



# Calibration of Time-Series Forecasting:

## Detecting and Adapting Context-Driven Distribution Shift

Mouxiong Chen<sup>1\*</sup>, Lefei Shen<sup>1\*</sup>, Han Fu<sup>1</sup>, Zhuo Li<sup>2</sup>, Jianling Sun<sup>1</sup>, Chenghao Liu<sup>3</sup>



<sup>1</sup>Zhejiang University, <sup>2</sup>State Street Technology (Zhejiang) Ltd, <sup>3</sup>Salesforce Research Asia, \*equal contribution

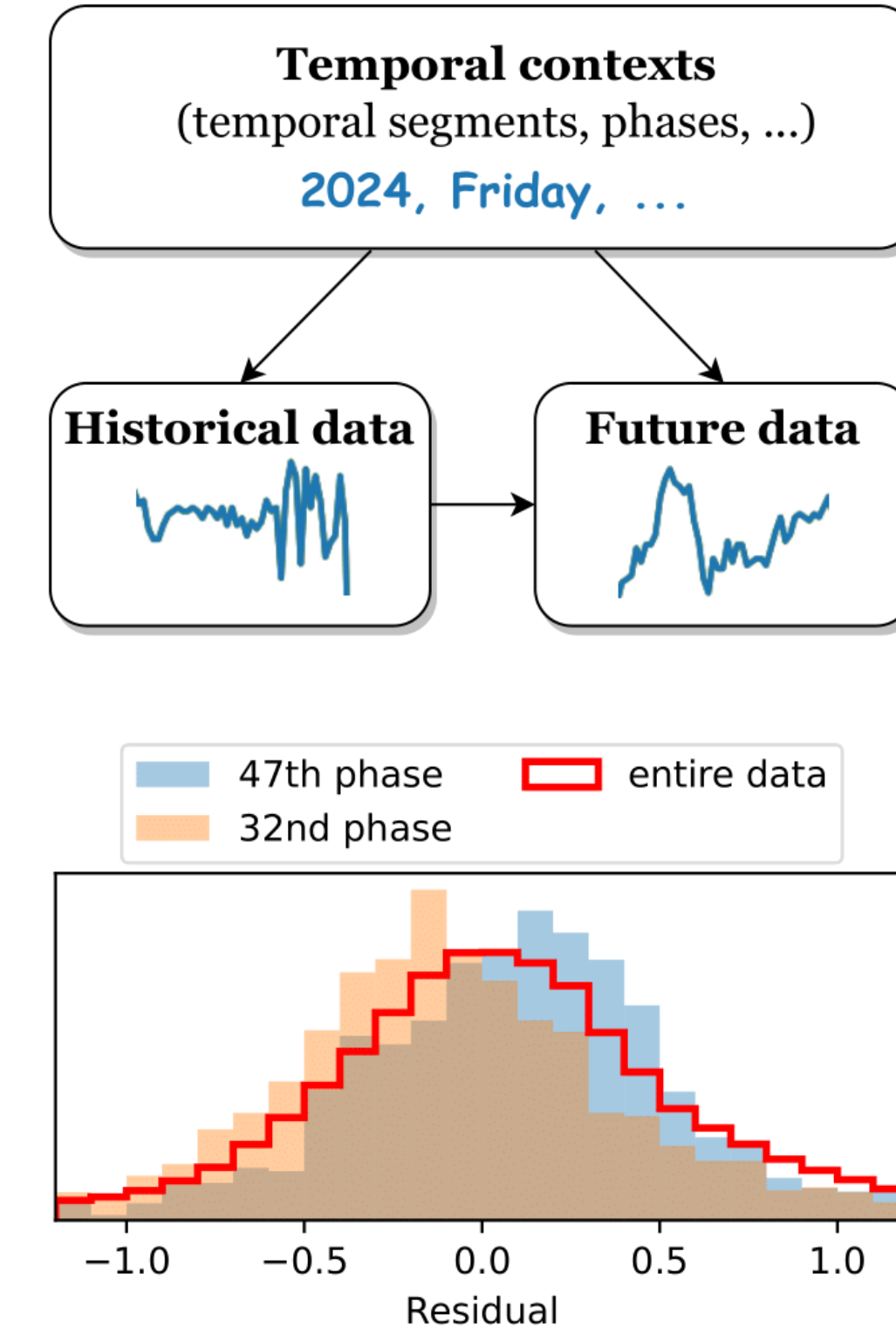
### 1. Background: CDS Problem

➤ **1.1 What is CDS?:** **Distribution shift** refers to the changing distribution and statistical properties of time series over time. Specifically, this shift is driven by some external **contexts**, such as **temporal segment** and **periodic phase**. Such phenomenon is called **Context-Driven Distribution Shift (CDS)**.

➤ **1.2 Impact of CDS?:** Contexts function as **confounders**, which simultaneously influence historical and future data. Also, the model's prediction residuals on overall data are unbiased, while those under specific contexts are biased, showing that models **struggling to achieve optimal performance across each individual contexts**.

➤ **1.3 How to solve it?:** We propose a model-agnostic "**detection and adaptation**" framework for model calibration, including:

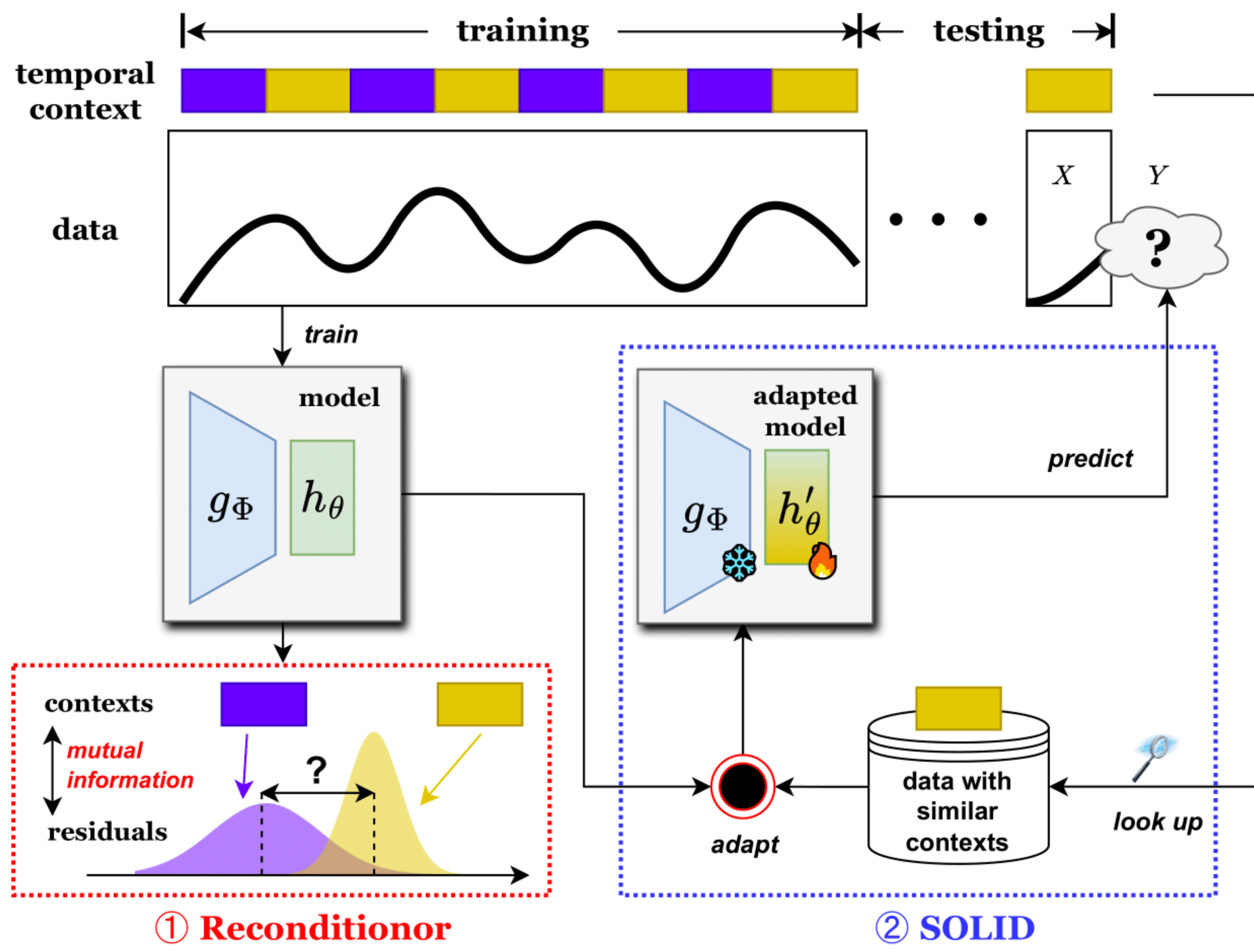
- **Reconditionor:** Detects and Quantifies the model's sensitivity to CDS.
- **SOLID:** Fine-tunes the model for each testing sample to calibrate the prediction.



◀ Causal graph in the presence of CDS.

◀ Model's prediction residuals on overall data and on different periodic phases.

### 2. Methodology



◀ Architecture of "detection and adaptation" calibration framework.

- (1) **Reconditionor:** Residual-based Context-driven Distribution Shift Detector
  - We calculate **KL divergence** between the prediction residual distributions under specific contexts and the overall residual distribution.
  - A higher value indicates a stronger impact of CDS on the model.

$$\delta = MI(\Delta Y; C) = E_C[D_{KL}(P(\Delta Y|C)||P(\Delta Y))]$$

- (2) **SOLID:** Sample-level Contextualized Adapter
  - For each test sample, we **construct a contextualized subset** with similar contexts, and **fine-tune the prediction layer** using this subset, to calibrate model's predictions.
  - Contextualized subset includes samples with **small time intervals, close periodic phases, and high sample similarity** to the test sample, corresponding to temporal segments, periodic phases, and other contexts.
  - Theoretical analysis proves that SOLID achieves a **bias-variance balance**, compared to not fine-tuning or fully retraining the prediction layer.

### 3. Algorithms

➤ **Alg1:** Calculation of Reconditionor

**Algorithm 1:** Algorithm for Reconditionor

**Input:** Model  $f$ , training data with  $K$  contexts  
 $\mathcal{D}^{\text{train}} = \{(X_{t-L:t}, X_{t:t+T}, c_t) : t < t_{\text{train}}, c_t \in [K]\}$ .

**Output:**  $\delta \in [0, 1]$  indicating  $f$ 's susceptibility to CDS.

```

1  $R \leftarrow \emptyset$ ;
2  $R_1, \dots, R_K \leftarrow \emptyset, \dots, \emptyset$ ;
3 for  $L \leq t < t_{\text{train}}$  do
4    $r \leftarrow f(X_{t-L:t}, X_{t:t+T}) - X_{t:t+T}$ ;
5    $R \leftarrow R \cup r$ ;
6    $R_{c_t} \leftarrow R_{c_t} \cup r$ ;
7 end
8  $\mu, \sigma \leftarrow \text{Mean}(R), \text{Standard-Deviation}(R)$ ;
9  $\delta \leftarrow 0$ ;
10 for  $c \in [K]$  do
11    $\mu_c, \sigma_c \leftarrow \text{Mean}(R_c), \text{Standard-Deviation}(R_c)$ ;
12    $\delta \leftarrow \delta + \frac{|R_c|}{|R|} \text{KL}(\mathcal{N}(\mu_c, \sigma_c^2) \parallel \mathcal{N}(\mu, \sigma^2))$ ;
13 end
14 return  $\delta$ ;

```

➤ **Alg2:** Algorithm for SOLID

**Algorithm 2:** Algorithm for SOLID

**Input:** Model  $f = (g_\Phi, h_\Theta)$ , test sample  $X_{t-L:t}$ , preceding data  $\{(X_{t'-L:t'}, X_{t':t'+T}) : t' + T \leq t\}$ , similarity metric  $S(\cdot, \cdot)$ , periodic length  $T^*$  computed by Eq.(3), hyperparameters  $\lambda_T, \lambda_P, \lambda_N$  and  $l_r$ .

**Output:** Prediction for the test sample:  $\hat{X}_{t:t+T}$

```

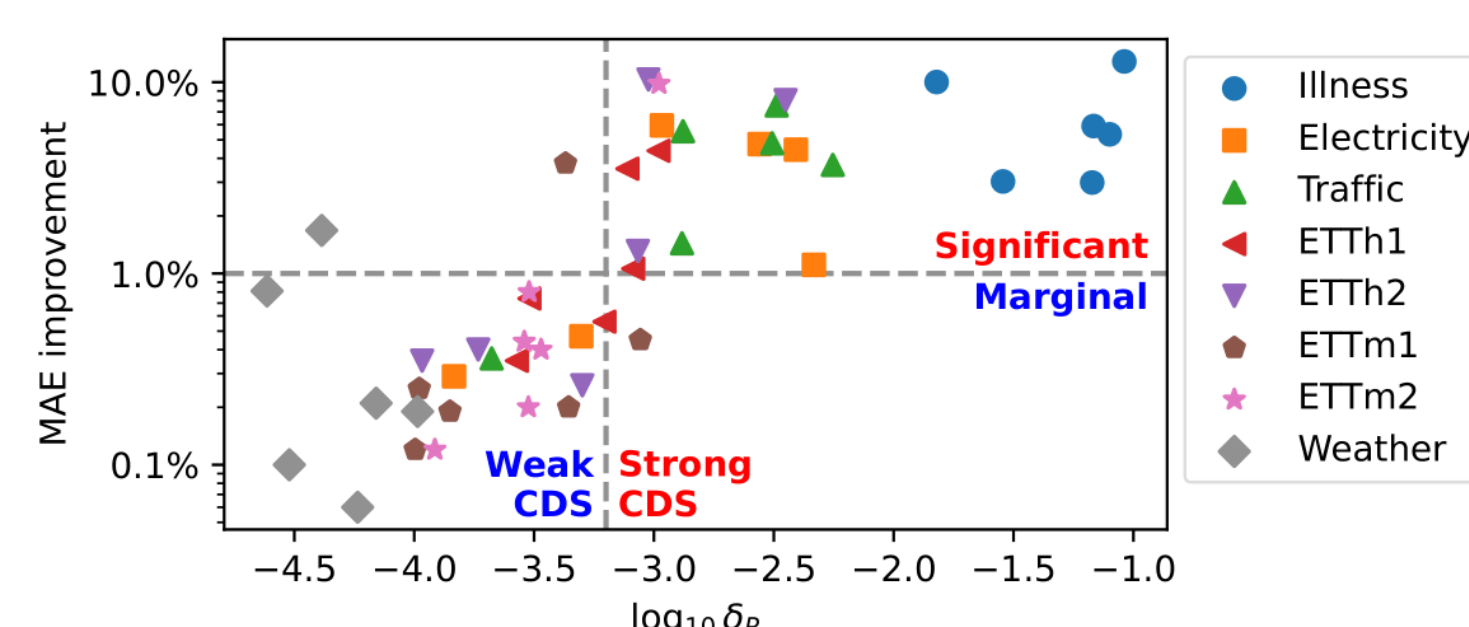
1  $\mathcal{T} \leftarrow \emptyset$ ;
2 for  $t - \lambda_T \leq t' \leq t - T$  do
3    $\Delta_P \leftarrow \left| \frac{t \bmod T^* - t' \bmod T^*}{T^*} \right|$ ;
4   if  $\Delta_P < \lambda_P$  then
5      $\mathcal{T} \leftarrow \mathcal{T} \cup \{t'\}$ ;
6   end
7 end
8  $\mathcal{T}_{\text{ctx}} \leftarrow \arg\text{Top-}\lambda_N(S(X_{t'-L:t'}, X_{t-L:t}));$ 
9  $\mathcal{D}_{\text{ctx}} \leftarrow \{(X_{t'-L:t'}, X_{t':t'+T}) \mid t' \in \mathcal{T}_{\text{ctx}}\}$ ;
10  $h'_\Theta \leftarrow \text{fine-tune } h_\Theta \text{ using } \mathcal{D}_{\text{ctx}} \text{ with a learning rate } l_r$ ;
11  $\hat{X}_{t:t+T} \leftarrow h'_\Theta(g_\Phi(X_{t-L:t}))$ ;
12 return  $\hat{X}_{t:t+T}$ ;

```

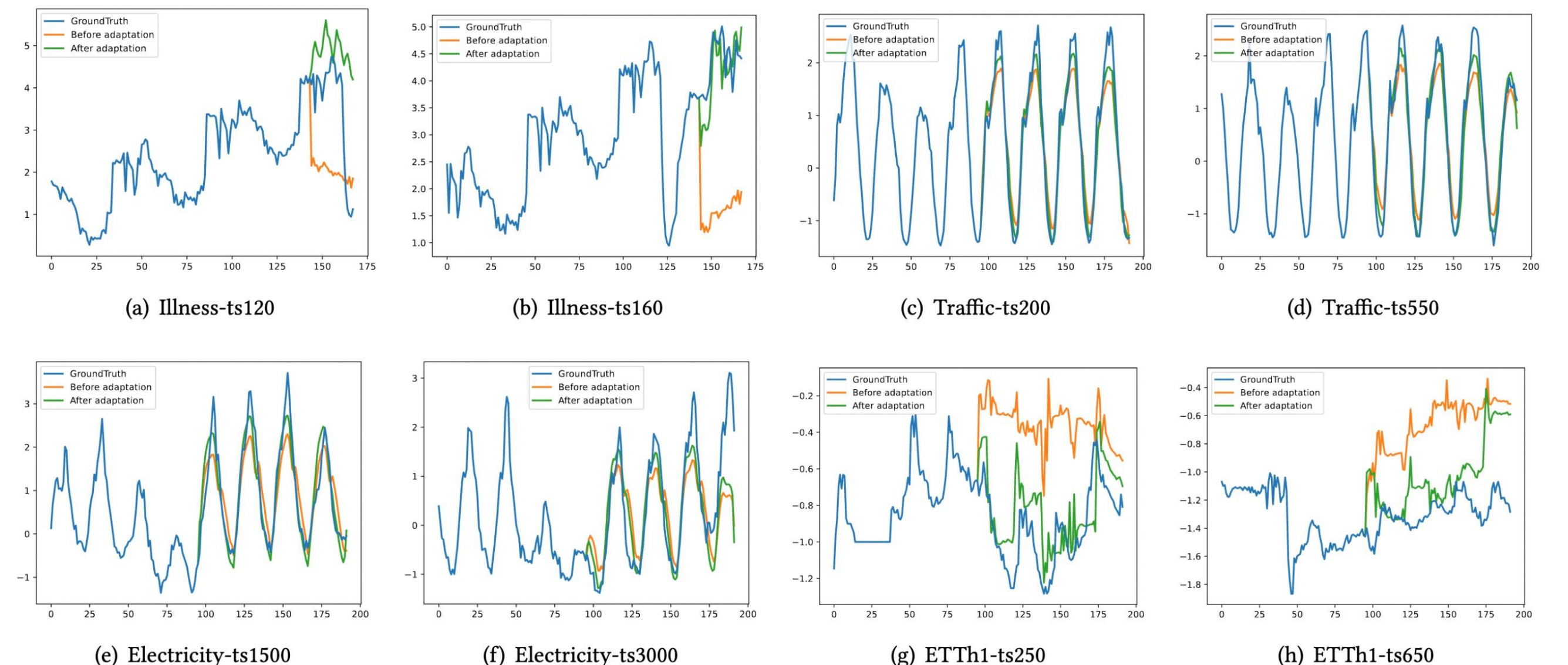
### 4. Experiments

▼ MSE & MAE are averaged from prediction length of 24/36/48/60 for Illness and 96/192/336/720 for others. "↑": average improvements by SOLID. "δ": Reconditioner value periodic phases ( $\delta_P$ ) and temporal segments ( $\delta_T$ ), reported in the form of " $\log_{10} \delta_P$  &  $\log_{10} \delta_T$ ". **RED** denotes a strong CDS in periodic phases ( $\log_{10} \delta_P \geq -3.2$ ), while **BLUE** denotes a weak CDS.

Dataset	Illness				Electricity				Traffic				ETTh1				ETTh2				
Method	/ +SOLID				/ +SOLID				/ +SOLID				/ +SOLID				/ +SOLID				
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
Informer	avg	5.140	1.554	3.206	1.239	0.352	0.442	0.277	0.379	0.816	0.457	0.725	0.449	1.036	0.797	0.792	0.644	4.648	1.770	2.272	1.181
	δ	-1.096 & -1.148				-2.975 & -2.593				-2.762 & -2.238				-3.83 & -1.883				-3.021 & -1.513			
Autoformer	avg	2.887	1.131	2.605	1.064	0.239	0.347	0.225	0.332	0.655	0.408	0.599	0.389	0.485	0.480	0.480	0.478	0.442	0.454	0.441	0.453
	δ	-1.166 & -1.016				-2.408 & -2.134				-2.507 & -2.304				-3.572 & -2.321				-3.967 & -1.921			
FEDformer	avg	2.787	1.124	2.546	1.064	0.210	0.324	0.194	0.304	0.608	0.374	0.549	0.360	0.436	0.456	0.431	0.452	0.439	0.451	0.438	0.450
	δ	-1.099 & -0.917				-2.967 & -2.545				-2.254 & -2.265				-3.52 & -2.463				-3.733 & -1.943			
ETSformer	avg	2.471	0.993	2.310	0.963	0.211	0.325	0.197	0.309	0.615	0.390	0.502	0.360	0.547	0.510	0.541	0.507	0.437	0.455	0.429	0.449
	δ	-1.544 & -1.001				-2.559 & -2.724				-2.488 & -2.201				-3.207 & -3.019				-3.067 & -1.079			
Crossformer	avg	3.443	1.231	2.621	1.073	0.229	0.334	0.226	0.330	0.537	0.302	0.482	0.285	0.443	0.462	0.417	0.446	1.152	0.778	0.829	0.716
	δ	-1.038 & -0.917				-2.333 & -2.42				-2.879 & -2.171				-3.111 & -2.767				-2.451 & -1.149			
DLinear	avg	2.192	1.046	1.842	0.941	0.167	0.264	0.166	0.263	0.434	0.295	0.429	0.291	0.467	0.468	0.441	0.447	0.448	0.453	0.364	0.406
	δ	-1.821 & -1.301				-3.303 & -2.645				-2.883 & -2.589				-2.981 & -2.321				-3.023 & -1.891			
PatchTST	avg	1.542	0.827	1.497	0.802	0.164	0.256	0.162	0.255	0.205	0.274	0.204	0.273	0.414	0.423	0.408	0.419	0.331	0.380	0.330	0.379
	δ	-1.172 & -1.054				-3.834 & -2.878				-3.677 & -2.901				-3.087 & -2.529				-3.299 & -1.851			



▲  $\delta_P$  and MAE improvement shows positive correlation. This indicates that models detected by **Reconditionor** as more sensitive to CDS exhibit greater performance improvements after adaptation by **SOLID**.



▲ **Blue** lines are ground-truth, **orange** lines are predictions before SOLID, and **green** lines are predictions after SOLID.