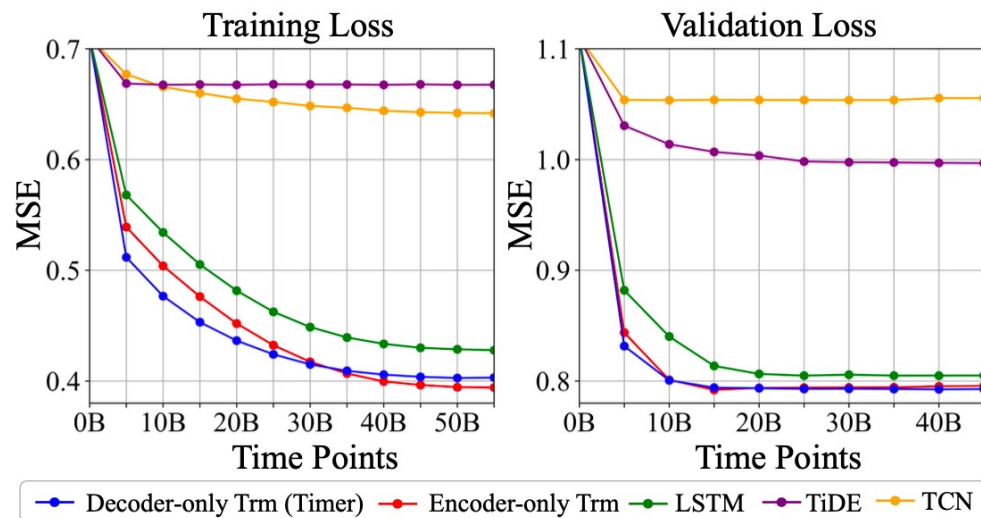
Loss Curve of Sequence Models on UTSD



**Transformer exhibits model capacity as the
scalable architecture for LTSM**

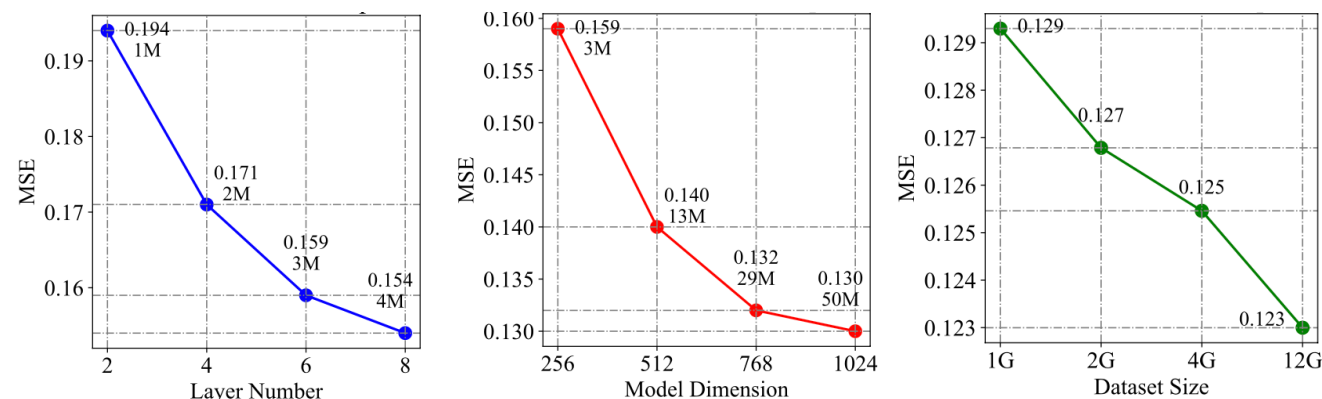
Scalability

Loss Curve of Sequence Models on UTSD



Transformer exhibits model capacity as the scalable architecture for LSTM

Scaling Model/Data Consistently Improves Performance



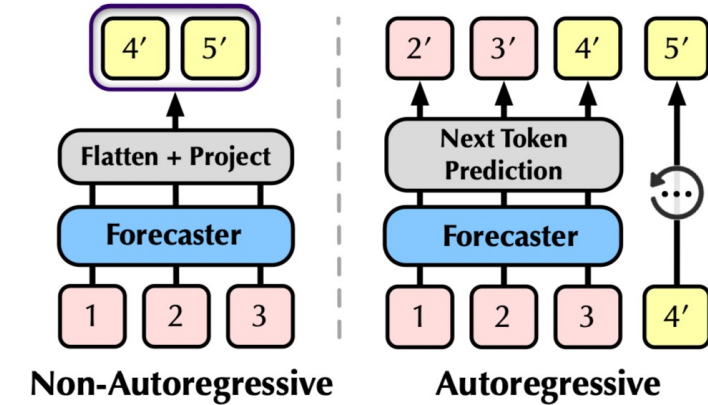
Scaling Timer achieves MSE: 0.194 \rightarrow 0.123 (-36.6%) under data scarcity, surpassing the state-of-the-art (0.129) model with full samples

Autoregressive Model

Variable Lookback Length

- Small models are constrained on fixed input/output lengths
- Similar to LLMs, Decoder is flexible on the context length

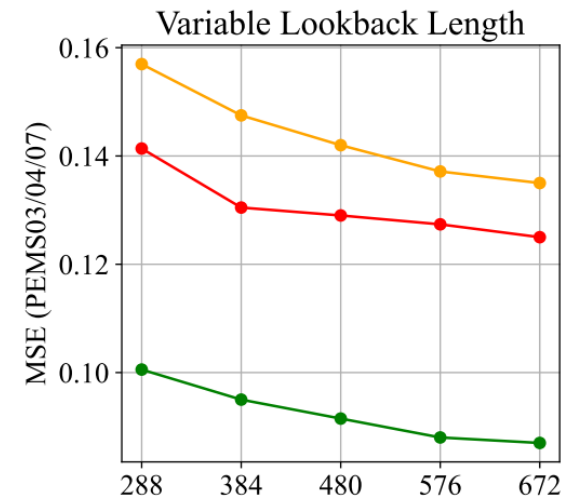
(a) Forecasting Approach



Autoregressive Model

Variable Lookback Length

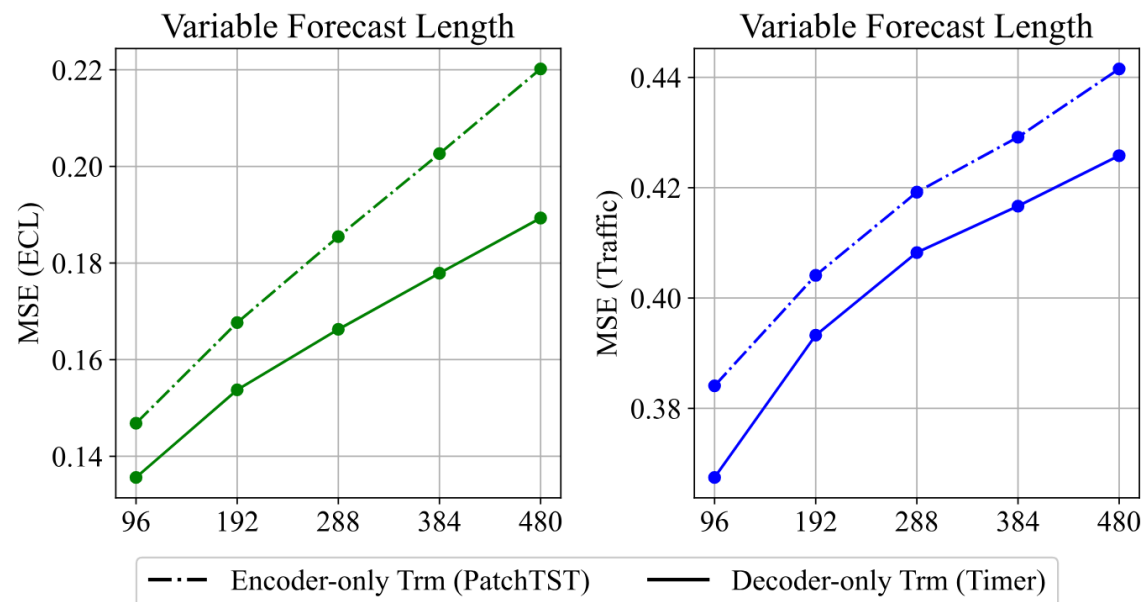
- Small models are constrained on fixed input/output lengths
- Similar to LLMs, Decoder is flexible on the context length
- Increasing the lookback window leads to stable accuracy growth



Iterative Multi-step Prediction

- Token-wise supervision
alleviates error accumulation

$$\mathcal{L}_{\text{MSE}} = \frac{1}{NS} \sum ||\mathbf{s}_i - \hat{\mathbf{s}}_i||_2^2, i = 2, \dots, N + 1.$$



Few-shot Generalization

Superwisely Trained from Scratch

Table 1. Downstream forecasting results under different data scarcity of the encoder-only and decoder-only Transformer respectively pre-trained on UTST-12G. Datasets are ordered by the oversaturation in Figure 1. Full results of PEMS and ETT can be found in Table 14.

SCENARIO	1% TARGET				5% TARGET				20% TARGET			
ARCHITECTURE	ENCODER		DECODER		ENCODER		DECODER		ENCODER		DECODER	
PRE-TRAINED	NONE	12G	NONE	12G	NONE	12G	NONE	12G	NONE	12G	NONE	12G
PEMS (AVG)	0.286	0.246	0.328	0.180	0.220	0.197	0.215	0.138	0.173	0.164	0.153	0.126
ECL	0.183	0.168	0.215	0.140	0.150	0.147	0.154	0.132	0.140	0.138	0.137	0.134
TRAFFIC	0.442	0.434	0.545	0.390	0.392	0.384	0.407	0.361	0.367	0.363	0.372	0.352
ETT (AVG)	0.367	0.317	0.340	0.295	0.339	0.303	0.321	0.285	0.309	0.301	0.297	0.288
WEATHER	0.224	0.165	0.246	0.166	0.182	0.154	0.198	0.151	0.153	0.149	0.166	0.151

Encoder will outperform when training samples are insufficient

Decoder necessitates substantial samples in end-to-end settings



Few-shot Generalization

Generalization on Downstream Tasks

Table 1. Downstream forecasting results under different data scarcity of the encoder-only and decoder-only Transformer respectively pre-trained on UTST-12G. Datasets are ordered by the oversaturation in Figure 1. Full results of PEMS and ETT can be found in Table 14.

SCENARIO	1% TARGET				5% TARGET				20% TARGET			
ARCHITECTURE	ENCODER		DECODER		ENCODER		DECODER		ENCODER		DECODER	
PRE-TRAINED	NONE	12G	NONE	12G	NONE	12G	NONE	12G	NONE	12G	NONE	12G
PEMS (AVG)	0.286	0.246	0.328	0.180	0.220	0.197	0.215	0.138	0.173	0.164	0.153	0.126
ECL	0.183	0.168	0.215	0.140	0.150	0.147	0.154	0.132	0.140	0.138	0.137	0.134
TRAFFIC	0.442	0.434	0.545	0.390	0.392	0.384	0.407	0.361	0.367	0.363	0.372	0.352
ETT (AVG)	0.367	0.317	0.340	0.295	0.339	0.303	0.321	0.285	0.309	0.301	0.297	0.288
WEATHER	0.224	0.165	0.246	0.166	0.182	0.154	0.198	0.151	0.153	0.149	0.166	0.151

In terms of Pre-training->Adaptation

Better performance can be achieved by Decoder Trm (Timer)



Evaluations of LTSMs

Quality Assessments

METHOD	TIMER (OURS)	MOIRAI (2024)	MOMENT (2024)	CHRONOS (2024)	LAG-LLAMA (2023)	TIMESFM (2023B)	TIMEGPT-1 (2023)
ARCHITECTURE	DECODER	ENCODER	ENCODER DECODER	ENCODER DECODER	DECODER	DECODER	ENCODER DECODER
MODEL SIZE	29M, 50M, 67M	14M, 91M, 311M	40M, 125M 385M	20M, 46M, 200M, 710M	200M	17M, 70M, 200M	UNKNOWN
SUPPORTED TASKS	FORECAST IMPUTATION DETECTION	FORECAST	FORECAST IMPUTATION CLASSIFICATION DETECTION	FORECAST	FORECAST	FORECAST	FORECAST DETECTION
PRE-TRAINING SCALE	28B	27.65B	1.13B	84B	0.36B	100B	100B
TOKEN TYPE	SEGMENT	SEGMENT	SEGMENT	POINT	POINT	SEGMENT	SEGMENT
CONTEXT LENGTH	≤1440	≤5000	= 512	≤512	≤1024	≤512	UNKNOWN
VARIABLE LENGTH	TRUE	TRUE	FALSE	TRUE	TRUE	TRUE	TRUE
PROBABILISTIC	FALSE	TRUE	FALSE	TRUE	TRUE	TRUE	TRUE

Future Directions

- Generalization
- Longer Context
- Probabilistic
- More Tasks
-

²<https://huggingface.co/AutonLab/MOMENT-1-large>

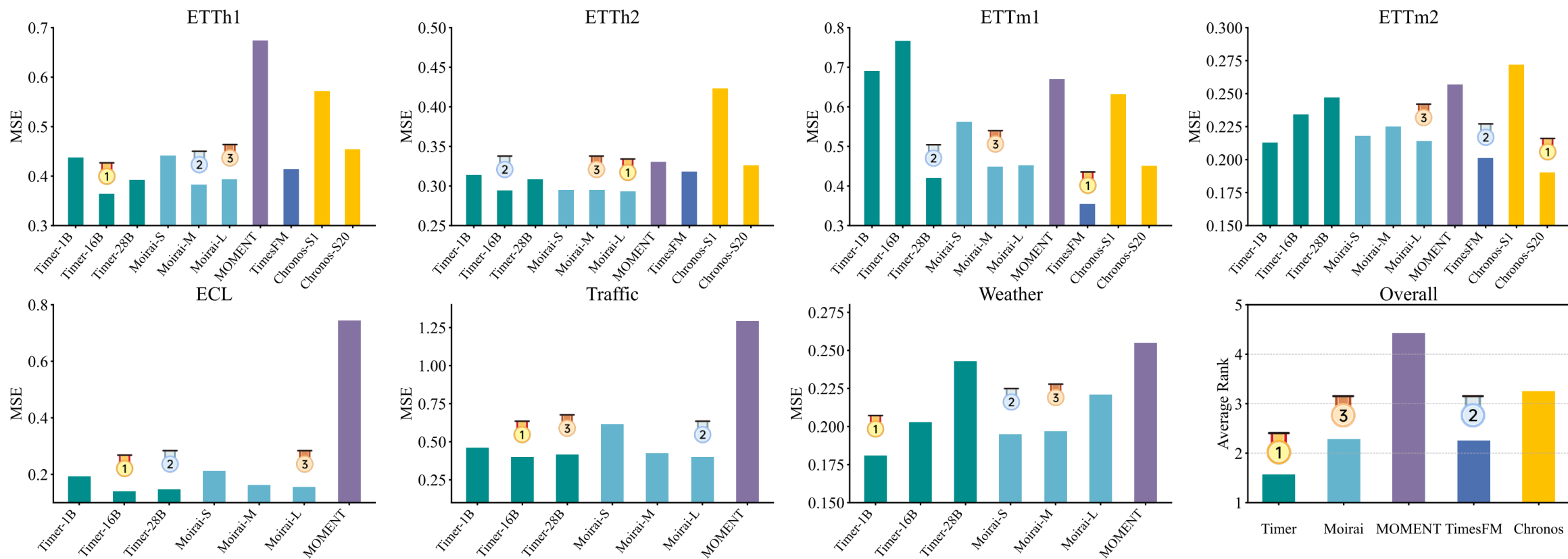
³<https://huggingface.co/amazon/chronos-t5-large>

⁴<https://huggingface.co/google/timesfm-1.0-200m>

⁵<https://huggingface.co/collections/Salesforce/moirai-10-r-models-65c8d3a94c51428c300e0742>

Benchmarks of LSTMs

Quantitative Evaluations (Zero-shot Forecasting)

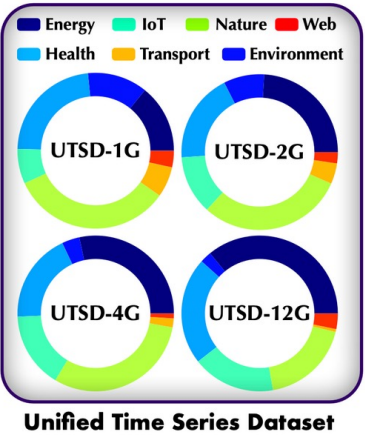


**We provided the average rank, where the lower is better,
to measure LSTMs as a *general-purpose zero-shot forecaster***

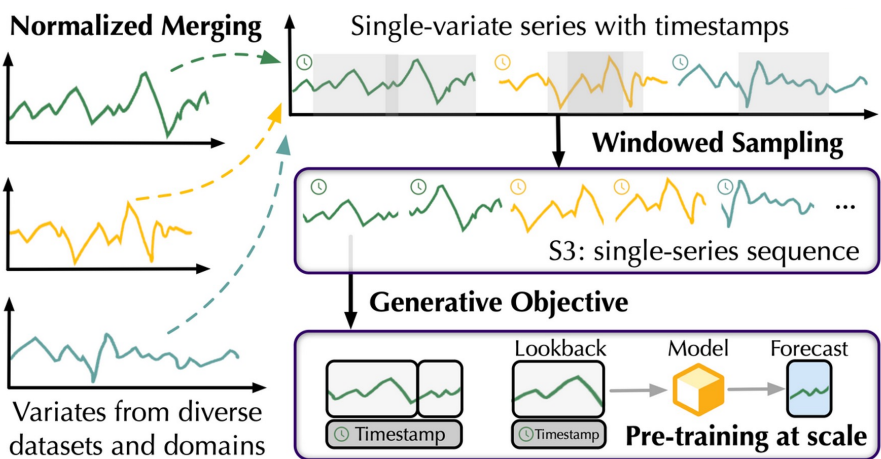
Summary



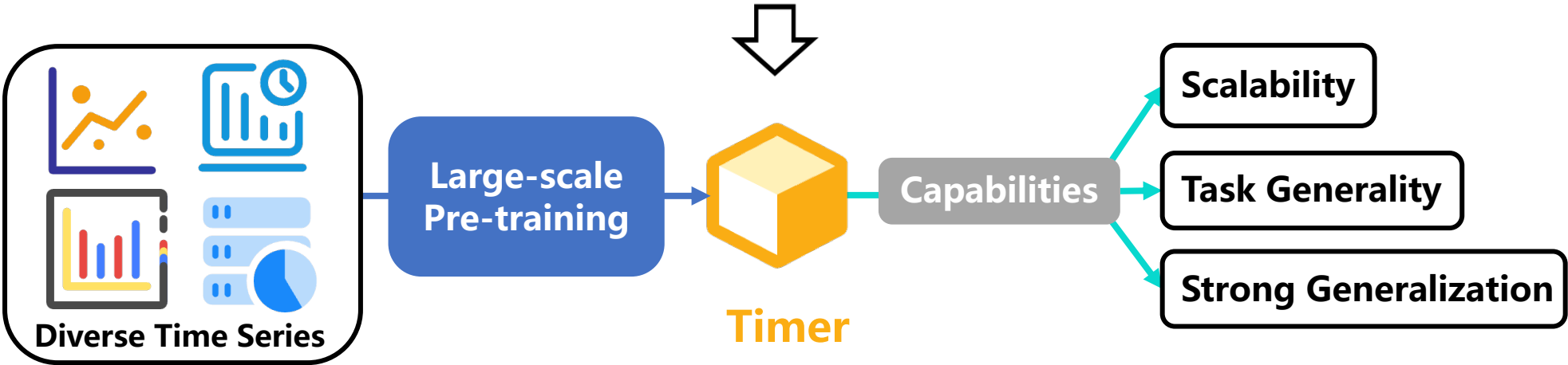
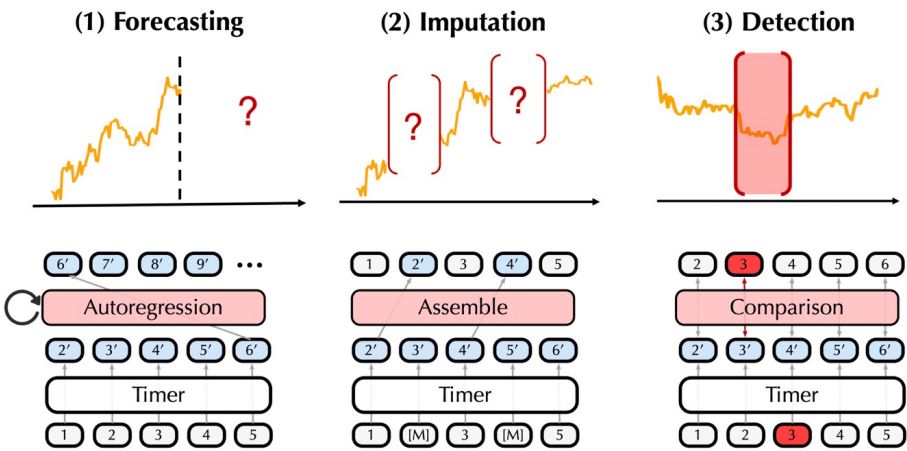
Dataset



Pre-training



Adaptation





Thank You!



GitHub: <https://github.com/thuml/Large-Time-Series-Model>

