

---

# Supplementary material: *Adaptive Normalization for Non-stationary Time Series Forecasting: A Temporal Slice Perspective*

---

Anonymous Author(s)

Affiliation

Address

email

## 1 Effects of SAN on Non-stationary Time Series Forecasting

### 1.1 Discussions

As illustrated in the main paper, our proposed SAN is a compact plug-and-play framework. We will first give a brief discussion in this section on how SAN can be effective.

It is of utmost importance that SAN can well alleviate the impact of the non-stationary nature of time series data. Forecasting models may encounter a non-i.i.d problem with non-stationary data, that is, the marginal distribution of each input instance can be different, which may lead to a huge difference between the distribution of the training set and the test set. Thus the models can not generalize well in future predictions. However, SAN will normalize all the input instances into a standard normal distribution and force the mean and variance of the training and test data distributions to be identical. In this way, all the data instances are from the same distribution, therefore the forecasting task is simplified as the models can get rid of the noises caused by non-stationary factors and only focus on mining the time-invariant patterns. Moreover, compared to existing normalization methods for forecasting, our modeling of the non-stationary property in a time slice view is more in-depth and realistic, so SAN can better remove the non-stationary factors in input sequences while keeping their instinct information in the normalization phase. Hence, SAN is theoretically expected to perform better in non-stationary time series forecasting.

Another part that contributes to the effectiveness of SAN is the statistics prediction module and the two-stage training schema. With the statistics prediction module independently modeling the evolving trends of statistical properties, SAN adopts more precise statistics for adaptive denormalization than existing solutions. Moreover, the proposed two-stage strategy actually simplifies the original forecasting task by divide and conquer: In the first stage we try to learn the general direction and dispersion of the future data, which is easy to fit and is conducted by the light statistics prediction module. Next, we utilize the powerful backbone model to discover the scale-free periodic-like features to estimate future values under the guidance of the well-trained statistics prediction module. Therefore, backbone models in SAN are actually responsible for an easier subtask. Considering that SAN can usually give reliable estimations on future distributions, SAN is expected to perform well on non-stationary time series forecasting by splitting the task into two simpler subtasks.

### 1.2 Theoretical Analysis

Using the same notation in the paper, we prove that all the inputs after SAN's normalization follow a standard normal distribution, validating SAN's capability to remove the non-stationary factors theoretically.

In detail, for arbitrary input sequence  $\mathbf{x}^i$ , SAN first split it into  $M$  non-overlapping slices  $\{\mathbf{x}_j^i\}_{j=1}^M$  and normalizes them according to their statistics. Therefore we will get:

$$\forall i, j \mathbb{E}[\bar{\mathbf{x}}_j^i] = 0, \text{Var}[\bar{\mathbf{x}}_j^i] = I \quad (1)$$

And as for the statistics of normalized input  $\bar{\mathbf{x}}^i$ , it satisfies the following equations:

$$\begin{aligned} \mathbb{E}[\bar{\mathbf{x}}^i] &= \mathbb{E}_j[\mathbb{E}[\bar{\mathbf{x}}_j^i]] \\ &= \mathbb{E}_j[0] \\ &= 0 \end{aligned} \quad (2)$$

$$\begin{aligned} \text{Var}[\bar{\mathbf{x}}^i] &= \frac{\sum_{t=0}^{L_{in}} (\bar{\mathbf{x}}_{:,t}^i - \mathbb{E}[\bar{\mathbf{x}}^i])^2}{L_{in}} \\ &= \frac{\sum_{t=0}^{L_{in}} (\bar{\mathbf{x}}_{:,t}^i)^2}{MT} \\ &= \frac{1}{M} * \left( \frac{\sum_{t=0}^T (\bar{\mathbf{x}}_{:,t}^i)^2}{T} + \frac{\sum_{t=T}^{2T} (\bar{\mathbf{x}}_{:,t}^i)^2}{T} + \dots + \frac{\sum_{t=(M-1)T}^{MT} (\bar{\mathbf{x}}_{:,t}^i)^2}{T} \right) \\ &= \mathbb{E}_j[\text{Var}[\bar{\mathbf{x}}_j^i]] \\ &= I \end{aligned} \quad (3)$$

Here  $\bar{\mathbf{x}}_{:,t}^i \in R^{V*1}$  denotes all the normalized variables in time step  $t$ . From the above equations, we can learn that any input sequence follows a standard normal distribution after the normalization operation of SAN, which meets our expectations.

## 2 Supplementary Experiments

### 2.1 Full Benchmark on the ETT Dataset

Table 1: Multivariate forecasting results on full ETT dataset.

Methods	Metric	DLinear		+ SAN		FEDformer		+ SAN		Autoformer		+ SAN		SCINet		+ SAN	
		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1	96	<b>0.377</b>	<b>0.399</b>	0.383	<b>0.399</b>	<b>0.371</b>	0.411	0.383	<b>0.409</b>	<b>0.458</b>	<b>0.448</b>	0.488	0.464	0.470	0.479	<b>0.391</b>	<b>0.405</b>
	192	<b>0.417</b>	0.426	0.419	<b>0.419</b>	<b>0.420</b>	0.443	0.431	<b>0.438</b>	<b>0.481</b>	0.474	0.498	<b>0.472</b>	0.541	0.520	<b>0.438</b>	<b>0.433</b>
	336	0.464	0.461	<b>0.437</b>	<b>0.432</b>	<b>0.446</b>	0.459	0.471	<b>0.456</b>	<b>0.508</b>	<b>0.485</b>	0.530	0.498	0.643	0.587	<b>0.477</b>	<b>0.451</b>
	720	0.493	0.505	<b>0.446</b>	<b>0.459</b>	<b>0.482</b>	0.495	0.504	<b>0.488</b>	<b>0.525</b>	0.516	0.555	<b>0.514</b>	0.774	0.669	<b>0.489</b>	<b>0.474</b>
ETTh2	96	0.292	0.356	<b>0.277</b>	<b>0.338</b>	0.341	0.382	<b>0.300</b>	<b>0.355</b>	0.384	0.420	<b>0.316</b>	<b>0.366</b>	0.690	0.625	<b>0.294</b>	<b>0.347</b>
	192	0.383	0.418	<b>0.340</b>	<b>0.378</b>	0.426	0.436	<b>0.392</b>	<b>0.413</b>	0.457	0.454	<b>0.413</b>	<b>0.426</b>	0.991	0.742	<b>0.374</b>	<b>0.398</b>
	336	0.473	0.477	<b>0.356</b>	<b>0.398</b>	0.481	0.479	<b>0.459</b>	<b>0.462</b>	0.468	0.473	<b>0.446</b>	<b>0.457</b>	1.028	0.759	<b>0.412</b>	<b>0.430</b>
	720	0.708	0.599	<b>0.396</b>	<b>0.435</b>	<b>0.458</b>	0.477	0.462	<b>0.472</b>	0.473	0.485	<b>0.471</b>	<b>0.474</b>	1.363	0.885	<b>0.437</b>	<b>0.461</b>
ETTM1	96	0.301	0.344	<b>0.288</b>	<b>0.342</b>	0.362	0.408	<b>0.311</b>	<b>0.355</b>	0.493	0.470	<b>0.343</b>	<b>0.378</b>	0.444	0.464	<b>0.321</b>	<b>0.360</b>
	192	0.335	0.366	<b>0.323</b>	<b>0.363</b>	0.395	0.427	<b>0.351</b>	<b>0.383</b>	0.546	0.498	<b>0.390</b>	<b>0.400</b>	0.491	0.500	<b>0.347</b>	<b>0.380</b>
	336	0.370	0.387	<b>0.357</b>	<b>0.384</b>	0.441	0.454	<b>0.390</b>	<b>0.407</b>	0.658	0.543	<b>0.415</b>	<b>0.418</b>	0.572	0.556	<b>0.385</b>	<b>0.403</b>
	720	0.425	0.421	<b>0.409</b>	<b>0.415</b>	0.488	0.481	<b>0.456</b>	<b>0.444</b>	0.626	0.532	<b>0.476</b>	<b>0.453</b>	0.728	0.654	<b>0.450</b>	<b>0.441</b>
ETTM2	96	0.169	0.263	<b>0.166</b>	<b>0.258</b>	0.191	0.283	<b>0.175</b>	<b>0.266</b>	0.261	0.329	<b>0.236</b>	<b>0.317</b>	0.303	0.404	<b>0.176</b>	<b>0.267</b>
	192	0.232	0.310	<b>0.223</b>	<b>0.302</b>	0.261	0.326	<b>0.246</b>	<b>0.315</b>	0.282	0.339	<b>0.260</b>	<b>0.329</b>	0.568	0.569	<b>0.240</b>	<b>0.311</b>
	336	0.303	0.361	<b>0.272</b>	<b>0.330</b>	0.327	0.365	<b>0.315</b>	<b>0.362</b>	0.350	0.378	<b>0.330</b>	<b>0.376</b>	0.793	0.689	<b>0.300</b>	<b>0.351</b>
	720	0.403	0.424	<b>0.360</b>	<b>0.384</b>	0.428	0.423	<b>0.412</b>	<b>0.422</b>	0.438	<b>0.428</b>	<b>0.417</b>	<b>0.428</b>	1.200	0.851	<b>0.391</b>	<b>0.405</b>

We provide the full multivariate forecasting results on the ETT dataset in Table 1, which includes the hourly datasets ETTh1&ETTh2 and the 15-minutes datasets ETTm1&ETTM2. It is obvious that SAN also achieves significant improvements on these datasets on various backbone models.

### 2.2 Univariate Forecasting Results

Following the same settings of our main experiment, we provide the univariate forecasting results in Table 2. Similar to the results of multivariate forecasting, SAN can boost the performance of mainstream forecasting models in most cases. On average of all the benchmark settings, DLinear enhanced by SAN reduces MSE by **6.04%** (from 0.230 to 0.214). The improvements for FEDformer, Autoformer and SCINet are **15.40%**, **29.27%** and **36.29%** respectively.

Table 2: Univariate forecasting results. The bold values indicate better performance.

Methods Metric		DLinear		+ SAN		FEDformer		+ SAN		Autoformer		+ SAN		SCINet		+ SAN	
		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Electricity	96	<b>0.203</b>	<b>0.315</b>	0.204	0.317	0.302	0.413	<b>0.248</b>	<b>0.363</b>	0.442	0.490	<b>0.283</b>	<b>0.386</b>	0.364	0.435	<b>0.321</b>	<b>0.412</b>
	192	<b>0.233</b>	<b>0.336</b>	0.238	0.341	0.377	0.459	<b>0.278</b>	<b>0.379</b>	0.555	0.550	<b>0.296</b>	<b>0.393</b>	0.345	0.419	<b>0.328</b>	<b>0.412</b>
	336	<b>0.268</b>	<b>0.363</b>	0.278	0.371	0.673	0.636	<b>0.324</b>	<b>0.411</b>	0.617	0.620	<b>0.359</b>	<b>0.440</b>	0.368	<b>0.435</b>	<b>0.363</b>	0.436
	720	0.330	0.425	<b>0.325</b>	<b>0.420</b>	0.575	0.575	<b>0.502</b>	<b>0.514</b>	0.645	0.624	<b>0.443</b>	<b>0.503</b>	0.420	0.478	<b>0.410</b>	<b>0.477</b>
Exchange	96	<b>0.108</b>	<b>0.254</b>	0.138	0.288	0.134	0.272	<b>0.113</b>	<b>0.252</b>	0.155	0.305	<b>0.097</b>	<b>0.233</b>	0.167	0.332	<b>0.090</b>	<b>0.226</b>
	192	<b>0.193</b>	<b>0.350</b>	0.287	0.436	<b>0.290</b>	0.418	0.307	<b>0.404</b>	0.405	0.495	<b>0.208</b>	<b>0.358</b>	0.486	0.552	<b>0.185</b>	<b>0.335</b>
	336	0.428	<b>0.511</b>	<b>0.416</b>	0.523	0.490	0.542	<b>0.431</b>	<b>0.501</b>	0.874	0.728	<b>0.401</b>	<b>0.495</b>	0.579	0.608	<b>0.396</b>	<b>0.484</b>
	720	1.137	0.848	<b>0.859</b>	<b>0.719</b>	1.302	0.883	<b>1.188</b>	<b>0.835</b>	1.193	0.845	<b>1.071</b>	<b>0.787</b>	<b>0.853</b>	<b>0.740</b>	1.106	0.797
Traffic	96	0.124	<b>0.197</b>	<b>0.123</b>	0.199	0.179	0.282	<b>0.144</b>	<b>0.236</b>	0.265	0.375	<b>0.172</b>	<b>0.273</b>	0.352	0.430	<b>0.267</b>	<b>0.364</b>
	192	0.125	<b>0.200</b>	<b>0.124</b>	<b>0.200</b>	0.211	0.316	<b>0.141</b>	<b>0.232</b>	0.266	0.372	<b>0.211</b>	<b>0.316</b>	0.291	0.377	<b>0.240</b>	<b>0.338</b>
	336	<b>0.126</b>	<b>0.206</b>	0.228	0.269	0.369	0.458	<b>0.207</b>	<b>0.318</b>	0.284	0.371	<b>0.164</b>	<b>0.259</b>	<b>0.298</b>	<b>0.387</b>	0.347	0.396
	720	0.141	0.226	<b>0.138</b>	<b>0.223</b>	<b>0.300</b>	<b>0.407</b>	0.477	0.526	0.260	0.369	<b>0.179</b>	<b>0.286</b>	0.339	0.417	<b>0.311</b>	<b>0.384</b>
Weather	96	0.004	0.047	<b>0.002</b>	<b>0.032</b>	<b>0.002</b>	<b>0.037</b>	0.003	0.042	0.004	0.047	<b>0.002</b>	<b>0.038</b>	0.005	0.060	<b>0.003</b>	<b>0.039</b>
	192	0.005	0.057	<b>0.002</b>	<b>0.037</b>	0.005	0.058	<b>0.004</b>	<b>0.049</b>	<b>0.003</b>	0.045	<b>0.003</b>	<b>0.047</b>	0.006	0.065	<b>0.002</b>	<b>0.036</b>
	336	0.006	0.068	<b>0.003</b>	<b>0.047</b>	<b>0.003</b>	<b>0.045</b>	0.004	0.052	0.008	0.068	<b>0.003</b>	<b>0.046</b>	0.007	0.068	<b>0.004</b>	<b>0.049</b>
	720	0.007	0.070	<b>0.004</b>	<b>0.050</b>	0.011	0.080	<b>0.004</b>	<b>0.048</b>	0.058	0.176	<b>0.004</b>	<b>0.049</b>	0.007	0.070	<b>0.003</b>	<b>0.045</b>
ILI	24	0.741	0.681	<b>0.663</b>	<b>0.626</b>	0.910	0.825	<b>0.798</b>	<b>0.688</b>	0.865	0.800	<b>0.765</b>	<b>0.721</b>	6.336	2.130	<b>0.707</b>	<b>0.665</b>
	36	0.570	0.634	<b>0.552</b>	<b>0.599</b>	0.873	0.823	<b>0.697</b>	<b>0.691</b>	0.984	0.855	<b>0.660</b>	<b>0.693</b>	6.159	1.998	<b>0.743</b>	<b>0.706</b>
	48	0.740	0.742	<b>0.647</b>	<b>0.669</b>	1.027	0.904	<b>0.820</b>	<b>0.761</b>	1.105	0.925	<b>0.753</b>	<b>0.752</b>	6.597	2.082	<b>0.783</b>	<b>0.744</b>
	60	0.911	0.848	<b>0.765</b>	<b>0.743</b>	1.221	1.002	<b>0.981</b>	<b>0.839</b>	1.222	0.982	<b>1.024</b>	<b>0.904</b>	7.556	2.418	<b>0.902</b>	<b>0.801</b>
ETTh1	96	0.058	<b>0.180</b>	<b>0.056</b>	0.181	0.097	0.241	<b>0.067</b>	<b>0.195</b>	0.093	0.241	<b>0.062</b>	<b>0.188</b>	0.110	0.262	<b>0.057</b>	<b>0.180</b>
	192	0.078	0.216	<b>0.076</b>	<b>0.212</b>	0.109	0.257	<b>0.081</b>	<b>0.215</b>	0.121	0.290	<b>0.082</b>	<b>0.216</b>	0.152	0.312	<b>0.075</b>	<b>0.209</b>
	336	0.099	0.246	<b>0.092</b>	<b>0.240</b>	0.103	0.251	<b>0.098</b>	<b>0.240</b>	0.115	0.271	<b>0.089</b>	<b>0.232</b>	0.183	0.350	<b>0.093</b>	<b>0.238</b>
	720	0.158	0.322	<b>0.092</b>	<b>0.240</b>	0.130	0.290	<b>0.103</b>	<b>0.248</b>	0.108	0.259	<b>0.106</b>	<b>0.249</b>	0.252	0.432	<b>0.096</b>	<b>0.245</b>
ETTh2	96	<b>0.132</b>	<b>0.280</b>	0.133	0.281	0.145	0.301	<b>0.141</b>	<b>0.286</b>	0.181	0.332	<b>0.141</b>	<b>0.288</b>	0.149	0.306	<b>0.129</b>	<b>0.274</b>
	192	0.177	0.330	<b>0.174</b>	<b>0.327</b>	0.188	0.339	<b>0.184</b>	<b>0.331</b>	0.213	0.371	<b>0.196</b>	<b>0.350</b>	0.187	0.340	<b>0.178</b>	<b>0.326</b>
	336	0.207	0.366	<b>0.200</b>	<b>0.359</b>	<b>0.220</b>	0.380	0.224	<b>0.371</b>	0.232	0.391	<b>0.221</b>	<b>0.370</b>	0.236	0.385	<b>0.222</b>	<b>0.374</b>
	720	0.301	0.447	<b>0.237</b>	<b>0.391</b>	0.279	0.427	<b>0.257</b>	<b>0.407</b>	0.267	0.417	<b>0.289</b>	<b>0.431</b>	0.326	0.468	<b>0.272</b>	<b>0.421</b>
ETTm1	96	0.027	<b>0.123</b>	<b>0.026</b>	<b>0.123</b>	0.060	0.193	<b>0.028</b>	<b>0.125</b>	0.059	0.193	<b>0.027</b>	<b>0.125</b>	0.065	0.204	<b>0.032</b>	<b>0.135</b>
	192	0.045	0.156	<b>0.040</b>	<b>0.151</b>	0.065	0.202	<b>0.044</b>	<b>0.159</b>	0.083	0.231	<b>0.042</b>	<b>0.155</b>	0.198	0.342	<b>0.049</b>	<b>0.168</b>
	336	0.059	0.178	<b>0.055</b>	<b>0.176</b>	0.066	0.199	<b>0.059</b>	<b>0.189</b>	0.069	0.205	<b>0.057</b>	<b>0.181</b>	0.221	0.382	<b>0.068</b>	<b>0.199</b>
	720	0.081	0.212	<b>0.077</b>	<b>0.208</b>	<b>0.084</b>	<b>0.230</b>	0.098	0.234	0.095	0.243	<b>0.081</b>	<b>0.213</b>	0.303	0.466	<b>0.093</b>	<b>0.231</b>
ETTm2	96	<b>0.063</b>	<b>0.183</b>	<b>0.063</b>	0.186	0.097	0.244	<b>0.060</b>	<b>0.183</b>	0.128	0.278	<b>0.068</b>	<b>0.195</b>	0.073	0.200	<b>0.069</b>	<b>0.193</b>
	192	<b>0.093</b>	<b>0.229</b>	<b>0.093</b>	0.230	0.129	0.281	<b>0.093</b>	<b>0.233</b>	0.145	0.298	<b>0.099</b>	<b>0.240</b>	0.107	0.248	<b>0.103</b>	<b>0.240</b>
	336	0.120	<b>0.263</b>	<b>0.119</b>	0.264	0.174	0.326	<b>0.129</b>	<b>0.276</b>	0.148	0.303	<b>0.123</b>	<b>0.269</b>	0.163	0.314	<b>0.135</b>	<b>0.281</b>
	720	0.173	<b>0.318</b>	<b>0.171</b>	0.319	0.201	0.354	<b>0.193</b>	<b>0.337</b>	0.208	0.359	<b>0.174</b>	<b>0.320</b>	0.325	0.441	<b>0.191</b>	<b>0.337</b>

## 2.3 Validation on Various Input Lengths

The input length plays an essential role in time series forecasting tasks as it determines how much historical temporal information the model can mine. One may hope that for powerful deep models, the longer the input length, the better the forecasting results. However, a recent study on this question reveals that deep Transformer-based models are not capable of capturing temporal dependencies in the long-term input sequences [6]. That is, the performance of these deep models stays stable or even degrades when the input length increases.

Apart from the design of these deep models, we hold that such a phenomenon can be raised by the non-stationary property of time series. As the input length increases, the variance among input sequences grows larger and ultimately makes it harder for deep models to discover the time-invariant patterns. Therefore, by removing the non-stationary factors in the input by SAN, deep models are expected to exhibit a steady decline in metrics with longer input lengths.

To evidence our thoughts, we conduct long-term forecasting experiments, i.e.,  $L_{out} = 720$ , with various input lengths  $L_{in} \in \{24, 48, 72, 96, 120, 144, 168, 192, 336, 504, 672, 720\}$  on the Transformer-based models. Here we choose Transformer [4], Informer [7], Autoformer [5] and FEDformer [8] as the backbone models. The MSE evaluations are plotted in Fig. 1. Note that we omit large values in the line chart to better demonstrate the trend of the overall results. From the figure, we can see that with the assistance of SAN, the performance of deep models with long sequence input is largely improved. When the input length is set to 720 on the Electricity dataset, the performance of Informer has been boosted by **77.83%** (from 0.9426 to 0.2090), and the average improvement on four backbones under the same setting is **52.55%**. Moreover, all of the backbones enhanced by SAN tend to produce more

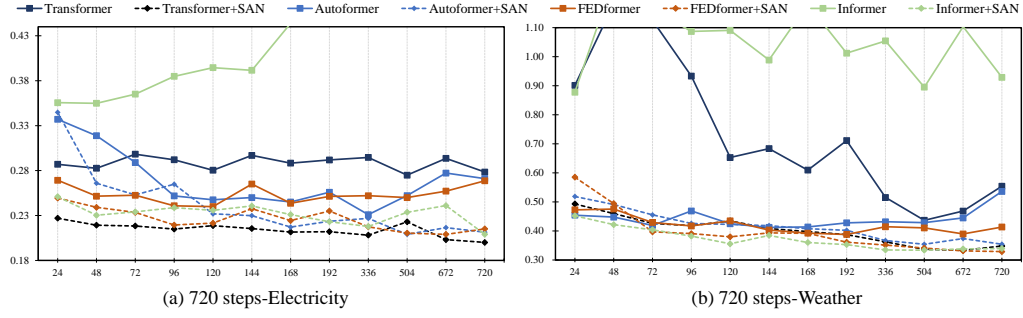


Figure 1: The long-term forecasting MSE evaluations of different Transformer-based models under various input lengths. Large values are discarded to illustrate the overall trend better.

72 accurate forecasting as the length increases. To be specific, on the Weather dataset, Transformer  
 73 achieves a reduction on MSE of **29.40%** when prolonging input from 24 steps to 720 steps, and the  
 74 average improvement on four backbones is **33.11%**. These results greatly meet our expectations and  
 75 also validate the effectiveness of SAN on various input lengths.

## 76 2.4 Additional Prediction Showcases

77 We provide the additional comparison between SAN and other normalization methods in Fig. 2 with  
 78 FEDformer [8] on various datasets. Clearly, SAN can better estimate the future distribution so as to  
 79 help the backbone model to achieve superior performance, where the forecasting results are better  
 aligned with the groundtruth.

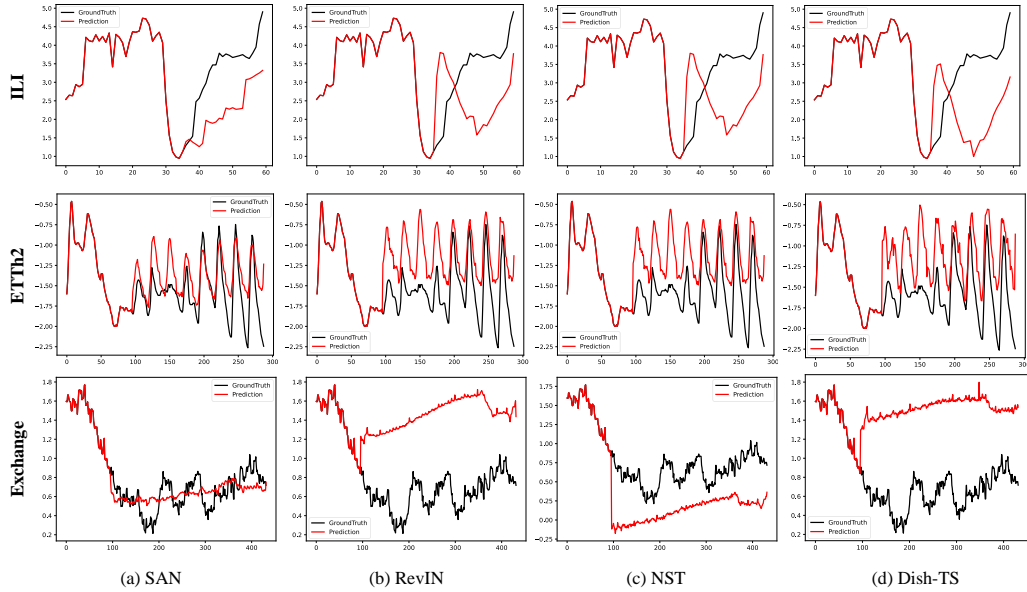


Figure 2: Illustration of the additional prediction showcases comparing SAN and baseline models. The experiment is conducted on the ILI, ETTh2, and Exchange dataset. Following the same input sequence length setting in our main experiments, the target sequence length is set to 24, 192, and 336 respectively.

## 2.5 Ablation Study

**Statistic Prediction Module** In this section, we aim to analyze the effectiveness of our designs in the statistic prediction module. We instantiate our method and its variants on Autoformer and test their performance on two typical non-stationary datasets: Exchange and ETTh2. Similarly, we repeat the experiments three times with fixed seed and report the evaluations with standard deviation in Table 3.

Table 3: Forecasting errors under the multivariant setting with respect to variants of SAN. The best performance are highlighted in **bold**.

Variants	Metric	SAN		w/o individual		w/o residual		w/o SAN	
		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Exchange	96	<b>0.082±0.001</b>	<b>0.208±0.001</b>	0.089±0.005	0.209±0.007	0.135±0.003	0.264±0.001	0.152±0.006	0.283±0.007
	192	<b>0.157±0.001</b>	<b>0.296±0.003</b>	0.184±0.009	0.306±0.010	0.331±0.044	0.416±0.025	0.369±0.055	0.437±0.033
	336	<b>0.262±0.004</b>	<b>0.385±0.002</b>	0.340±0.001	0.422±0.001	0.658±0.044	0.593±0.024	0.534±0.130	0.544±0.066
	720	<b>0.689±0.043</b>	<b>0.629±0.020</b>	0.982±0.001	0.753±0.002	1.456±0.011	0.882±0.009	1.222±0.099	0.848±0.021
ETTh2	96	<b>0.316±0.001</b>	<b>0.366±0.001</b>	0.321±0.013	0.367±0.008	0.383±0.020	0.413±0.012	0.384±0.021	0.420±0.013
	192	<b>0.413±0.013</b>	0.426±0.007	0.414±0.023	<b>0.422±0.012</b>	0.463±0.030	0.469±0.020	0.457±0.020	0.454±0.014
	336	<b>0.446±0.004</b>	0.457±0.003	0.448±0.003	<b>0.453±0.001</b>	0.586±0.025	0.541±0.009	0.468±0.010	0.473±0.005
	720	<b>0.471±0.009</b>	<b>0.474±0.005</b>	0.483±0.012	0.477±0.005	0.889±0.007	0.682±0.001	0.473±0.005	0.485±0.005

Obviously, with the proposed two techniques combined, the statistic prediction module can achieve the best accuracy, leading to optimal forecasting performance. Besides, both *residual learning* and *individual preference* contribute positive effects and the former one is much more important, without which SAN can even bring negative effects to the backbone model. These results validate the rationality of our thoughts about the characteristics of the mean value and also reveal the importance of accurate modeling of future statistics to SAN. Besides, SAN without individual modeling performs well on the ETT2 dataset but performs poorly on the Exchange dataset. Such a phenomenon reveals that the evolving trends of different scenarios vary, and it is required to model the complex relationships among multiple variables individually. Moreover, since we only incorporate the properties of mean values into a simple MLP network, how to design a proper mechanism or network architecture for statistics modeling is a promising direction for optimizing our method, and we leave such explorations for future work.

**Slicing Length** The slicing length is a key parameter of SAN. We aim to study the effect of different slicing lengths on our method. Ablation experiments are conducted by using SCINet as the backbone model under the long-term forecasting setting ( $L_{out} = 60$  for the ILI dataset and  $L_{out} = 720$  for the rest datasets). Each experiment is conducted three times with a fixed random seed. The forecasting errors and the corresponding standard deviation are presented in Table 4.

Table 4: Forecasting errors under the multivariant setting with respect to different slicing lengths. The best performance are highlighted in **bold**.

Slicing Length	Metric	6		12		24		48	
		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Electricity		0.210±0.002	0.305±0.002	0.207±0.002	<b>0.305±0.002</b>	<b>0.206±0.004</b>	0.307±0.003	0.208±0.002	0.307±0.001
Exchange		<b>0.892±0.028</b>	<b>0.712±0.013</b>	0.895±0.037	<b>0.712±0.017</b>	0.901±0.005	0.715±0.002	0.898±0.037	0.714±0.015
Traffic		0.612±0.001	<b>0.376±0.001</b>	0.608±0.002	<b>0.373±0.001</b>	<b>0.607±0.001</b>	0.381±0.001	0.611±0.002	0.382±0.002
Weather		<b>0.338±0.002</b>	0.366±0.002	<b>0.338±0.001</b>	<b>0.365±0.002</b>	0.340±0.001	0.367±0.001	0.339±0.001	0.366±0.001
ILI		<b>2.487±0.034</b>	<b>1.063±0.008</b>	2.680±0.055	1.118±0.015	n/a	n/a	n/a	n/a
ETTh1		0.491±0.002	0.475±0.001	<b>0.488±0.001</b>	0.474±0.001	0.489±0.004	<b>0.473±0.001</b>	0.492±0.004	0.474±0.002
ETTh2		0.440±0.001	0.465±0.001	<b>0.435±0.002</b>	0.460±0.002	0.437±0.007	<b>0.459±0.006</b>	0.443±0.007	0.462±0.004
ETThm1		0.495±0.043	0.469±0.024	<b>0.450±0.001</b>	<b>0.441±0.001</b>	0.611±0.218	0.503±0.084	0.463±0.006	0.448±0.003
ETThm2		<b>0.391±0.001</b>	0.406±0.001	<b>0.391±0.001</b>	<b>0.405±0.001</b>	0.392±0.001	<b>0.405±0.001</b>	0.403±0.009	0.415±0.006

Our heuristic selection of slicing length appears to be effective among the candidates, indicating that both artificially defined and actual periods are useful in selecting the optimal setting. Additionally, there were no significant performance differences observed under various settings, suggesting that SAN is resilient to changes in slicing length.

## 2.6 Detailed Results of the Comparison between SAN and Normalization Methods

In Table 5, we provide the detailed experimental results of the comparison between SAN and state-of-the-art normalization methods for non-stationary time series forecasting: RevIN [2], NST [3] and Dish-TS [1]. We re-implement the former two methods and Dish-TS is implemented by its official code<sup>1</sup>.

The table clearly shows that SAN outperforms existing approaches in most cases, except for the Weather dataset. Considering that the Weather dataset is the most stationary dataset, the results suggest that SAN can better remove the non-stationary factors in the raw data, even leading to an over-stationary issue that degrades the performance.

Besides, Dish-TS performs poorly in the benchmark. While it addresses the distribution shift between input and horizon series, it fails to optimize both the coefficient network and backbone network for overlooking the intrinsic bi-level optimization target of distribution estimation and forecasting tasks. By adopting a joint training schema, Dish-TS disturbs both networks and results in poor performance in certain cases. On the opposite, SAN benefits from the proposed two-stage schema which decouples the two tasks. This allows for proper optimization of each component and leads to improved overall performance.

Table 5: Detailed results of the comparison between SAN and normalization methods. The best results are highlighted in **bold**.

Methods	Metric	FEDformer								Autoformer							
		+ SAN		+ RevIN		+ NST		+ Dish-TS		+ SAN		+ RevIN		+ NST		+ Dish-TS	
		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Electricity	96	<b>0.164</b>	<b>0.272</b>	0.172	0.278	0.172	0.279	0.175	0.284	<b>0.172</b>	<b>0.281</b>	0.179	0.286	0.179	0.285	0.179	0.290
	192	<b>0.179</b>	<b>0.286</b>	0.185	0.289	0.187	0.291	0.188	0.296	<b>0.195</b>	<b>0.300</b>	0.216	0.316	0.209	0.309	0.215	0.318
	336	<b>0.191</b>	<b>0.299</b>	0.200	0.304	0.202	0.307	0.209	0.319	<b>0.211</b>	<b>0.316</b>	0.233	0.331	0.246	0.335	0.244	0.343
	720	<b>0.230</b>	0.334	0.243	0.337	<b>0.230</b>	<b>0.326</b>	0.239	0.343	<b>0.236</b>	<b>0.335</b>	0.246	0.341	0.252	0.345	0.286	0.370
Exchange	96	<b>0.079</b>	<b>0.205</b>	0.148	0.279	0.145	0.275	0.131	0.263	<b>0.082</b>	<b>0.208</b>	0.166	0.295	0.177	0.304	0.225	0.341
	192	<b>0.156</b>	<b>0.295</b>	0.266	0.377	0.274	0.383	0.538	0.523	<b>0.157</b>	<b>0.296</b>	0.299	0.404	0.275	0.385	0.760	0.610
	336	<b>0.260</b>	<b>0.384</b>	0.428	0.484	0.437	0.488	0.667	0.591	<b>0.262</b>	<b>0.385</b>	0.448	0.496	0.442	0.490	0.707	0.628
	720	<b>0.697</b>	<b>0.633</b>	1.056	0.789	1.064	0.787	1.480	0.954	<b>0.689</b>	<b>0.629</b>	1.068	0.791	1.049	0.784	2.341	1.063
Traffic	96	<b>0.536</b>	<b>0.330</b>	0.613	0.347	0.612	0.348	0.613	0.350	<b>0.569</b>	<b>0.350</b>	0.643	0.354	0.645	0.354	0.652	0.363
	192	<b>0.565</b>	<b>0.345</b>	0.637	0.356	0.641	0.357	0.644	0.362	<b>0.594</b>	<b>0.364</b>	0.659	0.373	0.643	0.367	0.669	0.374
	336	<b>0.580</b>	<b>0.354</b>	0.652	0.363	0.654	0.363	0.659	0.370	<b>0.591</b>	<b>0.363</b>	0.662	0.371	0.665	<b>0.363</b>	0.683	0.376
	720	<b>0.607</b>	<b>0.367</b>	0.686	0.382	0.688	0.380	0.693	0.388	<b>0.623</b>	0.380	0.700	0.384	0.667	<b>0.373</b>	0.703	0.392
Weather	96	<b>0.179</b>	0.239	0.187	<b>0.234</b>	0.187	<b>0.234</b>	0.244	0.317	<b>0.194</b>	0.256	0.212	0.257	0.211	<b>0.254</b>	0.268	0.338
	192	0.234	0.296	<b>0.235</b>	<b>0.272</b>	<b>0.235</b>	<b>0.272</b>	0.320	0.380	<b>0.258</b>	0.316	0.264	<b>0.300</b>	0.265	0.301	0.376	0.421
	336	0.304	0.348	<b>0.287</b>	<b>0.307</b>	0.289	0.308	0.424	0.452	0.329	0.367	0.309	0.329	<b>0.303</b>	<b>0.324</b>	0.476	0.486
	720	0.400	0.404	0.361	0.353	<b>0.359</b>	<b>0.352</b>	0.604	0.553	0.440	0.438	0.377	0.367	<b>0.366</b>	<b>0.357</b>	0.612	0.560
ILI	24	<b>2.461</b>	<b>1.075</b>	3.152	1.141	3.190	1.145	2.829	1.074	<b>2.548</b>	<b>1.098</b>	3.623	1.244	3.652	1.269	3.283	1.174
	36	<b>2.095</b>	<b>0.937</b>	2.498	0.990	2.615	1.026	2.595	0.972	<b>2.102</b>	<b>0.942</b>	2.767	1.074	2.394	1.005	2.792	1.054
	48	<b>2.107</b>	<b>0.936</b>	2.430	0.977	2.526	1.003	2.547	0.969	<b>2.103</b>	<b>0.932</b>	2.585	1.039	2.303	0.990	2.401	0.969
	60	<b>2.234</b>	<b>0.981</b>	2.822	1.084	2.891	1.109	2.866	1.066	<b>2.313</b>	<b>0.994</b>	2.693	1.068	2.489	1.008	2.681	1.042
ETTh1	96	<b>0.383</b>	<b>0.409</b>	0.392	0.413	0.394	0.414	0.390	0.424	0.522	0.474	0.491	0.463	0.550	0.503	<b>0.456</b>	<b>0.454</b>
	192	<b>0.431</b>	<b>0.438</b>	0.443	0.444	0.441	0.442	0.441	0.458	0.498	<b>0.472</b>	0.513	0.478	0.530	0.492	<b>0.495</b>	0.480
	336	<b>0.471</b>	<b>0.456</b>	0.495	0.467	0.485	0.466	0.495	0.486	0.571	0.509	0.528	0.485	<b>0.524</b>	<b>0.484</b>	0.539	0.496
	720	<b>0.504</b>	<b>0.488</b>	0.520	0.498	0.505	0.496	0.519	0.509	0.555	0.514	0.543	0.510	<b>0.510</b>	<b>0.491</b>	0.563	0.522
ETTh2	96	<b>0.300</b>	<b>0.355</b>	0.380	0.402	0.381	0.403	0.806	0.589	<b>0.316</b>	<b>0.366</b>	0.411	0.410	0.394	0.398	1.100	0.670
	192	<b>0.392</b>	<b>0.413</b>	0.457	0.443	0.478	0.453	0.936	0.659	<b>0.413</b>	<b>0.426</b>	0.478	0.450	0.473	0.450	0.976	0.672
	336	<b>0.459</b>	<b>0.462</b>	0.515	0.479	0.561	0.499	1.039	0.702	<b>0.446</b>	<b>0.457</b>	0.545	0.493	0.528	0.490	1.521	0.783
	720	<b>0.462</b>	<b>0.472</b>	0.507	0.487	0.502	0.481	1.237	0.759	<b>0.471</b>	<b>0.474</b>	0.523	0.490	0.524	0.498	1.105	0.745
ETTh1	96	<b>0.311</b>	<b>0.355</b>	0.340	0.385	0.336	0.382	0.348	0.397	<b>0.343</b>	<b>0.378</b>	0.458	0.446	0.468	0.448	0.477	0.460
	192	<b>0.351</b>	<b>0.383</b>	0.390	0.411	0.386	0.409	0.406	0.428	<b>0.390</b>	<b>0.400</b>	0.560	0.491	0.526	0.468	0.545	0.488
	336	<b>0.390</b>	<b>0.407</b>	0.432	0.436	0.438	0.441	0.438	0.450	<b>0.415</b>	<b>0.418</b>	0.607	0.508	0.786	0.559	0.650	0.533
	720	<b>0.456</b>	<b>0.444</b>	0.497	0.466	0.483	0.460	0.497	0.481	<b>0.476</b>	<b>0.453</b>	0.623	0.526	0.564	0.501	0.595	0.518
ETTh2	96	<b>0.175</b>	<b>0.266</b>	0.192	0.272	0.191	0.272	0.394	0.395	0.236	0.317	<b>0.233</b>	<b>0.307</b>	0.253	0.323	0.976	0.572
	192	<b>0.246</b>	<b>0.315</b>	0.270	0.320	0.270	0.321	0.552	0.472	<b>0.260</b>	<b>0.329</b>	0.288	0.337	0.289	0.335	0.532	0.485
	336	<b>0.315</b>	<b>0.362</b>	0.348	0.367	0.353	0.371	0.808	0.601	<b>0.330</b>	0.376	0.345	0.370	0.339	<b>0.365</b>	0.795	0.592
	720	<b>0.412</b>	0.422	0.430	<b>0.415</b>	0.445	0.422	1.282	0.771	<b>0.417</b>	0.428	0.434	<b>0.419</b>	0.426	0.432	1.271	0.768

<sup>1</sup><https://github.com/weifant/Dish-TS>

## 3 Implementation Details

### 3.1 Architecture of Statistic Prediction Module

The computation of  $\text{MLP}(x_1, x_2)$  in our paper can be summarized as follows:

$$\begin{aligned} x_1 &= \text{act}_1(W_1 * x_1) \\ x_2 &= \text{act}_1(W_2 * x_2) \\ x &= [x_1; x_2] \\ \text{output} &= \text{act}_2(W_3 * x) \end{aligned} \tag{4}$$

Here, the symbol  $[*; *]$  represents the concatenate operation. We set  $\text{act}_1(), \text{act}_2() = \text{Relu}(), \text{Relu}()$  for standard deviation and the activate function of the mean is set to  $\text{Tanh}(), \text{Identity}()$  respectively.  $W_1, W_2, W_3$  are learnable transformation matrices with hidden sizes of  $\{512, 512, 1024\}$ .

### 3.2 Algorithm of The Two-stage Training Schema

To apply SAN to backbone forecasting models, we propose a two-stage training schema to tackle the challenge of the bi-level optimization target. The statistics prediction module is first trained into convergence, which is then frozen and treated as a plugin during the second stage of training the forecasting model. We provide the pseudo-code of such a procedure in Alg. 1.

---

#### Algorithm 1 Two-stage Training Schema.

---

**Require:** Input series  $X = \{x^i\}_{i=1}^N$ ; Horizon series  $Y = \{y^i\}_{i=1}^N$ ; Slicing length  $T$

- 1: Initialize parameters  $\phi, \theta$
- 2: **while** not converge **do**
- 3:   **for all** input  $x^i \in X$ , horizon  $y^i \in Y$  **do**
- 4:     Compute input statistics  $\mu_j^i, \sigma_j^i$  by Eq. 1 with  $T$
- 5:     Predict future statistics  $\hat{\mu}^i, \hat{\sigma}^i$  by Eq. 3 using  $f_\phi(*)$
- 6:     Update  $\phi$  using loss function  $l_{sp}$
- 7:   **end for**
- 8: **end while** ▷ Training of the statistics prediction module
- 9:
- 10: **while** not converge **do**
- 11:   **for all** input  $x^i \in X$ , horizon  $y^i \in Y$  **do**
- 12:     Compute input statistics  $\mu_j^i, \sigma_j^i$  by Eq. 1 with  $T$
- 13:     Normalize input series to  $\bar{x}^i$  by Eq. 2
- 14:     Forecast  $\bar{y}^i = g_\theta(\bar{x}^i)$
- 15:     Predict future statistics  $\hat{\mu}^i, \hat{\sigma}^i$  by Eq. 3 using  $f_\phi(*)$
- 16:      $\hat{\mu}^i.\text{detach}(), \hat{\sigma}^i.\text{detach}()$  ▷ Stop-gradient, freeze the statistics prediction module
- 17:     Denormalize  $\bar{y}^i$  to  $\hat{y}^i$  by Eq. 4
- 18:     Update  $\theta$  using loss function  $l_{fc}$
- 19:   **end for**
- 20: **end while** ▷ Training of the forecasting model

---

## 4 Limitations

Though SAN shows promising performance on the benchmark dataset, there are still some limitations of this method. First is that we mainly select the slicing length heuristically or search in predefined candidates and the current design cannot handle indivisible length or the multi-period characteristic of time series. Such a solution works for the experiments but lacks generality in real-world applications. Second is that SAN may lead to an over-stationary issue, leading to sub-optimal performance. Therefore, a more flexible solution with automatic slicing length selection and normalization intensity control will be our exploring direction.

## References

- [1] Wei Fan, Pengyang Wang, Dongkun Wang, Dongjie Wang, Yuanchun Zhou, and Yanjie Fu. 2023. Dish-TS: A General Paradigm for Alleviating Distribution Shift in Time Series Forecasting. *arXiv preprint arXiv:2302.14829* (2023).
- [2] Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. 2021. Reversible instance normalization for accurate time-series forecasting against distribution shift. In *International Conference on Learning Representations*.
- [3] Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. 2022. Non-stationary Transformers: Rethinking the Stationarity in Time Series Forecasting. *arXiv preprint arXiv:2205.14415* (2022).
- [4] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [5] Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. 2021. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in Neural Information Processing Systems* 34 (2021), 22419–22430.
- [6] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. 2023. Are Transformers Effective for Time Series Forecasting? *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [7] Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. 2021. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 11106–11115.
- [8] Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. 2022. FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *Proc. 39th International Conference on Machine Learning (ICML 2022)* (Baltimore, Maryland).