

# [SUPPLEMENTARY MATERIAL] ENSURING FAIR COMPARISONS IN TIME SERIES FORECASTING: ADDRESSING QUALITY ISSUES IN THREE BENCHMARK DATASETS

**Anonymous authors**

Paper under double-blind review

## A MULTIVARIATE TIME SERIES (MTS) DATASETS FOUND IN LITERATURE

Table 1 lists the multivariate time series (MTS) datasets used in the following time series forecasting (TSF) papers: [1] LogTrans (Li et al., 2019), [2] Informer (Zhou et al., 2021), [3] Autoformer (Wu et al., 2021), [4] Pyraformer (Liu et al., 2022), [5] FEDformer (Zhou et al., 2022), [6] Triformer (Cirstea et al., 2022), [7] RevIn (Kim et al., 2022), [8] Preformer (Du et al., 2023), [9] ETSformer (Woo et al., 2023), [10] Crossformer (Zhang & Yan, 2023), [11] D·NLinear (Zeng et al., 2023), [12] TimesNet (Wu et al., 2023), [13] PatchTST (Nie et al., 2023) [14] RLinear (Li et al., 2024) and [15] iTransformer (Liu et al., 2024).

	Dataset																				
	electricity-f (fine)	electricity-c (coarse)	traffic-f (fine)	traffic-c (coarse)	Solar	Wind	M4-Hourly	ETT1 & m1	ETT2	ETTm2	ECL	Weather	App Flow	Electricity	Electricity	Weather	Exchange	ILI	Traffic	PeMS-Bay	Market
[1]	✓	✓	✓	✓	✓	✓	✓														
[2]								✓	✓		✓	✓									
[3]								✓	✓	✓					✓	✓	✓	✓	✓		
[6]								✓	✓		✓	✓									
[4]						✓							✓	✓							
[5]								✓	✓	✓					✓	✓	✓	✓	✓	✓	
[7]								✓	✓		✓										
[8]								✓	✓	✓					✓	✓	✓	✓	✓	✓	
[9]								✓	✓	✓					✓	✓	✓	✓	✓	✓	
[10]								✓	✓		✓	✓									
[11]								✓	✓	✓					✓	✓	✓	✓	✓	✓	
[12]							✓	✓	✓	✓					✓	✓	✓	✓	✓	✓	
[13]								✓	✓	✓					✓	✓			✓	✓	
[14]								✓	✓	✓	✓					✓					
[15]					✓			✓	✓	✓	✓					✓	✓			✓	✓

Table 1: Table listing all the datasets used in the selected papers.

## B SELECTED DATASET DESCRIPTIONS

In the previous table, it is evident that some datasets share the same name, such as *Electricity* and *Weather*. However, these datasets are actually different either when looking at the number of features or looking at the splitting. Furthermore, all “electricity” datasets— *electricity-f*, *electricity-c*, *ECL* and both *Electricity*— are actually variants of the same dataset: UCI electricity load diagrams (ELD).

In this study, we examined three real-world datasets for inconsistencies: (1) *Weather* from Informer (Zhou et al., 2021), which includes 12 meteorological indicators collected hourly at a Surface Weather Station in the U.S. from 2010 to 2013; (2) *Weather* from Autoformer (Wu et al., 2021) that comprises 21 meteorological variables collected every 10 minutes in 2020 from one of the Weather Station at the Max-Planck-Institute of Biogeochemistry; (3) *ELD* from UCI first introduced in (Li et al., 2019), which records the hourly electricity consumption of 370 clients from 2011 to 2014. These datasets were selected to clarify potential confusion in the multivariate time series forecasting (MTSF) literature.

To align with existing discussions (Han et al., 2024; Zhao & Shen, 2024), dataset variables or features (i.e., weather indicators or electricity clients) is referred as *channels* throughout this document.

### B.1 WEATHER FROM INFORMER

This dataset is derived from the local climatological data (LCD) dataset, which originally includes weather observations of 20 indicators from multiple worldwide stations. The Informer subset represents data from a single U.S. station collected between 2010 and 2013 (a more detailed description is provided in Appendix J).

We propose three revised versions of this dataset: (a) **LCDWf\_1H\_4Y\_USUNK**, where identified inconsistencies have been corrected; (b) **LCDWi\_1H\_4Y\_USUNK**, where inconsistencies have been corrected and the **usual** target channel has been rounded to “integer” values, **consistent with** other temperature channels in the dataset; and (c) **LCDWr\_1H\_4Y\_USUNK**, where inconsistencies have been corrected and the **usual** target channel is the **actual Fahrenheit value converted from the Celsius value using the known relation**.

The latter two versions aim to evaluate whether rounding to integer or direct conversion impacts predictive performance.

	From Informer (Zhou et al., 2021)	Corrected Proposal		
Dataset	Weather	LCDWf_1H_4Y_USUNK	LCDWi_1H_4Y_USUNK	LCDWr_1H_4Y_USUNK
Granularity	1H	1H		
Number of time steps	35 064	35 064		
Dataset Start Date	2010-1-1 00:00	2010-1-1 00:00		
Dataset End Date	2013-12-31 23:00	2013-12-31 23:00		
Number of channels	12	12 + 6 inconsistency identifier		
Target	WetBulbCelsius	WetBulbCelsius	WetBulbCelsiusInt	RealWetBulbCelsius

Table 2: Details of the Weather Dataset from Informer against the proposed corrected version.

### B.2 WEATHER FROM AUTOFORMER

This dataset is derived from the Max-Planck-Institute (MPI) dataset, initially comprising weather observations from three stations at the Max-Planck-Institute of Biogeochemistry in Germany. The Autoformer subset uses data collected from the station located on the roof of the building during 2020 (a more detailed description is provided in Appendix K).

We propose three revised versions of this dataset: (a) **MPIW\_10T\_1Y\_R** where identified inconsistencies have been corrected; (b) **MPIW\_10T\_4Y\_R** which extends the dataset to four years with a 10-minute resolutions and where identified inconsistencies have been corrected; and (c) **MPIW\_1H\_4Y\_R** which is the hourly resolution version of our extended revision.

	From Autoformer (Wu et al., 2021)	Corrected Proposal		
Dataset	Weather	MPIW_10T_1Y_R	MPIW_10T_4Y_R	MPIW_1H_4Y_R
Granularity	10T	10T	10T	1H
Number of time steps	52 696	52 705	210 284	35 064
Dataset Start Date	2020-1-1 00:10	2020-1-1 00:10	2020-1-1 00:10	2020-1-1 00:00
Dataset End Date	2021-1-1 00:00	2021-1-1 00:00	2024-1-1 00:00	2023-12-31 23:00
Number of channels	21	21 + 5 inconsistency identifier		
Target	OT	CO2 (ppm)		

Table 3: Details of the Weather Dataset from Autoformer against the proposed corrected version.

### B.3 ECL

This dataset is derived from the ELD dataset, originally providing 15-minute electricity consumption data for 370 clients collected between 2011 and 2014. The version used in many paper aggregates this to an hourly resolution and focuses on 321 clients from 2012 to 2014 – excluding clients with excessive zero data in the first year (a more detailed description is provided in Appendix L).

	From Zhou et al. (2021)	Corrected Proposal
<b>Dataset</b>	ECL	PELD_1H_3Y_308
<b>Granularity</b>	1H	1H
<b>Number of time steps</b>	35 064	35 064
<b>Dataset Start Date</b>	2011-1-1 00:00	2011-1-1 00:00
<b>Dataset End Date</b>	2013-12-31 23:00	2013-12-31 23:00
<b>Number of channels</b>	321	308
<b>Target</b>	MT_320	MT_320

Table 4: **Details** of the ECL from Informer against the proposed corrected version.

We propose a revised version of this dataset: **PELD\_1H\_3Y\_308**, which further reduces the **dataset to 308 clients** by removing those with **unusual** profiles and the remaining clients with excessive missing data.

## C EXPERIMENT DETAILS

### C.1 SETUP

We utilized the ADAM optimizer with an initial learning rate of 0.0001 and L2 loss for model optimization. Each experiment is run three times, with a total of 25 epochs and **an early stopping patience** set to 5.

### C.2 IMPLEMENTATION

We used the original PyTorch implementations of Informer <sup>1</sup>, Autoformer <sup>2</sup>, NLinear and DLinear <sup>3</sup> as well as iTransformer <sup>4</sup>. All experiments were conducted using the default parameter values **outlined** in Table 5. Each iteration used a **unique seed selected** from the following set: {24, 1024, 2024}.

Parameters	Informer	Autoformer	iTransformer	xLinear
d_model		512		
n_heads		8		-
e_layers		2		-
d_layers		1		-
s_layers	"3,2,1"		-	-
d_ff		2048		-
moving_avg	-		25	-
factor	5		3	-
padding	0		-	-
distil		True		-
dropout		0.05		0.1
attn	"prob"		-	-
embed		"timeF"		-
activation		gelu		-
output_attention		False		-
mix	"store_false"		-	-
num_workers	0		10	-
batch_size		32		-
learning_rate		0.0001		-
des		"Exp"		-
loss		mse		-
lradj		"type1"		-
channel_independence	-		False	-
class_strategy	-		Projection	-

Table 5: List of the default parameters used in our experiments

### C.3 PLATFORM

All experiments were **executed** on one NVIDIA DGX-1 **system equipped** with Tesla P100 GPUs.

<sup>1</sup><https://github.com/zhouhaoyi/Informer2020>

<sup>2</sup><https://github.com/thuml/Autoformer/tree/main>

<sup>3</sup><https://github.com/honeywell21/DLinear>

<sup>4</sup><https://github.com/thuml/iTransformer/tree/main>

D RESULTS (COMPLETE)

Dataset	Horizon	Zhou et al. (2021)	Wu et al. (2021)	Zeng et al. (2023)		Liu et al. (2024)		
		InF	AutoF	NL	DL	DL	InF	iTransF
LCD (Informer weather)	96	-	-	-	-	-	-	-
	192	-	-	-	-	-	-	-
	336	0.297	-	-	-	-	-	-
	720	0.359	-	-	-	-	-	-
MPI (Autoformer weather)	96	-	0.266	0.182	<u>0.176</u>	0.196	0.300	<u>0.174</u>
	192	-	0.307	0.225	<u>0.220</u>	0.237	0.598	<u>0.221</u>
	336	-	0.359	<u>0.271</u>	<u>0.265</u>	0.283	0.578	0.278
	720	-	0.419	<u>0.338</u>	<u>0.323</u>	0.345	1.059	0.358
ECL	96	-	0.201	<u>0.141</u>	<u>0.140</u>	0.197	0.274	0.148
	192	-	0.222	<u>0.154</u>	<u>0.153</u>	0.196	0.296	0.162
	336	0.489	0.231	<u>0.171</u>	<u>0.169</u>	0.209	0.300	0.178
	720	0.540	0.254	<u>0.210</u>	<u>0.203</u>	0.245	0.373	0.225

Table 6: MSE performances for multivariate-to-multivariate (M2M) predictions reported in the different literature papers.

Dataset	Horizon	Zhou et al. (2021)	Wu et al. (2021)	Zeng et al. (2023)		Liu et al. (2024)		
		InF	AutoF	NL	DL	DL	InF	iTransF
LCD (Informer weather)	96	-	-	-	-	-	-	-
	192	-	-	-	-	-	-	-
	336	0.416	-	-	-	-	-	-
	720	0.466	-	-	-	-	-	-
MPI (Autoformer weather)	96	-	0.336	<u>0.232</u>	0.237	0.255	0.384	<u>0.214</u>
	192	-	0.367	<u>0.269</u>	0.282	0.296	0.544	<u>0.254</u>
	336	-	0.395	<u>0.301</u>	0.319	0.335	0.523	<u>0.296</u>
	720	-	0.428	<u>0.348</u>	0.362	0.381	0.741	<u>0.349</u>
ECL	96	-	0.317	<u>0.237</u>	<u>0.237</u>	0.282	0.368	<u>0.240</u>
	192	-	0.334	<u>0.248</u>	<u>0.249</u>	0.285	0.386	0.253
	336	0.528	0.338	<u>0.265</u>	<u>0.267</u>	0.301	0.394	0.269
	720	0.571	0.361	<u>0.297</u>	<u>0.301</u>	0.333	0.439	0.317

Table 7: MAE performances for M2M predictions reported in the different literature papers.

Table 6 and 7 list the M2M prediction performance reported in associated papers. Notably, the reproduced results for Informer and DLinear by the iTransformer authors deviate from the originally published results, which could potentially be attributed to differences in the random seed used. However, despite running each experiment three times, our reproduced results, presented in Tables 8 through 16, **closely align with those reported by the iTransformer authors**. This **consistency** suggests that our findings are **in line with** recent literature and that variations in performance might stem from differences in datasets or data-splitting strategies.

In these tables, **bold and underline values** represent the best performance (**lowest value per row**) for each model across different prediction horizons and datasets. Values highlighted in **blue** [resp. **purple**] denote the best [resp. second-best] performances (**lowest** value per column and prediction horizon) obtained for a given dataset among all considered models.

D.1 PORTUGUESE ELECTRICITY LOAD DIAGRAMS

Table 8 presents **predictions** results for **the ECL dataset** using **both** the Informer version (ECL 321) and our revised version (PELD\_1H\_3Y\_308), **evaluated** with a 7:1:2 ratio split. At first glance, our revised version appears to perform better than the Informer version –as depicted by **the underlined values**. However, this **apparent improvement** should be **viewed with caution**. The reduced number of channels in our **revised dataset may contribute** for the lower error rates. Specifically, for DLinear and iTransformer, the mean average error (MAE) for some prediction horizons shows similar average errors and standard deviations for both datasets, **with iTransformer occasionally performing worse**. These observations could be interpreted in two ways: (i) fewer channels lead to worse predictions overall, or (ii) the removed channels helped lower the error (i.e., removing some complexity and making prediction easier). Despite this uncertainty, we **argue that** our revised version offers a fairer comparison of model performance. Ultimately, with our revised dataset, the ranking trend remains consistent with the literature: iTransformer **outperforms** DLinear, which in turn **surpasses** Informer.

Table 9 **reports model** performance with cycle-inclusive splitting, where each training, validation, and evaluation set covers one year of data chronologically. With this splitting **strategy**, our revised dataset **continues to yield** the best performances. However, errors for Informer with this strategy

216  
217  
218  
219  
220  
221  
222  
223  
224  
225

		7:1:2 (767.2/109.6/219.2 days)						Splitting	
		Published		Produced				Results	
		Reduced (ECL 321)				Revised (PELD.1H.3Y_308)		Dataset	
		F	MSE	MAE	MSE	MAE	MSE	MAE	Metric
DLinear	96	0.140	0.237	<u>0.195±0.0001</u>	<u>0.278±0.0001</u>	<u>0.192±0.0001</u>	<u>0.277±0.0001</u>		
	192	0.153	0.249	<u>0.194±0.0000</u>	<u>0.280±0.0000</u>	<u>0.191±0.0000</u>	<u>0.280±0.0000</u>		
	336	0.169	0.267	<u>0.207±0.0001</u>	<u>0.296±0.0004</u>	<u>0.202±0.0000</u>	<u>0.295±0.0000</u>		
	720	0.203	0.301	<u>0.242±0.0001</u>	<u>0.329±0.0003</u>	<u>0.235±0.0002</u>	<u>0.326±0.0003</u>		
Informer	96	-	-	0.286±0.0037	0.381±0.0023	<u>0.249±0.0022</u>	<u>0.355±0.0019</u>		
	192	-	-	0.293±0.0012	0.385±0.0026	<u>0.250±0.0049</u>	<u>0.356±0.0044</u>		
	336	0.489	0.528	0.305±0.0091	0.396±0.0071	<u>0.274±0.0084</u>	<u>0.376±0.0046</u>		
	720	0.540	0.571	0.332±0.0151	0.410±0.0102	<u>0.279±0.0087</u>	<u>0.381±0.0078</u>		
iTrans.	96	0.148	0.240	<u>0.163±0.0002</u>	<u>0.253±0.0001</u>	<u>0.161±0.0003</u>	<u>0.254±0.0002</u>		
	192	0.162	0.253	<u>0.175±0.0002</u>	<u>0.263±0.0001</u>	<u>0.172±0.0001</u>	<u>0.264±0.0001</u>		
	336	0.178	0.269	<u>0.192±0.0001</u>	<u>0.280±0.0001</u>	<u>0.188±0.0003</u>	<u>0.280±0.0002</u>		
	720	0.225	0.317	/	/	/	/		

Table 8: Results with PELD (electricity dataset) for multivariate-to-multivariate predictions and a ratio splitting (7:1:2). Our experiments are run three times, both the average error and standard deviation are reported in this table.

226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240

increased significantly compared to the ratio splitting. This result may be due to the lack of samples in the training set or the **larger** number of samples in the evaluation set. A similar trend is observed with DLinear and iTransformer, **though the impact is less severe**, suggesting that **these models are less sensitive to the training and evaluation set size**. Conversely, it may indicate that Informer, which focuses on *temporal tokens*, struggles to capture channel relationships **effectively**, while iTransformer, which uses *variate tokens* and **processes** each **channel independently**, **delivers more robust performance**. To further **validate** this assumption, it would be **beneficial** to create a **four-year version** of the dataset. **This version would** increase the training sample size while significantly reducing the number of clients, allowing **for a more comprehensive** comparison between ratio-based and cycle-inclusive splitting strategies.

241  
242  
243  
244

Overall, these experiments suggest that iTransformer and DLinear outperform Informer for spatiotemporal MTS datasets, with iTransformer achieving the best performance. These preliminary findings should be extended to other spatiotemporal MTS datasets, such as Traffic or Weather (datasets using multiple monitoring locations but focusing on only one observation).

245  
246  
247  
248  
249  
250  
251  
252  
253  
254

		1/1/1 year (366/365/365 days)				Splitting	
		F	Reduced (ECL 321)		Revised (PELD.1H.3Y_308)		Dataset
			MSE	MAE	MSE	MAE	Metric
DLinear	96		<u>0.208±0.0004</u>	<u>0.288±0.0006</u>	<u>0.197±0.0000</u>	<u>0.283±0.0001</u>	
	192		<u>0.207±0.0004</u>	<u>0.290±0.0007</u>	<u>0.194±0.0002</u>	<u>0.284±0.0005</u>	
	336		<u>0.226±0.0002</u>	<u>0.310±0.0002</u>	<u>0.207±0.0001</u>	<u>0.301±0.0001</u>	
	720		<u>0.277±0.0023</u>	<u>0.350±0.0019</u>	<u>0.246±0.0009</u>	<u>0.337±0.0010</u>	
Informer	96	0.624±0.0151	0.542±0.0087	0.480±0.0257	0.510±0.0192		
	192	0.639±0.0166	<u>0.563±0.0108</u>	0.555±0.0450	0.558±0.0299		
	336	0.680±0.0544	0.583±0.0284	<u>0.520±0.0204</u>	<u>0.535±0.0144</u>		
	720	0.879±0.1037	0.662±0.0500	<u>0.589±0.0121</u>	<u>0.571±0.0068</u>		
iTrans.	96	<u>0.179±0.0003</u>	<u>0.260±0.0001</u>	<u>0.171±0.0002</u>	<u>0.259±0.0001</u>		
	192	<u>0.189±0.0002</u>	<u>0.269±0.0001</u>	<u>0.180±0.0001</u>	<u>0.268±0.0001</u>		
	336	<u>0.208±0.0001</u>	<u>0.287±0.0001</u>	<u>0.197±0.0001</u>	<u>0.285±0.0001</u>		
	720	/	/	/	/		

Table 9: Results with PELD (electricity dataset) for multivariate-to-multivariate predictions and a cycle splitting (1/1/1 year). Our experiments are run three times, both the average error and standard deviation are reported in this table.

255  
256  
257  
258  
259  
260  
261  
262

## D.2 LOCAL CLIMATOLOGICAL DATA

### D.2.1 M2M

263  
264  
265  
266  
267  
268  
269

Table 10 presents the results for LCD and M2M predictions using the Informer version of the dataset, our corrected version (LCDWf\_1H.4Y\_USUNK), and our corrected version with the target in “integer” form (LCDWi\_1H.4Y\_USUNK) using a ratio splitting (7:1:2). Our corrected versions exhibit better performance across all the models used (as depicted by underlined results). The float version (LCDWf\_1H.4Y\_USUNK version), which keeps *WetBulbCelsius* as a *float*, generally performs better than the integer version than the integer version (LCDWi\_1H.4Y\_USUNK), suggesting that our corrections enhance model performance and dataset understanding. Therefore, such corrected versions should be preferred for fairer TSF model comparisons. iTransformer consistently outperforms

		Published		Produced						Results
				7:1:2 (i.e., 24544.8 / 3506.4 / 7012.8 days)						Splitting
		Original		LCDWf.1H.4Y_USUNK		LCDWi.1H.4Y_USUNK		Dataset		
		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	Metric
NLinear	F									
	96	-	-	0.520±0.0012	0.498±0.0006	<b>0.504</b> ±0.0015	<b>0.491</b> ±0.0009	0.504±0.0015	0.492±0.0009	
	192	-	-	0.589±0.0002	<b>0.542</b> ±0.0001	<b>0.572</b> ±0.0013	<b>0.535</b> ±0.0006	<b>0.572</b> ±0.0013	0.536±0.0006	
	336	-	-	<b>0.624</b> ±0.0006	<b>0.565</b> ±0.0001	<b>0.606</b> ±0.0006	<b>0.558</b> ±0.0001	<b>0.606</b> ±0.0006	<b>0.558</b> ±0.0001	
Informer	720	-	-	0.688±0.0002	0.601±0.0001	<b>0.669</b> ±0.0002	<b>0.595</b> ±0.0000	<b>0.669</b> ±0.0002	<b>0.595</b> ±0.0000	
	96	-	-	<b>0.482</b> ±0.0030	<b>0.490</b> ±0.0022	0.472±0.0031	0.483±0.0048	<b>0.466</b> ±0.0021	<b>0.483</b> ±0.0043	
	192	-	-	<b>0.586</b> ±0.0104	0.548±0.0116	0.567±0.0025	0.540±0.0044	<b>0.562</b> ±0.0076	<b>0.533</b> ±0.0044	
	336	0.702	0.620	0.627±0.0067	0.586±0.0086	<b>0.610</b> ±0.0102	<b>0.579</b> ±0.0102	0.610±0.0103	0.580±0.0106	
iTrans.	720	0.831	0.731	<b>0.623</b> ±0.0137	<b>0.586</b> ±0.0091	<b>0.598</b> ±0.0103	<b>0.575</b> ±0.0067	<b>0.599</b> ±0.0094	<b>0.576</b> ±0.0059	
	96	-	-	<b>0.509</b> ±0.0041	<b>0.487</b> ±0.0022	<b>0.492</b> ±0.0044	<b>0.480</b> ±0.0022	<b>0.492</b> ±0.0043	<b>0.481</b> ±0.0021	
	192	-	-	<b>0.577</b> ±0.0026	<b>0.533</b> ±0.0008	<b>0.559</b> ±0.0024	<b>0.526</b> ±0.0005	<b>0.559</b> ±0.0024	<b>0.526</b> ±0.0005	
	336	-	-	<b>0.609</b> ±0.0029	<b>0.555</b> ±0.0034	0.591±0.0026	<b>0.548</b> ±0.0032	<b>0.591</b> ±0.0025	<b>0.548</b> ±0.0031	
	720	-	-	<b>0.655</b> ±0.0033	<b>0.583</b> ±0.0026	<b>0.636</b> ±0.0041	<b>0.576</b> ±0.0029	<b>0.636</b> ±0.0042	<b>0.576</b> ±0.0029	

Table 10: Results with LCD (informer weather dataset) for multivariate-to-multivariate predictions and a ratio splitting (7:1:2). Our experiments are run three times, both the average error and standard deviation are reported in this table.

other models, with Informer often achieving the second-best results and occasionally surpassing iTransformer for specific prediction horizons.

		24/12/12 months (i.e., 17520 / 8784 / 8760 days)						Splitting
		Original		LCDWf.1H.4Y_USUNK		LCDWi.1H.4Y_USUNK		Dataset
		MSE	MAE	MSE	MAE	MSE	MAE	Metric
NLinear	F							
	96	0.582±0.0000	0.535±0.0000	<b>0.566</b> ±0.0001	<b>0.528</b> ±0.0000	<b>0.566</b> ±0.0001	<b>0.528</b> ±0.0000	
	192	0.660±0.0001	0.581±0.0001	<b>0.644</b> ±0.0001	<b>0.575</b> ±0.0001	<b>0.644</b> ±0.0001	<b>0.575</b> ±0.0001	
	336	0.680±0.0001	<b>0.597</b> ±0.0001	<b>0.663</b> ±0.0001	<b>0.591</b> ±0.0001	<b>0.663</b> ±0.0001	<b>0.591</b> ±0.0001	
Informer	720	0.741±0.0000	0.634±0.0000	<b>0.725</b> ±0.0000	<b>0.628</b> ±0.0000	<b>0.725</b> ±0.0000	<b>0.628</b> ±0.0000	
	96	<b>0.545</b> ±0.0064	<b>0.532</b> ±0.0050	<b>0.535</b> ±0.0099	0.529±0.0063	0.544±0.0170	<b>0.529</b> ±0.0058	
	192	<b>0.624</b> ±0.0013	0.570±0.0027	0.622±0.0047	0.571±0.0021	<b>0.620</b> ±0.0004	<b>0.571</b> ±0.0016	
	336	<b>0.661</b> ±0.0030	0.611±0.0056	<b>0.639</b> ±0.0004	<b>0.600</b> ±0.0046	0.639±0.0011	0.601±0.0042	
iTrans.	720	<b>0.673</b> ±0.0128	<b>0.619</b> ±0.0105	<b>0.650</b> ±0.0084	<b>0.609</b> ±0.0085	<b>0.648</b> ±0.0100	<b>0.608</b> ±0.0092	
	96	<b>0.562</b> ±0.0004	<b>0.520</b> ±0.0015	<b>0.546</b> ±0.0002	<b>0.513</b> ±0.0012	<b>0.546</b> ±0.0002	<b>0.513</b> ±0.0012	
	192	<b>0.644</b> ±0.0014	<b>0.572</b> ±0.0028	<b>0.627</b> ±0.0016	<b>0.565</b> ±0.0025	<b>0.627</b> ±0.0015	<b>0.565</b> ±0.0025	
	336	<b>0.669</b> ±0.0010	<b>0.590</b> ±0.0010	<b>0.652</b> ±0.0012	<b>0.584</b> ±0.0010	<b>0.652</b> ±0.0014	<b>0.584</b> ±0.0010	
	720	<b>0.723</b> ±0.0030	<b>0.623</b> ±0.0029	<b>0.702</b> ±0.0031	<b>0.612</b> ±0.0025	<b>0.703</b> ±0.0029	<b>0.612</b> ±0.0025	

Table 11: Results with LCD (informer weather dataset) for multivariate-to-multivariate predictions and a cycle splitting (24/12/12 months). Our experiments are run three times, both the average error and standard deviation are reported in this table.

Table 11 lists the models’ performance on M2M predictions with cycle-inclusive splitting. Here, the training set spans approximately two years, while validation and evaluation sets each cover one year chronologically. Although the metrics are worse compared to ratio splitting –likely due to reduced training samples and increased evaluation samples– the corrected dataset versions still perform better. With the cycle-inclusive splitting, Informer consistently offers lower mean squared error (MSE), while iTransformer provides the second-best results. For MAE, both models show competitive performance, but iTransformer generally performs better.

These findings suggest that Informer might be more suitable for MTS datasets with direct relations among channels, such as the electricity transformer temperature (ETT) dataset. Moreover, using a cycle-inclusive split challenges iTransformer previous superiority.

### D.2.2 UNIVARIATE-TO-UNIVARIATE (U2U)

Table 12 demonstrates that for LCD with ratio splitting and for U2U predictions, iTransformer clearly takes the lead over the other models. However, for the longest prediction horizon (i.e., 720), Informer achieves the best performance.

Table 13 indicates that cycle-inclusive splitting for U2U predictions also challenges iTransformer’s superiority. On average, Informer performs better, although iTransformer shows the best performance on our corrected dataset (LCDWf.1H.4Y\_USUN) in terms of MAE.

As a conclusion, these experiments suggest that with ratio splits, iTransformer is the leading model for both M2M and U2U predictions. Contrary to previous studies, Informer outperforms NLinear and even surpasses iTransformer for the 720 prediction horizon, suggesting that *ProbSparse* attention may be particularly beneficial for long prediction horizons. Further experiments comparing iTransformer, Informer, and inverse Informer for very large prediction horizons ( $\geq 720$ ) are re-

		Published		Produced						Results
				7:1:2 (i.e., 24544.8 / 3506.4 / 7012.8 days)						Splitting
		Original		LCDWf.IH.4Y.USUNK		LCDWf.IH.4Y.USUNK		Dataset		
F		MSE	MAE	MSE	MAE	MSE	MAE	Metric		
NLinear	96	-	-	0.196±0.0003	0.316±0.0004	0.180±0.0003	0.304±0.0004	0.183±0.0004	0.309±0.0005	
	192	-	-	0.252±0.0001	0.363±0.0000	0.235±0.0002	0.350±0.0002	0.238±0.0002	0.354±0.0002	
	336	-	-	0.297±0.0003	0.398±0.0002	0.279±0.0002	0.386±0.0002	0.283±0.0002	0.390±0.0002	
	720	-	-	0.393±0.0001	0.465±0.0001	0.377±0.0002	0.455±0.0001	0.380±0.0002	0.458±0.0001	
Informer	96	-	-	0.206±0.0087	0.340±0.0182	0.184±0.0123	0.321±0.0200	0.193±0.0071	0.328±0.0147	
	192	-	-	0.246±0.0095	0.370±0.0059	0.221±0.0136	0.349±0.0164	0.223±0.0157	0.351±0.0169	
	336	0.297	0.416	0.268±0.0206	0.398±0.0172	0.258±0.0195	0.390±0.0170	0.259±0.0154	0.391±0.0128	
	720	0.359	0.466	0.247±0.0081	0.385±0.0110	0.237±0.0073	0.378±0.0108	0.238±0.0078	0.379±0.0109	
iTrans.	96	-	-	0.191±0.0012	0.310±0.0016	0.176±0.0018	0.295±0.0015	0.178±0.0014	0.299±0.0013	
	192	-	-	0.235±0.0017	0.351±0.0025	0.217±0.0012	0.337±0.0018	0.221±0.0013	0.341±0.0019	
	336	-	-	0.265±0.0032	0.373±0.0005	0.245±0.0042	0.358±0.0014	0.248±0.0039	0.362±0.0017	
	720	-	-	0.292±0.0043	0.394±0.0021	0.263±0.0030	0.378±0.0025	0.266±0.0034	0.382±0.0012	

Table 12: Results with LCD (informer weather dataset) for univariate-to-univariate predictions and a ratio splitting (7:1:2). Our experiments are run three times, both the average error and standard deviation are reported in this table.

		24/12/12 months (i.e., 17520 / 8784 / 8760 days)						Splitting
		Original		LCDWf.IH.4Y.USUNK		LCDWf.IH.4Y.USUNK		Dataset
F		MSE	MAE	MSE	MAE	MSE	MAE	Metric
NLinear	96	0.251±0.0000	0.358±0.0000	0.233±0.0004	0.345±0.0004	0.237±0.0004	0.350±0.0004	
	192	0.310±0.0001	0.407±0.0001	0.292±0.0001	0.396±0.0001	0.296±0.0001	0.400±0.0001	
	336	0.352±0.0000	0.438±0.0000	0.334±0.0000	0.427±0.0000	0.339±0.0000	0.431±0.0000	
	720	0.442±0.0000	0.505±0.0000	0.426±0.0000	0.495±0.0000	0.430±0.0000	0.498±0.0000	
Informer	96	0.258±0.0041	0.367±0.0018	0.231±0.0028	0.345±0.0012	0.236±0.0030	0.351±0.0025	
	192	0.299±0.0005	0.408±0.0035	0.283±0.0041	0.397±0.0052	0.287±0.0027	0.401±0.0024	
	336	0.317±0.0067	0.425±0.0040	0.305±0.0047	0.416±0.0024	0.308±0.0040	0.418±0.0035	
	720	0.305±0.0182	0.419±0.0128	0.283±0.0064	0.408±0.0078	0.289±0.0083	0.411±0.0085	
iTrans.	96	0.257±0.0037	0.361±0.0033	0.237±0.0016	0.344±0.0016	0.241±0.0016	0.349±0.0016	
	192	0.304±0.0030	0.401±0.0023	0.286±0.0019	0.387±0.0021	0.290±0.0018	0.391±0.0020	
	336	0.346±0.0041	0.430±0.0023	0.324±0.0023	0.414±0.0009	0.336±0.0117	0.422±0.0052	
	720	0.362±0.0059	0.445±0.0041	0.337±0.0037	0.431±0.0032	0.340±0.0054	0.434±0.0041	

Table 13: Results with LCD (informer weather dataset) for univariate-to-univariate predictions and a cycle splitting (24/12/12 months). Our experiments are run three times, both the average error and standard deviation are reported in this table.

quired to investigate this finding. In addition, the results indicate that cycle-inclusive splits can re-define model rankings, with iTransformer being second-best to Informer for both M2M and U2U predictions. To confirm these observations, extending the study to other MTS datasets with direct relations among channels, such as ETT, is recommended.

### D.3 MAX-PLANCK-INSTITUTE

#### D.3.1 M2M

		7:1:2 (8.4 / 1.2 / 2.4 months)								Splitting
		Original				Simple		Corrected		Results
		Published		Produced		Produced		Produced		Dataset
F		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	Metric
DLinear	96	0.176	0.237	0.195±0.0002	0.255±0.0020	0.242±0.0005	0.299±0.0012	0.252±0.0006	0.303±0.0009	
	192	0.220	0.282	0.237±0.0008	0.296±0.0013	0.293±0.0048	0.350±0.0082	0.306±0.0048	0.357±0.0092	
	336	0.265	0.319	0.285±0.0015	0.336±0.0024	0.341±0.0013	0.387±0.0026	0.356±0.0032	0.396±0.0032	
	720	0.323	0.362	0.349±0.0027	0.387±0.0045	0.412±0.0011	0.446±0.0010	0.424±0.0021	0.445±0.0023	
Auto.	96	0.266±0.0007	0.336±0.0006	0.262±0.0094	0.340±0.0094	NA	NA	0.328±0.0107	0.389±0.0116	
	192	0.307±0.024	0.367±0.022	0.341±0.0154	0.396±0.0109	NA	NA	0.392±0.0110	0.428±0.0099	
	336	0.359±0.035	0.395±0.031	0.375±0.0275	0.413±0.0259	NA	NA	0.461±0.0220	0.476±0.0254	
	720	0.419±0.017	0.428±0.014	0.501±0.0350	0.492±0.0245	NA	NA	0.568±0.0312	0.542±0.0222	
iTrans.	96	0.174	0.214	0.174±0.0005	0.215±0.0015	0.218±0.0021	0.258±0.0015	0.227±0.0028	0.263±0.0022	
	192	0.221	0.254	0.225±0.0014	0.257±0.0008	0.278±0.0002	0.306±0.0003	0.291±0.0012	0.313±0.0008	
	336	0.278	0.296	0.281±0.0014	0.299±0.0006	0.340±0.0010	0.351±0.0014	0.351±0.0012	0.357±0.0012	
	720	0.358	0.349	0.360±0.0003	0.351±0.0004	0.426±0.0004	0.407±0.0006	0.441±0.0024	0.414±0.0013	

Table 14: Results with mpiw! (autoformer weather dataset) for multivariate-to-multivariate predictions and a ratio splitting (7:1:2). Our experiments are run three times, both the average error and standard deviation are reported in this table.

Table 14 presents the M2M prediction results for MPI using ratio splitting (7:1:2). We compare three versions of the dataset: the original version from Autoformer, a simple version where failure values (-9999) are replaced by 0 (Simple), and our corrected version using linear interpolation or context-aware imputation (MPIW\_10T\_1Y\_R). Our corrected version performs the worst among

these datasets, and even the Simple version underperforms compared to the original dataset, which retains the failure values. iTransformer outperforms other models for both the original and corrected datasets, with DLinear providing the second-best results.

		24/12/12 months				Splitting
		MPIW_10T_4Y_R		MPIW_1H_4Y_R		Granularity
	F	MSE	MAE	MSE	MAE	Metric
DLinear	96	0.417±0.0001	0.392±0.0002	0.504±0.0000	0.472±0.0000	
	192	0.478±0.0000	0.436±0.0001	0.562±0.0000	0.507±0.0000	
	336	0.542±0.0000	0.479±0.0002	0.601±0.0000	0.534±0.0001	
	720	0.615±0.0001	0.525±0.0001	0.664±0.0000	0.570±0.0000	
Auto.	96	0.416±0.0056	0.409±0.0030	0.598±0.0097	0.537±0.0072	
	192	0.551±0.0069	0.500±0.0047	0.642±0.0190	0.561±0.0091	
	336	0.613±0.0288	0.534±0.0150	0.669±0.0241	0.578±0.0105	
	720	0.668±0.0055	0.568±0.0028	0.713±0.0234	0.599±0.0099	
iTrans.	96	0.363±0.0005	0.336±0.0004	0.521±0.0015	0.470±0.0014	
	192	0.443±0.0002	0.394±0.0004	0.591±0.0002	0.510±0.0005	
	336	0.531±0.0020	0.449±0.0016	0.638±0.0012	0.540±0.0011	
	720	0.637±0.0021	0.512±0.0014	0.716±0.0004	0.578±0.0003	

Table 15: Results with MPI (autofomer weather dataset) for multivariate-to-multivariate predictions and a cycle splitting (24/12/12 months). Our experiments are run three times, both the average error and standard deviation are reported in this table.

Table 15 shows the performance with cycle-inclusive splitting and extended dataset versions: MPIW\_10T\_4Y\_R (with a 10-minute resolution) and MPIW\_1H\_4Y\_R (with an hourly resolution). Here, the training set spans approximately two years, while validation and evaluation sets each cover one year chronologically. Results with cycle-inclusive splitting are significantly worse than with ratio splitting, likely due to the significant increased sample size in the evaluation set and its comprehensive coverage of all seasons. This suggests potential overfitting in models trained on only 8.5 months and evaluated on 2.5 months. We note that DLinear performs better with the hourly dataset, whereas iTransformer excels with the 10-minute resolution, indicating DLinear’s difficulty with lower resolution cycles. Future work should verify if model performance varies across different seasons within the evaluation set. Notably, the hourly dataset performs worse than the 10-minute version, implying that our process for creating the hourly dataset may need revision.

Overall, these experiments suggest that iTransformer is the best model for MTS datasets. Extending these preliminary results to similar MTS datasets like Exchange would be valuable.

### D.3.2 U2U

		7:1:2 (8.4 / 1.2 / 2.4 months)								Splitting
		Original		Simple		Corrected		Dataset	Results	
	F	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	Metric
DLinear	96	-	-	0.005±0.0003	0.056±0.0027	0.387±0.0145	0.429±0.0071	0.555±0.0089	0.514±0.0029	
	192	-	-	0.006±0.0001	0.064±0.0006	0.476±0.0033	0.484±0.0021	0.651±0.0027	0.567±0.0012	
	336	-	-	0.006±0.0002	0.064±0.0019	0.527±0.0039	0.510±0.0025	0.743±0.0007	0.604±0.0002	
	720	-	-	0.006±0.0002	0.066±0.0021	0.595±0.0024	0.548±0.0012	0.947±0.0223	0.690±0.0093	
Auto.	96	-	-	0.003±0.0002	0.041±0.0017	NA	NA	0.767±0.0347	0.674±0.0182	
	192	-	-	0.004±0.0009	0.047±0.0063	NA	NA	0.767±0.0347	0.674±0.0182	
	336	-	-	0.004±0.0002	0.050±0.0016	NA	NA	0.940±0.0796	0.756±0.0353	
	720	-	-	0.004±0.0005	0.052±0.0030	NA	NA	1.205±0.0670	0.861±0.0259	
iTrans.	96	-	-	0.001±0.0000	0.027±0.0002	0.266±0.0020	0.360±0.0016	0.440±0.0103	0.456±0.0029	
	192	-	-	0.002±0.0000	0.029±0.0002	0.339±0.0007	0.414±0.0005	0.571±0.0043	0.532±0.0025	
	336	-	-	0.002±0.0000	0.031±0.0002	0.377±0.0028	0.444±0.0024	0.641±0.0157	0.573±0.0054	
	720	-	-	0.002±0.0000	0.035±0.0001	0.499±0.0046	0.516±0.0005	0.857±0.0124	0.671±0.0050	

Table 16: Results with MPI (autofomer weather dataset) for univariate-to-univariate predictions and a ratio splitting (7:1:2). Our experiments are run three times, both the average error and standard deviation are reported in this table.

For U2U predictions using ratio splitting, performance trends mirror those of M2M predictions: the corrected dataset yields worse results. The performance gap between the original and corrected datasets is significant, with the original dataset showing surprisingly low error values. This discrepancy likely arises from the impact of the failure value (-9999) on data normalization. Such an extreme value may distort z-normalization, affecting metrics calculated before reversing the normalization. Despite these issues, the corrected dataset maintains the same model ranking, with iTransformer performing best and DLinear second.



We believe our corrected dataset provides more accurate metric values, enabling fairer model comparisons.

		24/12/12 months				Splitting	
		MPIW_10T_4Y_R		MPIW_1H_4Y_R		Granularity	
		F	MSE	MAE	MSE	MAE	Metric
DLinear	96		<b>0.393</b> ±0.0001	<b>0.411</b> ±0.0002	0.425±0.0001	0.443±0.0001	
	192		<b>0.436</b> ±0.0001	<b>0.441</b> ±0.0001	0.476±0.0001	0.471±0.0000	
	336		<b>0.473</b> ±0.0000	<b>0.465</b> ±0.0001	0.514±0.0005	0.490±0.0002	
	720		<b>0.523</b> ±0.0001	<b>0.493</b> ±0.0001	0.568±0.0005	0.518±0.0003	
Auto.	96		<b>0.498</b> ±0.0181	<b>0.493</b> ±0.0116	0.501±0.0130	0.509±0.0118	
	192		0.695±0.0636	0.594±0.0324	<b>0.563</b> ±0.0100	<b>0.540</b> ±0.0042	
	336		0.715±0.0021	0.607±0.0032	<b>0.657</b> ±0.0373	<b>0.578</b> ±0.0047	
	720		0.717±0.0079	0.607±0.0031	<b>0.653</b> ±0.0268	<b>0.590</b> ±0.0127	
iTrans.	96		<b>0.330</b> ±0.0030	<b>0.363</b> ±0.0015	0.443±0.0032	0.455±0.0012	
	192		<b>0.402</b> ±0.0007	<b>0.414</b> ±0.0008	0.527±0.0020	0.502±0.0025	
	336		<b>0.461</b> ±0.0031	<b>0.451</b> ±0.0015	0.582±0.0031	0.528±0.0008	
	720		<b>0.530</b> ±0.0036	<b>0.490</b> ±0.0009	0.618±0.0040	0.551±0.0031	

Table 17: Results with MPI (autoformer weather dataset) for univariate-to-univariate predictions and a cycle splitting (24/12/12 months). Our experiments are run three times, both the average error and standard deviation are reported in this table.

Contrary to M2M predictions, U2U models trained with cycle-inclusive splitting outperform those using ratio splitting, as shown in Table 17. However, similarly to M2M predictions, DLinear excels with the hourly version, while iTransformer performs best with the 10-minute resolution dataset.

These findings suggest the need to revisit the generation of the hourly dataset and the temporal embedding implementation, which may influence model performance.

## E ADDITIONAL DISCUSSIONS

Our findings highlight the critical importance of clean datasets for improving model learning and ensuring fair model comparisons **across TSF models**. Notably, **our proposed cycle-inclusive splitting strategy** suggests that evaluating models over the **longest** temporal cycle **offers a more complete assessment of TSF model efficiency**. However, this outcome warrants further validation through experiments involving diverse datasets and alternative splitting strategies to fully assess the impact of varying sizes in training, validation, and evaluation sets.

Furthermore, our results suggest **that** no single model **consistently** excels in MTS forecasting. Instead, the optimal model or architecture may depend on the dataset’s **characteristics**. Models focusing on *variate tokens* tend to perform better on datasets **lacking** explicit inter-channel relationships (e.g., datasets monitoring different physical quantities **that are** not directly intertwined or spatiotemporal datasets where delays between channels **may occur**). In contrast, architectures based on *temporal tokens* appear more efficient when clear **and direct** relationships exist between channels. For example, despite **both** being weather datasets, the key difference between LCD and MPI is the explicitness of the relationships between weather indicators. **The LCD dataset** includes both Celsius and Fahrenheit **temperature readings**, providing explicit **interdependencies between channels**. Conversely, **The MPI dataset** may exhibit less explicit relationships, **where changes in one channel may influence others only after a delay**. Consequently, **models that effectively capture these relationships excel on datasets like LCD**. Therefore, **Informer, which prioritizes temporal tokenization, might captures “direct” inter-channel relationships more effectively, explaining its effectiveness with LCD**. On the other hand, iTransformer, which focuses on variate tokens, and linear-based models that treat each channel **independently, deliver superior performance on datasets like MPI and ELD, where inter-channel relationships are more complex**.

These insights **highlight the need for further experiments involving a broader range of transformer-based models and their variants**. Such studies could refine our understanding of model suitability for different dataset types, potentially guiding the development of tailored architectures for specific MTS forecasting tasks.

## F ADDITIONAL LIMITATIONS AND PERSPECTIVES

Beyond the limitations discussed in the main paper, additional issues must be addressed in the future.

486 Firstly, the current **approach for identifying errors on a per-time-step basis is inadequate**, particularly  
 487 when only a subset of channels **is affected by errors**. To improve this point, we plan to **create** separate  
 488 error masks **that pinpoint** error positions per time step and channel. Additionally, a dedicated file  
 489 containing only the proposed corrections should be created. This approach **would** simplify the use  
 490 of multiple correction versions, **eliminating** the need to **manage** multiple files **and reducing** storage  
 491 complexity and space requirements.

492 Secondly, we believe that the temporal embedding implementation, inherited from the Informer  
 493 model, also requires revision. Our experiments, particularly with the hourly and 10-minute reso-  
 494 lution versions of **the MPI dataset, reveal inconsistencies in its performance**. We **suggest revising**  
 495 **the encoding scheme to better capture cyclical patterns, which are prevalent** in TSF datasets. **A**  
 496 **more robust implementation could enhance the ability of models to represent and leverage temporal**  
 497 **dynamics effectively**.

## 499 G SOCIETAL IMPACTS

501 Time series forecasting (TSF) **plays a pivotal role in optimizing resource management and facilitat-**  
 502 **ing strategic economic planning** across various sectors. Accurate TSF **contributes to** (i) Enhanced  
 503 resource utilization, (ii) Reduced service disruptions and operational costs, and (iii) Better-informed  
 504 decision-making in domains such as energy, healthcare, finance, and logistics.

505 Our research underscores the necessity for clean datasets and **rigorous** model evaluation **methodolo-**  
 506 **gies**. **By advancing these aspects, we aim to improve TSF accuracy, foster a deeper understanding**  
 507 **of model strengths and limitations, and contribute to the development of more resilient, efficient**  
 508 **systems that benefit society as a whole**.

## 511 H HOSTING AND LICENSING

512 The following GitHub repository <sup>5</sup> is made available during the reviewing period **and contains** the  
 513 following **resources**:

- 515 • Code used for dataset analysis
- 516 • Code used for dataset correction
- 517 • Implementation of cycle-inclusive splitting dataloader
- 518 • CSV files of the revised versions of these datasets
- 519 • Experiment results in markdown format

522 The original dataset **used in this study are licensed as follows**:

- 524 • Electricity load diagrams (ELD) is available from UCI and **distributed** under a Creative  
 525 Commons Attribution 4.0 International (CC-BY-4.0) license;
- 526 • Local climatological data (LCD) is publicly available and according to the National  
 527 Oceanic and Atmospheric Administration (NOAA), it is “open and free to use. There are  
 528 no restrictions.”;
- 529 • Max-Planck-Institute (MPI) is publicly available and **distributed** under a Creative Com-  
 530 mons CC-BY-4.0 license.

532 The revised version of ELD and corrected versions of MPI adhere to the licensing **terms of their**  
 533 original datasets. The corrected version of LCD will be **distributed** under a Creative Commons  
 534 Attribution 4.0 International (CC-BY-4.0) license **to ensure consistency with open-access principles**.

535 **This repository aims to provide transparency, foster reproducibility, and encourage further research**  
 536 **in the field**. Upon acceptance of this paper, it would be important to include these dataset versions  
 537 on platforms such as HuggingFace <sup>6</sup> (as **updated** version of the existing datasets) or libraries such  
 538

539 <sup>5</sup><https://anonymous.4open.science/r/2392-NDBT-2AED/>

<sup>6</sup><https://huggingface.co/datasets?sort=trending>

as GluonTS <sup>7</sup>. This will increase their accessibility for future research and support comprehensive benchmarking of existing TSF models.

## I DATASETS PRESENTATION AND ANALYSIS METHOD

For AI models trained on data, presence of errors can severely hinder the learning of correlations and physical relationships, particularly if these errors are pervasive throughout the dataset. Furthermore, including time steps with inconsistencies in the evaluation set can significantly impair model assessment. A model performing well on an evaluation set that includes errors may either (i) excel on correct time steps while performing poorly on erroneous ones, demonstrating its ability to understand the data and its patterns, or (ii) it may perform moderately on both correct and erroneous time steps. However, the latter scenario does not necessarily indicate a robust model.

Benchmark datasets should ideally be free from such errors unless the objective explicitly **targets predictions** with erroneous data, tests model robustness to errors, or aims at anomaly detection. When evaluating TSF and comparing model performances, it is crucial to use datasets that are free from such errors, especially in the evaluation set. Therefore, it is essential to identify and annotate these problematic time steps, and to correct these errors **or at least** select appropriate metrics that account for them.

Our approach aims to address these concerns by proposing inconsistency-free dataset versions, **accompanied by** detailed annotations which will be beneficial for future research. The following sections present the inconsistencies found in each dataset, **the method used to identify them** and our proposed corrections.

### I.1 FREQUENCY

**To investigate the dominant frequencies in each dataset channel, we applied fast Fourier transform (FFT) using the following method:** (1) compute the trend of the considered channel, (2) apply the scipy FFT to the detrended channel, and (3) select the top K frequencies with the highest magnitude. Based on our experimentation, we adopted  $k = 3$  in this paper. **By** combining domain knowledge with frequency analysis, we can determine the overall (considering all channels) longest cycle for each dataset.

### I.2 DISTRIBUTION

For each dataset (original and revised versions), **distribution analyses were conducted for:** (i) the entire dataset, (ii) per longest cycle (mostly year), and (iii) per data splitting strategy. To better understand the impact of the data splitting in regard of seasonal **variations**, these distributions were plotted per solar season: *Spring, Summer, Autumn* and *Winter*. **When visually tractable, these distribution plots were performed for each channel. These plots** can help researcher understand the impact of splitting strategies that can introduce significant distribution differences between sets.

### I.3 CORRELATION

For each dataset, we performed four correlation analyses using *Pearson, Kendall, Spearman* and *Cosine similarity* methods. **To simplify the interpretation of the resulting heatmaps,** we focused on highly correlated channels by ignoring values between  $-0.75$  and  $0.75$ , which are represented as gray areas in the plots. Similarly to distribution plots, **correlation analyses were conducted for:** (i) the entire dataset, (ii) per longest cycle (considering all seasons), and (iii) per longest cycle and per solar season. Due to the **large** number of channels in ELD dataset and its variants, **correlation plots for these datasets were excluded from the analysis.**

<sup>7</sup><https://ts.gluon.ai/stable/index.html>

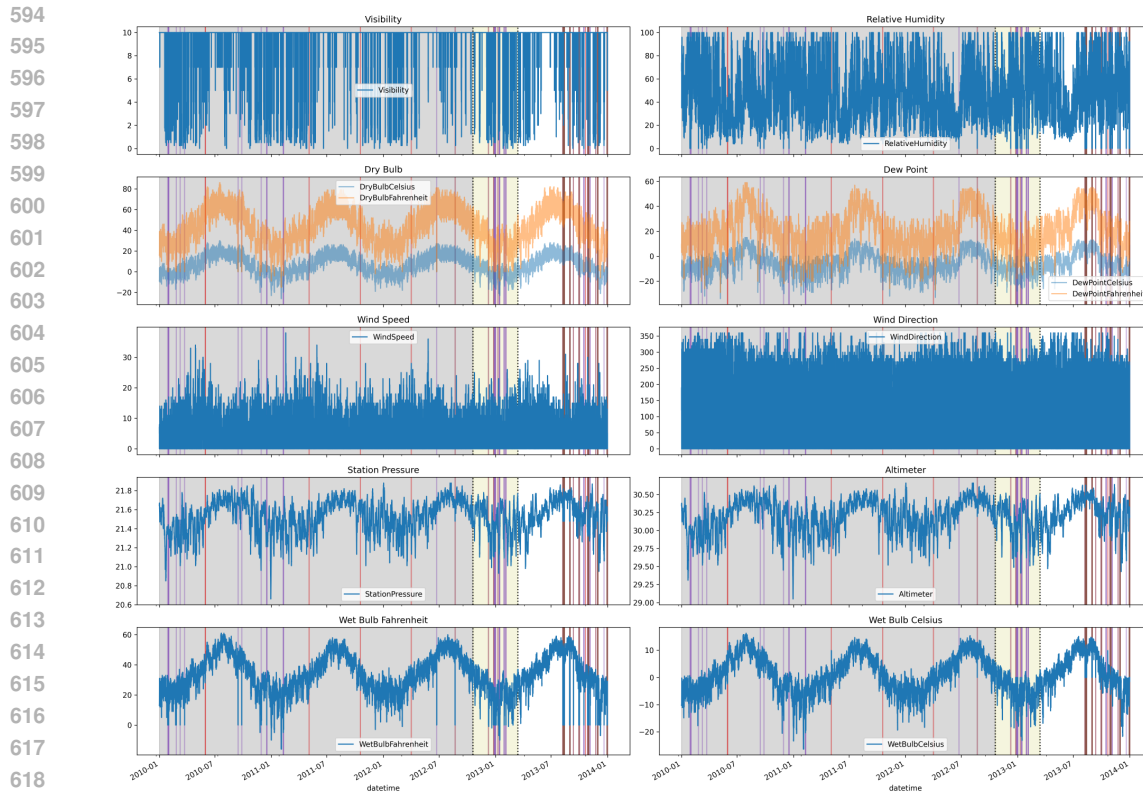


Figure 1: Overview of the weather indicators from the 4-year dataset used in Informer and collected from LCD. The gray background area represents the training period, and the yellow area represents the validation period as defined in the ratio splitting. Colored vertical lines indicate time steps where inconsistencies were found.

## J LOCAL CLIMATOLOGICAL DATA DATASET

### J.1 DESCRIPTION

The local climatological data (LCD) <sup>8</sup> dataset archives climatological data from approximately 20,000 stations worldwide, of which around 14,000 are active. For each station, surface observations are collected from various sources, including both manual and automated methods, and are managed by the National Centers for Environmental Information’s Integrated Surface Data (ISD). The dataset includes records of 20 weather indicators, such as dry bulb temperature in both Celsius and Fahrenheit, relative humidity, and more. Data in the archive spans from January 1<sup>st</sup>, 1901, to the present day, although the availability of data may vary significantly by station.

### J.2 ANALYSIS

The LCD dataset is a multi-variable spatiotemporal dataset consisting of observations from various weather stations. Researchers can utilize this dataset to explore the spatiotemporal relationships among the monitored physical quantities, investigating how different weather indicators interact over time. Additionally, it offers opportunities to analyze how artificial intelligence (AI) models learn and interpret fundamental unit conversions, such as the relationship between Celsius and Fahrenheit temperatures. Ultimately, the dataset facilitates studies aimed at predicting future values of individual or multiple weather indicators based on historical observations.

<sup>8</sup><https://www.ncei.noaa.gov/data/local-climatological-data/>

### 648 J.3 ORIGINAL VERSION

649 The version of the dataset selected by Zhou et al. (2021) introduce in the Informer paper provides  
650 hourly weather observations from a U.S. station over a 4-year period. It includes the following 12  
651 weather indicators:  
652

- 653 • Visibility (float)
- 654 • Dry bulb Temperature: Fahrenheit (integer), Celsius (integer)
- 655 • Wet bulb Temperature Fahrenheit (integer)
- 656 • Dew point Temperature: Fahrenheit (integer), Celsius (integer)
- 657 • Relative humidity (integer)
- 658 • Wind speed (integer)
- 659 • Wind direction (integer)
- 660 • Pressure (float)
- 661 • Altimeter (float)
- 662 • Wet bulb Temperature Celsius (float)

663 The dataset’s timestamp is unspecified regarding the time zone and spans from “2010-01-01  
664 00:00:00” to “2013-12-31 23:00:00” (included).  
665

#### 666 J.3.1 OVERALL ANALYSIS

667 Figure 1 displays grouped plots of the 12 weather indicators over the 4-year period. The gray  
668 and yellow areas represent the training and validation periods, respectively, as defined by the ratio  
669 splitting. At first glance, the dataset appears consistent; however, the vertical lines mark time steps  
670 where inconsistencies were identified.  
671

#### 672 J.3.2 FREQUENCY ANALYSIS

	Fundamental	2 <sup>nd</sup>	3 <sup>rd</sup>
673 Visibility	8766.0 (365.25)	389.6 (16.23)	313.1 (13.04)
674 DryBulbFahrenheit	8766.0 (365.25)	24.0 (1.00)	17532.0 (730.50)
675 DryBulbCelsius	8766.0 (365.25)	24.0 (1.00)	17532.0 (730.50)
676 WetBulbFahrenheit	8766.0 (365.25)	24.0 (1.00)	4383.0 (182.62)
677 DewPointFahrenheit	8766.0 (365.25)	4383.0 (182.62)	2922.0 (121.75)
678 DewPointCelsius	8766.0 (365.25)	4383.0 (182.62)	2922.0 (121.75)
679 RelativeHumidity	24.0 (1.00)	8766.0 (365.25)	4383.0 (182.62)
680 WindSpeed	24.0 (1.00)	12.0 (0.50)	4383.0 (182.62)
681 WindDirection	24.0 (1.00)	12.0 (0.50)	4383.0 (182.62)
682 StationPressure	8766.0 (365.25)	4383.0 (182.62)	407.7 (16.99)
683 Altimeter	8766.0 (365.25)	4383.0 (182.62)	407.7 (16.99)
684 WetBulbCelsius	8766.0 (365.25)	24.0 (1.00)	4383.0 (182.62)

685 Table 18: Frequency analysis of the original Weather dataset from Informer. The first value is the  
686 period in number of time steps the value in parentheses is the correspondence in days.  
687

688 This study reveals that most channels exhibit a primary cycle of one year (8766 time steps, eq.  
689 365.25 days). However, exceptions include *Relative Humidity*, *Wind Speed* and *Wind Direction*,  
690 which demonstrate a one-day cycle—an expected outcome for wind due to its inherently chaotic  
691 nature. The most prominent cycles identified in this dataset include one year, half a year, two years,  
692 one day, half a day, and approximately half a month. As a results, the longest dominant cycle across  
693 all channels is one year. Consequently, a cycle-inclusive splitting strategy should ensure that each  
694 set (training, validation, and evaluation) covers at least one full year to represent these temporal  
695 patterns effectively.  
696

### J.3.3 CORRELATION ANALYSIS

Figures 2 represents the channels correlation of the Weather dataset from Informer using the different methods mentioned in Appendix I.3. For all metrics, we can observe similar patterns:

1. **By row:** For each year, the correlations remain consistent or show only slight variations;
2. **By column:** Within a given period, when divided by solar seasons, the correlations between channels can vary significantly. For instance, *Winter* and *Spring* exhibit notable differences compared to *Summer* and *Autumn*. Additionally, while differences between *Winter* and *Spring*, as well as *Summer* and *Autumn*, are less pronounced, they are still evident.

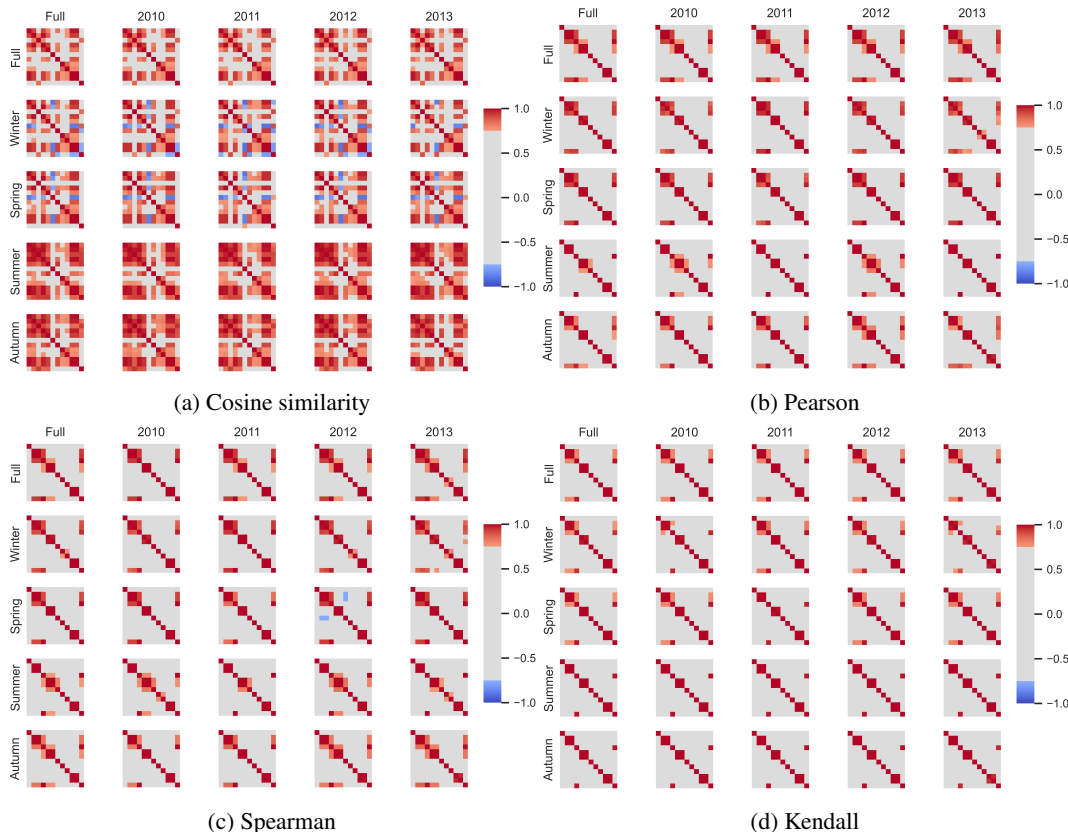


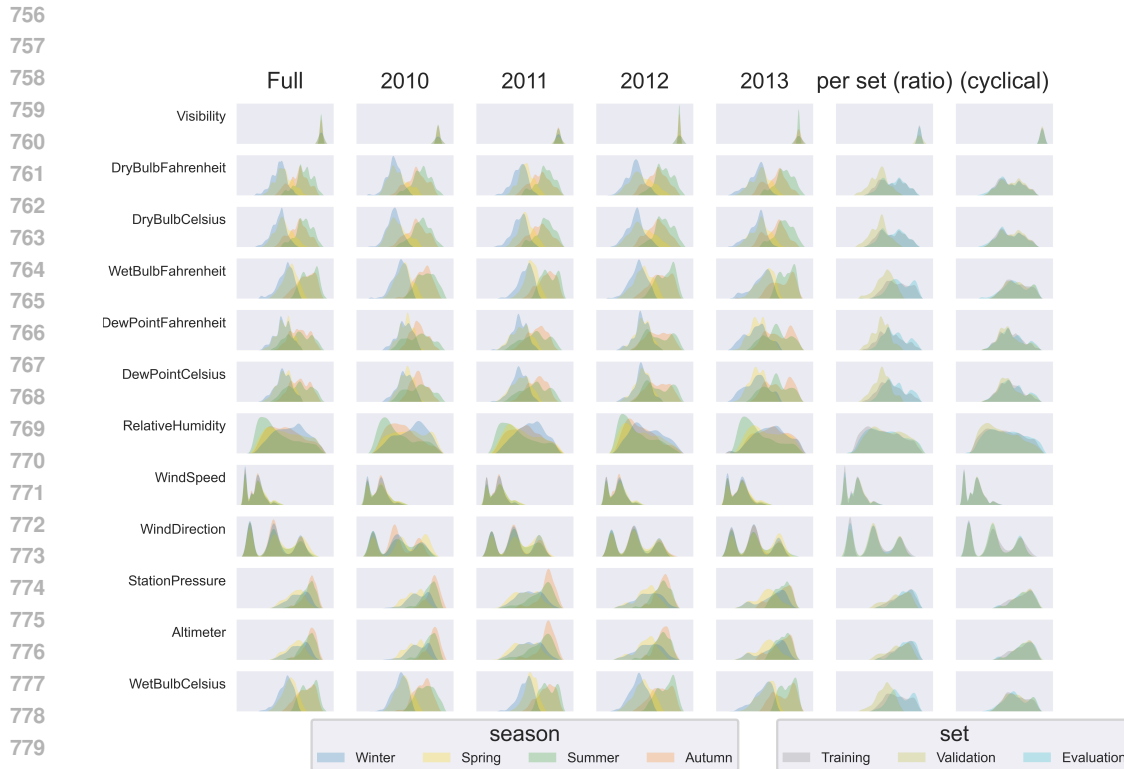
Figure 2: Weather Dataset from Informer - Channels correlation for the full dataset, per year and per season.

An efficient model for MTSF should be able to efficiently capture these seasonal variations and dynamically adapt the dependencies based on the input season.

### J.3.4 DATA DISTRIBUTION ANALYSIS

Figure 3 provides various distribution plots for the original dataset. As expected, most weather indicators exhibit distinct seasonal distributions, with the exception of *Visibility*, *Wind Speed* and *Wind Direction*. These seasonal fluctuations are especially significant for most of the channels. In addition, some variations can be observed across years, such as changes in the distribution of *Visibility* in 2012 and 2013 compared to 2010 and 2011. Any efficient MTSF model should be able to account for such differences and patterns in order to ensure robust performance.

Although the dataset provides enough data to consider a splitting strategy based on the longest cycle, Zhou et al. (2021) opted for a ratio splitting (7:1:2 ~ 28/10/10-month). This approach is not optimal for time series and chronological data because neither the validation nor the evaluation periods encompass a complete longest cycle, which, according to our frequency analysis, is one year.



782 Figure 3: Weather Dataset from Informer - Distribution plots per channel. The last two columns  
783 illustrate data distribution per splitting strategy: ratio and our proposal cycle-inclusive. The other  
784 columns illustrate the data distribution for the whole datasets and per year, with a differentiation per  
785 season.

787 Consequently, the training process is skewed to optimize performance for the selected validation  
788 period (i.e., *Winter*), while the evaluation period (i.e., *Spring*) does not fully test the model’s ability  
789 to generalize across the full cycle. **Our distribution and correlation analyses further highlight no-**  
790 **table differences between these periods, reinforcing the limitations of the ratio-based approach.** In  
791 addition, **Figure 3 demonstrates that** ratio splitting introduces significant distribution discrepancies  
792 between the training, validation, and evaluation sets. **In contrast, our cycle-inclusive splitting strat-**  
793 **egy mitigates these discrepancies, ensuring that the model is trained using a score that reflects the**  
794 **longest cycle and evaluated over a period covering an entire cycle.**

### 796 J.3.5 INCONSISTENCIES PRESENTATION

797 We identified several inconsistencies in the LCD dataset:

- 798  
799
- 800 1. **Missing Values Set to Zero:** Figure 4 highlights instances where missing values were in-  
801 appropriately set to zero. For example, it is not plausible for both Fahrenheit and Celsius  
802 values of the same indicator (e.g., Dew Point Temperature in the figure) to be zero simul-  
803 taneously at a given time step. Moreover, having the relative humidity also set to zero at  
804 this time step is inconsistent with surrounding values, which are close to 100%. Such an  
805 example advocates for missing data filled with zero.
  - 806 2. **Incorrect Fahrenheit to Celsius Conversion:** For the Wet Bulb Temperature feature,  
807 while the expected conversion from Fahrenheit to Celsius is affine, we observed signifi-  
808 cant errors. Figure 5a shows that for a Fahrenheit value of 32°F, the corresponding Celsius  
809 values range between  $-9.5^{\circ}\text{C}$  and  $9.9^{\circ}\text{C}$ , which is unacceptably wide and indicates a prob-  
lem with the data.

3. **Inconsistent Altimeter and Surface Pressure Relationship:** Figure 5b illustrates a somewhat staircase relationship between Altimeter and Surface Pressure. However, inconsistencies are evident when certain pressure values (e.g., 21.478686), where the altimeter values deviate significantly from the expected pattern. Such inconsistencies hinder the model’s ability to learn this relationship accurately.

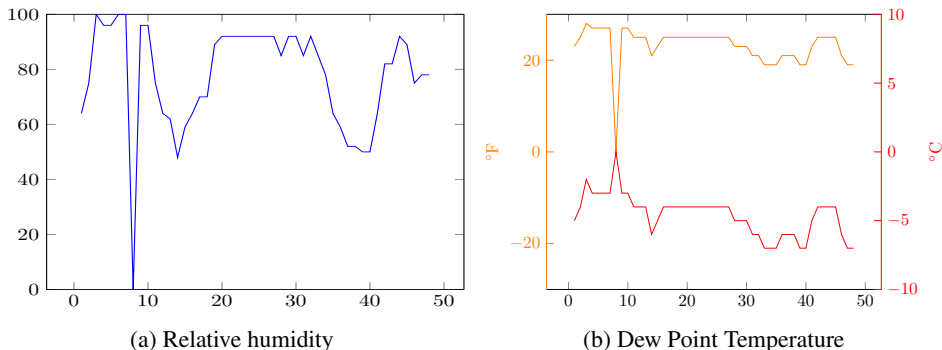


Figure 4: Visualization of LCD’s Relative Humidity and Dew Point Temperature for January 28-29, 2010. This figure highlights instances of missing values improperly set to zero, with both Relative Humidity and Dew Point Temperature showing simultaneous zero values, which are inconsistent with expected meteorological behavior.

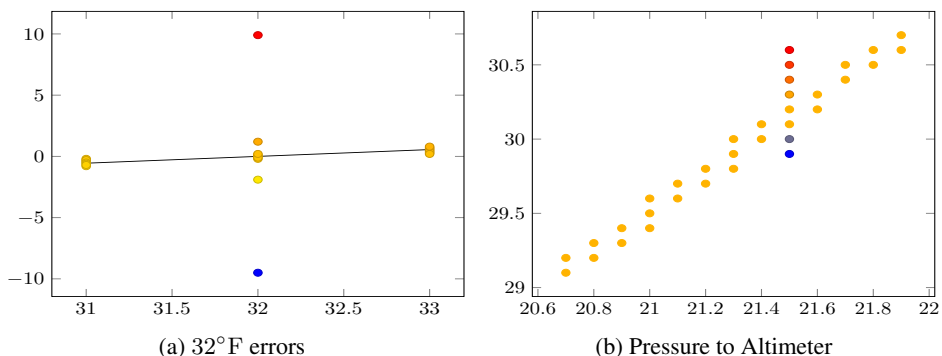


Figure 5: Visualization of Errors for 32°F Conversion and Altimeter-to-Pressure Relationship. In the left panel, the black line represents the affine function used for converting Fahrenheit to Celsius using the formula  $C = (F - 32) * 5/9$ . The red and blue points illustrate discovered inconsistencies in the dataset. In the right panel, the relationship between Altimeter and Surface Pressure is shown, highlighting deviations from the expected “staircase” pattern in Altimeter values for the pressure value 21.478686

### J.3.6 INCONSISTENCIES APPEARANCE

As shown by the colored vertical lines in Figure 1, inconsistencies are widespread throughout this dataset, but particularly present in the evaluation period.

The **red** vertical lines indicate time steps where errors in the 32°F values were identified, while the **purple** lines highlight time steps where missing data were inaccurately filled with zeros for multiple variables. **Pink** lines mark time steps where errors in Wet Bulb temperature conversions were found, and **brown** lines depict time steps where inconsistencies between pressure and altimeter values occurred.

### J.4 PROPOSED CORRECTION

To address the issues of missing data filled with zero and altimeter-to-pressure errors, we propose the following process outlined in the main paper: (i) replace erroneous values with NaN, (ii) apply



linear interpolation for isolated errors, and (iii) use either context-aware when possible or linear interpolation for consecutive errors.

Regarding the 32°F errors, we recommend replacing inconsistent values with values computed from the observed data and the well-known affine conversion. Specifically, 32°F converted values are identified as errors, if they deviate beyond the standard deviation of the correct converted data.

#### J.4.1 IDENTIFY INCONSISTENCIES

Six additional columns have been appended to the CSV file in order to identify the time step where inconsistencies were corrected:

- *32F\_errors*: identify time steps with 32°F errors
- *common\_conversion\_errors*: flags time steps where missing value were filled with zeros for a subset of variables.
- *wet\_conversion\_errors*: marks time steps with other conversion errors on Wet Bulb Temperature features.
- *pressure\_relation\_errors*: highlights time steps where altimeter-to-pressure errors were corrected.
- *is\_ts\_missing*: indicates time steps that were missing in the original dataset.
- *is\_ts\_modified*: logs all time steps where corrections were applied.

#### J.4.2 OVERALL ANALYSIS

Figure ?? displays grouped plots of the corrected 12 weather indicators over the 4-year period from the **LCDWf\_1H\_4Y\_USUNK** version. The gray and yellow areas represent the training and validation periods, respectively, as defined by the cycle-inclusive splitting. No data stand out which would imply that there are no errors in this version.

#### J.4.3 FREQUENCY ANALYSIS

	Fundamental	2 <sup>nd</sup>	3 <sup>rd</sup>
Visibility	8766.0 (365.25)	389.6 (16.23)	313.1 (13.04)
DryBulbFahrenheit	8766.0 (365.25)	24.0 (1.00)	17532.0 (730.50)
DryBulbCelsius	8766.0 (365.25)	24.0 (1.00)	17532.0 (730.50)
WetBulbFahrenheit	8766.0 (365.25)	24.0 (1.00)	4383.0 (182.62)
DewPointFahrenheit	8766.0 (365.25)	4383.0 (182.62)	2922.0 (121.75)
DewPointCelsius	8766.0 (365.25)	4383.0 (182.62)	2922.0 (121.75)
RelativeHumidity	24.0 (1.00)	8766.0 (365.25)	4383.0 (182.62)
WindSpeed	24.0 (1.00)	12.0 (0.50)	4383.0 (182.62)
WindDirection	24.0 (1.00)	12.0 (0.50)	4383.0 (182.62)
StationPressure	8766.0 (365.25)	4383.0 (182.62)	407.7 (16.99)
Altimeter	8766.0 (365.25)	4383.0 (182.62)	407.7 (16.99)
WetBulbCelsius	8766.0 (365.25)	24.0 (1.00)	4383.0 (182.62)

Table 19: **LCDWf\_1H\_4Y\_USUNK** - Frequency analysis. The first value is the period in number of time steps the value in parentheses is the equivalent in days.

The revised version **does not differ** from the original dataset in terms of dominant frequencies. Therefore, the longest cycle **remains** one year.

#### J.4.4 CORRELATION ANALYSIS

The correlation patterns observed in the revised dataset are consistent with those in the original dataset. This observation suggests that models still need to be capable of adapting dependencies based on seasonal variations.

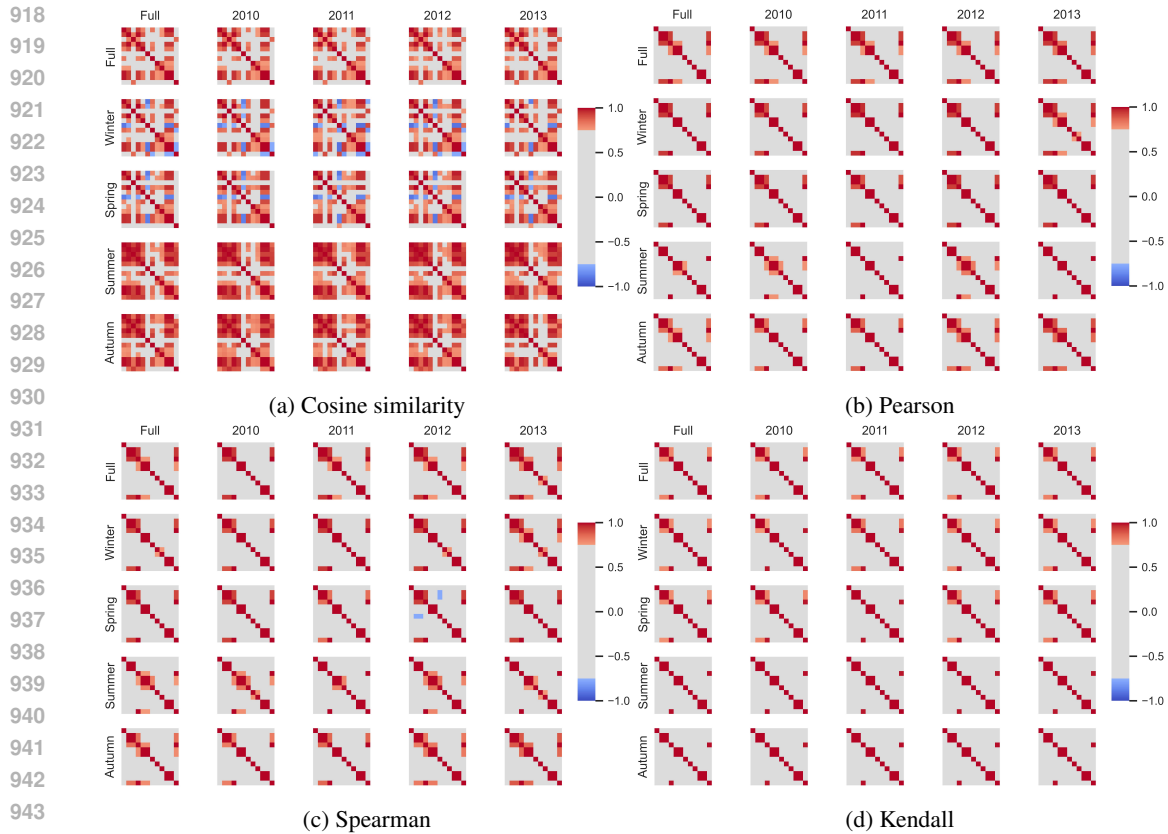


Figure 6: **LCDWf\_1H\_4Y\_USUNK** - Channels correlation for the full dataset, per year and per season.

#### J.4.5 DATA DISTRIBUTION ANALYSIS

Figure 7 presents the distribution plots for the revised dataset: **LCDWf\_1H\_4Y\_USUNK**. The corrections applied to address inconsistencies and errors have not altered the dataset’s inherent properties. While data distributions continue to vary significantly by season, our cycle-inclusive splitting strategy ensures better distributional similarity between the training, validation, and test sets. This strategy makes the dataset more suitable for benchmarking and facilitates more reliable model evaluations.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

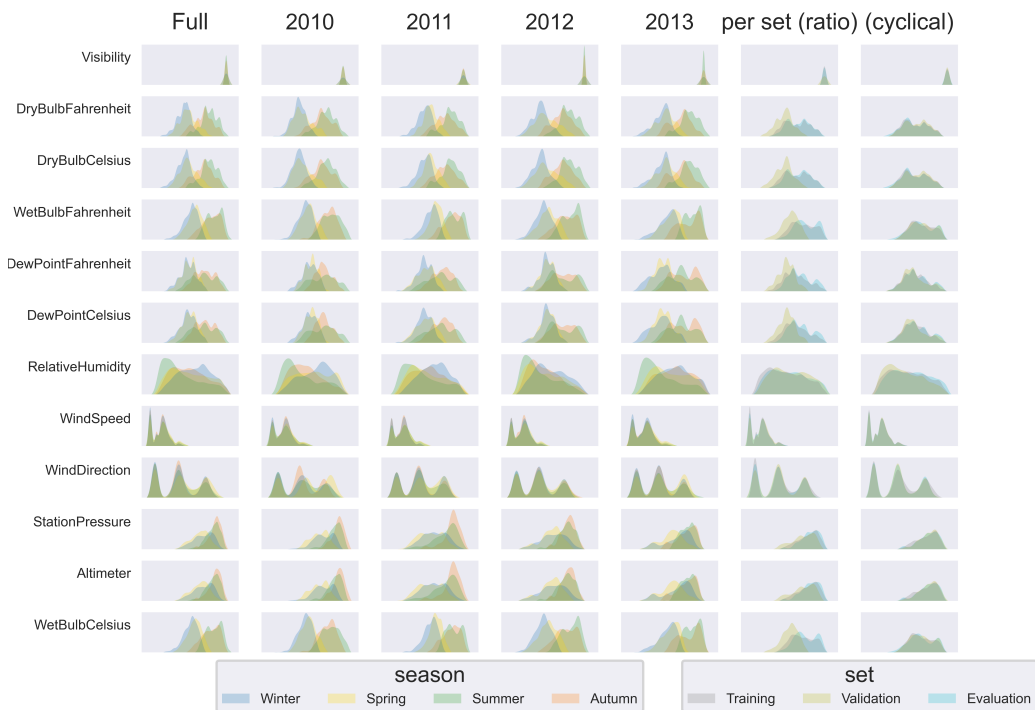


Figure 7: **LCDWf\_1H\_4Y\_USUNK** - Distribution plots per channel. The last two columns illustrate data distribution per splitting strategy: ratio and our proposal cycle-inclusive. The other columns illustrate the data distribution for the whole datasets and per year, with a differentiation per season.

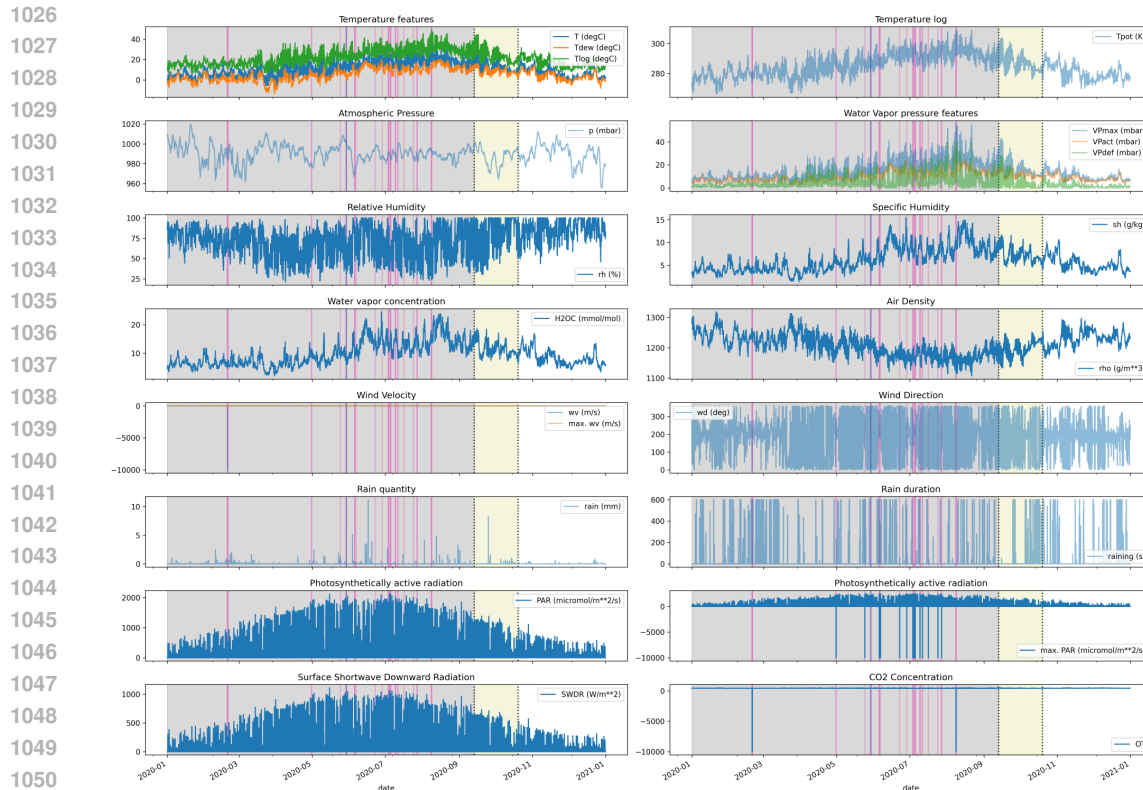


Figure 8: Overview of the weather indicators from the 1-year dataset used in Autoformer and collected from MPI. The gray background area represents the training period, while the yellow area denotes the validation period as defined in the ratio splitting. Colored vertical lines indicate time steps where inconsistencies were identified.

## K MAX-PLANCK-INSTITUTE DATASET

### K.1 DESCRIPTION

The Max-Planck-Institute (MPI)<sup>9</sup> dataset provides weather measurements collected from three distinct weather stations. One of these stations, *WS Beutenberg*, is located atop the building’s roof of the Max-Planck-Institute for Biogeochemistry. It comprises 21 weather indicators, including air temperature and humidity, recorded at 10-minute intervals. This dataset spans from “2003-11-24 16:00:00” to the present days.

### K.2 ANALYSIS

Similarly to LCD, the MPI dataset is a multi-variable spatiotemporal dataset. When focusing on data from a single station, the resulting dataset is a MTS dataset capturing observations from a specific location in Germany via various sensors. These observations exhibit variations intricately linked to Earth’s revolution (year, seasons) and rotation (day, hours). Other factors, such as human behavior and global warming, likely contribute to fluctuations in the recorded parameters. This dataset offers the opportunity for models to discern relationships between these indicators and leverage such insights to predict one or multiple variables.

<sup>9</sup><https://www.bgc-jena.mpg.de/wetter/>

### 1080 K.3 ORIGINAL VERSION

1081 Wu et al. (2021) selected a 1-year period—the year 2020—of the data available from *WS Beutenberg*,  
 1082 located on the roof of the Max-Planck-Institute for Biogeochemistry. It has a 10-minute resolution.  
 1083 This datasets includes weather observations of the following 21 indicators:  
 1084

- 1085 • Atmospheric Pressure (p (mbar))
- 1086 • Air Temperature (T (degC))
- 1087 • Potential Temperature (Tpot (K))
- 1088 • Dew Point Temperature (Tdew (degC))
- 1089 • Relative Humidity (rh (%))
- 1090 • Saturation Water Vapor Pressure (VPmax (mbar))
- 1091 • Actual Water Vapor Pressure (VPact (mbar))
- 1092 • Water Vapor Pressure Deficit (VPdef (mbar))
- 1093 • Specific Humidity (sh (g/kg))
- 1094 • Water Vapor Concentration ( $H_2O$  ( $\mu\text{mol/mol}$ ))
- 1095 • Air Density ( $\rho$  ( $\text{g/m}^3$ ))
- 1096 • Wind Velocity (wv (m/s))
- 1097 • Maximum Wind Velocity (max. wv (m/s))
- 1098 • Wind Direction (wd (deg))
- 1099 • Precipitation Amount (rain (mm))
- 1100 • Precipitation Duration (raining (s))
- 1101 • Surface Shortwave Downward Radiation (SWDR ( $\text{W/m}^2$ ))
- 1102 • Photosynthetic Active Radiation (PAR ( $\mu\text{mol/m}^2/\text{s}$ ))
- 1103 • Maximum Photosynthetic Active Radiation (max. PAR ( $\mu\text{mol/m}^2/\text{s}$ ))
- 1104 • Internal Logger Temperature (Tlog (degC))
- 1105 •  $\text{CO}_2$  concentration ( $\text{CO}_2$  (ppm))

1106 The timestamp are provided without any specific time zone. The dataset spans from “2020-01-01  
 1107 00:10:00” to “2021-01-01 00:10:00” (included).  
 1108

#### 1109 K.3.1 OVERALL ANALYSIS

1110 Figure 8 **presents** the plots of the different weather indicators. The gray area represents the training  
 1111 period, while the yellow area **indicates** the validation period, as defined by the ratio splitting strat-  
 1112 egy. The presence of errors is **particularly noticeable in plots** where the y-axis extends to values as  
 1113 extreme as  $-10000$ , which are clearly unrealistic for any of the weather indicators monitored.

1114 In addition, as this dataset spans only one year, the ratio splitting **approach trains on one part of the**  
 1115 **year and evaluates on another**, leading to a **highly** specific evaluation. This splitting method does  
 1116 not adequately represent the model’s ability to **produce accurate predictions across** the entire year,  
 1117 which poses problem for potential real-world applications.

#### 1118 K.3.2 FREQUENCY ANALYSIS

1119 **The frequency analysis indicates that 10 channels exhibit a dominant yearly cycle (52696 time steps,**  
 1120 **approximately 365.94 days). The remaining channels show dominant cycles of one day (7 channels),**  
 1121 **half a month (1 channel), two months (1 channel), two and a half months (1channel), and four**  
 1122 **months (1 channel). The most prominent cycles** in this dataset are one year, six months, four months,  
 1123 and one day.

	Fundamental	2 <sup>nd</sup>	3 <sup>rd</sup>
p (mbar)	10539.2 (73.19)	8782.7 (60.99)	4391.3 (30.50)
T (degC)	52696.0 (365.94)	144.0 (1.00)	26348.0 (182.97)
Tpot (K)	52696.0 (365.94)	144.0 (1.00)	26348.0 (182.97)
Tdew (degC)	52696.0 (365.94)	26348.0 (182.97)	17565.3 (121.98)
rh (%)	144.0 (1.00)	52696.0 (365.94)	17565.3 (121.98)
VPmax (mbar)	52696.0 (365.94)	144.0 (1.00)	26348.0 (182.97)
VPact (mbar)	52696.0 (365.94)	26348.0 (182.97)	10539.2 (73.19)
VPdef (mbar)	52696.0 (365.94)	144.0 (1.00)	143.6 (1.00)
sh (g/kg)	52696.0 (365.94)	26348.0 (182.97)	10539.2 (73.19)
H <sub>2</sub> OC ( $\mu$ mol/mol)	52696.0 (365.94)	26348.0 (182.97)	10539.2 (73.19)
rho (g/m <sup>3</sup> )	52696.0 (365.94)	144.0 (1.00)	26348.0 (182.97)
wv (m/s)	144.0 (1.00)	143.6 (1.00)	72.0 (0.50)
max. wv (m/s)	144.0 (1.00)	8782.7 (60.99)	143.6 (1.00)
wd (deg)	8782.7 (60.99)	52696.0 (365.94)	3293.5 (22.87)
rain (mm)	2107.8 (14.64)	17565.3 (121.98)	258.3 (1.79)
raining (s)	17565.3 (121.98)	893.2 (6.20)	958.1 (6.65)
SWDR (W/m <sup>2</sup> )	144.0 (1.00)	52696.0 (365.94)	143.6 (1.00)
PAR ( $\mu$ mol/m <sup>2</sup> /s)	144.0 (1.00)	52696.0 (365.94)	143.6 (1.00)
max. PAR ( $\mu$ mol/m <sup>2</sup> /s)	144.0 (1.00)	52696.0 (365.94)	143.6 (1.00)
Tlog (degC)	52696.0 (365.94)	144.0 (1.00)	26348.0 (182.97)
CO <sub>2</sub> (ppm)	144.0 (1.00)	4053.5 (28.15)	521.7 (3.62)

Table 20: Weather from Autoformer - Frequency analysis. The first value is the period in number of time steps the value in parentheses is the equivalent in days.

### K.3.3 CORRELATION ANALYSIS

Figure 9 represents the channels correlation of the Weather dataset from Autoformer using the different methods mentioned in Appendix I.3. Across all metrics, significant seasonal differences are observed:

- *Winter* and *Spring* exhibit correlations that differ substantially from those of *Summer* and *Autumn*
- Some smaller differences are also observed between *Winter* and *Spring*, as well as between *Summer* and *Autumn*.

An efficient MTSF model must effectively capture these seasonal variations and adapt the dependencies based on the input season.

### K.3.4 DATA DISTRIBUTION ANALYSIS

Similar to LCD, most weather indicators demonstrate distinct seasonal distributions, with significant fluctuations for several channels.

Figure 10 provides two data distribution plots for the original dataset: one per season and one per data splitting set. As expected, channels with inconsistencies or where failure values have been identified appear anomalous. Similarly to LCD, most weather indicators **demonstrate distinct seasonal distributions, with significant fluctuations for several channels.**

The lack of data **spanning multiple years** prevent from using a cycle-inclusive splitting strategy **with a one year dominant cycle**. Instead, Wu et al. (2021) **adopted** a ratio splitting ( $7:1:2 \sim (8.4/1.2/2.4)$  months). This approach implies that neither the validation nor the evaluation periods encompass a complete **longest cycle**. Consequently, the training process is skewed to optimize performance for the selected validation period (i.e., *Autumn*), while the evaluation (i.e., *Winter*) **fails to adequately test the model’s ability to generalize across the full cycle. As demonstrated by the distribution and correlation analyses, notable differences exist between these periods.** In addition, we observed in Figure 10 that the ratio splitting **strategy** implies significant distribution difference between training, validation and evaluation sets.

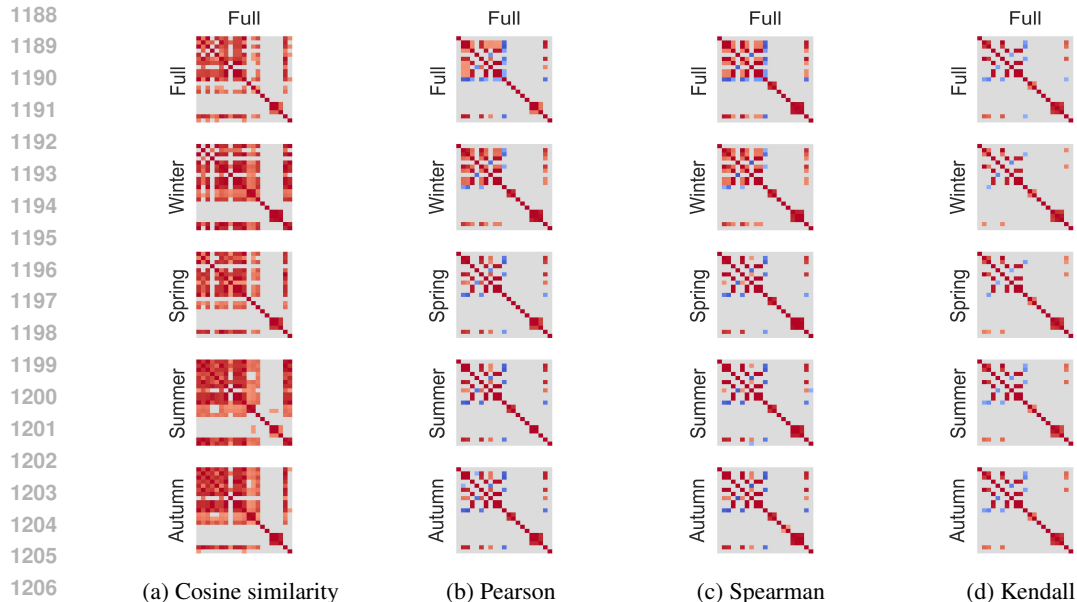


Figure 9: Weather Dataset from Autoformer - Channels correlation for the full dataset and per season.

### K.3.5 INCONSISTENCIES PRESENTATION

We identified three types of inconsistencies in the MPI dataset:

1. **Failure Values** ( $-9999$ ): Figure 8 shows instances where the value  $-9999$  appears throughout the dataset, likely indicating measurement failures or missing observations due to instrument errors.
2. **Duplicated entries**: We found duplicated entries with identical timestamp and values across all variables.
3. **Missing Time Step**: Certain time steps are missing from the original data.

These inconsistencies can present in the file provided by Autoformer as well as the data archives on the original website, as detailed in Table 21.

	Duplicated	Missing
Autoformer	1	9
2020a	1	9
2020b	0	0
2021a	0	0
2021b	0	0
2022a	6	0
2022b	1	82
2023a	0	3
2023b	143	0

Table 21: Count of duplicated entries and missing time steps found in the Autoformer CSV file and data archives from the Max-Planck-Institute original website.

### K.3.6 INCONSISTENCIES APPEARANCE

In Figure 8, colored vertical lines indicate time steps with inconsistencies. These errors occur only in the training period for the dataset introduced in the Autoformer paper.

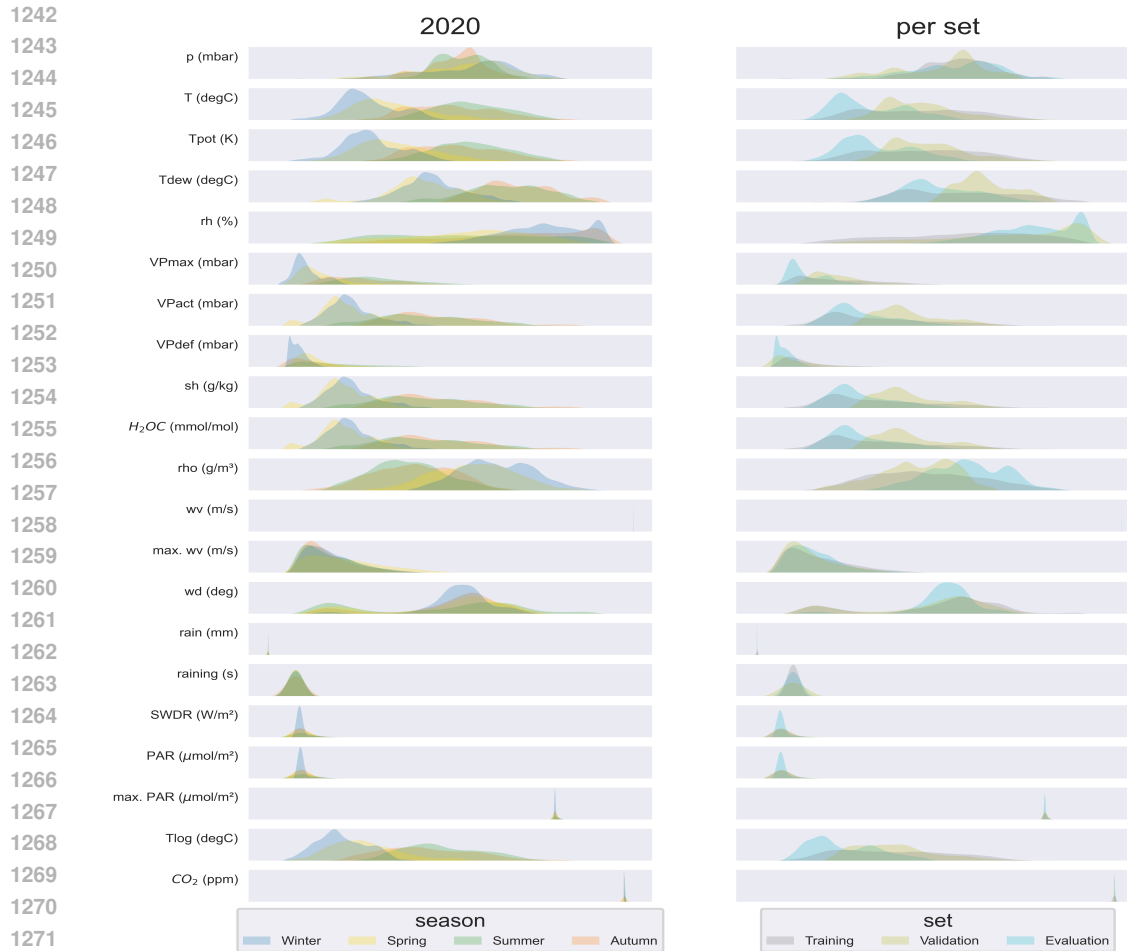


Figure 10: Weather Dataset from Autoformer - Distribution plots per channel. The last column illustrates data distribution with the ratio splitting strategy. The first column illustrates the data distribution for the whole datasets with a differentiation per season.

The **pink** vertical lines mark time steps where failure value appeared, while the **purple** lines denote missing time steps.

#### K.4 PROPOSED CORRECTION

To address these errors (failure values and missing time steps), we propose the following correction process as described in the main paper: (i) replace erroneous values with NaN, (ii) apply linear interpolation for isolated errors, and (iii) for consecutive errors, use context-aware when possible or linear interpolation.

The corrected dataset is visualized in Figure 11.

##### K.4.1 IDENTIFY INCONSISTENCIES

Five additional columns have been added to the CSV file in order to identify the time steps where inconsistencies were corrected:

- *is\_ww\_value\_error*: flags time steps where a failure value arose in *Wind Velocity*.
- *is\_maxPAR\_value\_error*: marks time steps where a failure value occurred in the *Maximum Photosynthetic Active Radiation* variable.



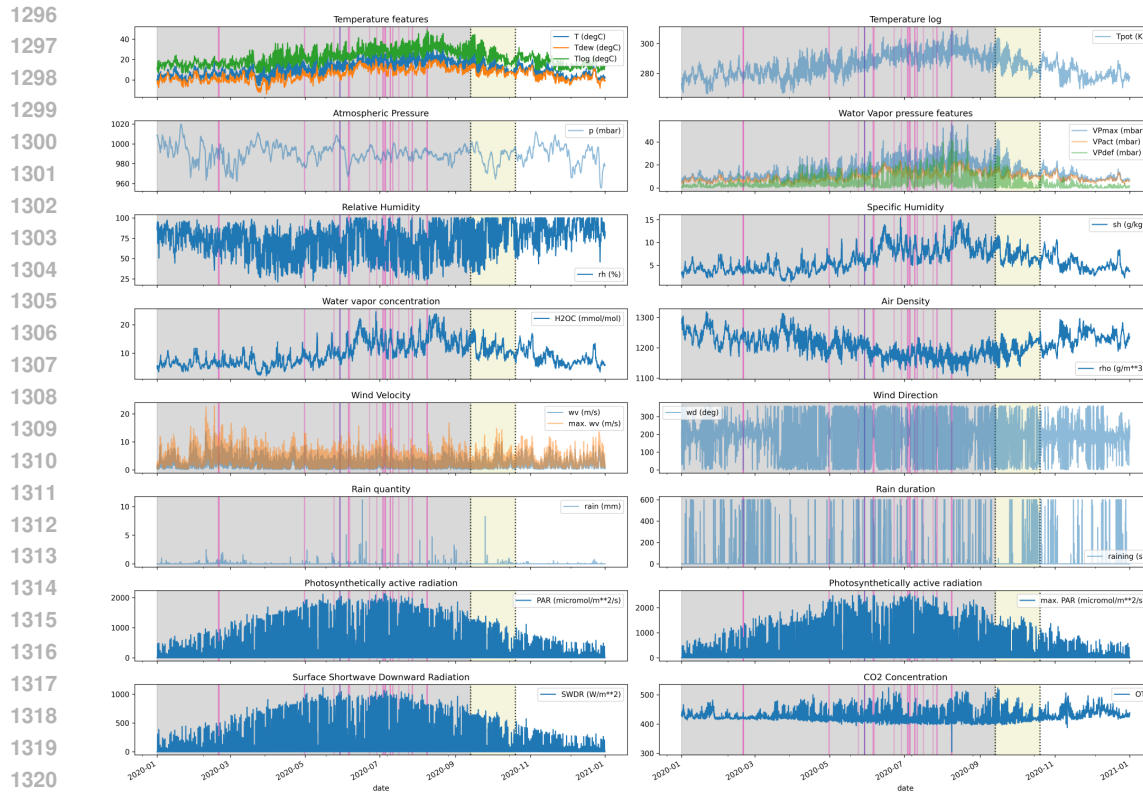


Figure 11: Overview of the weather indicators from the 1-year dataset used in Autoformer **after our correction process**. The gray [resp. yellow] background area denotes the training [resp. validation] period as defined in the ratio splitting. Colored vertical lines indicate time steps where inconsistencies were identified.

- *is\_OT\_value\_error*: identifies time steps where a failure value appeared in the *CO2 concentration* variable.
- *is\_ts\_missing*: indicates time steps that were missing in the original dataset.
- *is\_ts\_modified*: logs all time steps where corrections were applied.

#### K.4.2 OVERALL ANALYSIS

Figure 11 **shows** the plots of the different weather indicators **from the corrected version**. The gray area represents the training period, while the yellow area **indicates** the validation period as defined **by the ratio splitting strategy**. **The corrections appear to have effectively addressed the errors and inconsistencies.**

#### K.4.3 FREQUENCY ANALYSIS

**The frequency analysis of the revised dataset reveals slight differences from the original. While 10 channels still exhibit a dominant yearly frequency, the cycle now spans 52704 time steps (equivalent to 366 days), confirming that time steps were missing in the original dataset. The primary cycles in this dataset are now one year, six months, four months, and one day.**

#### K.4.4 CORRELATION ANALYSIS

**The correlation analysis for the corrected version closely resembles that of the original Autoformer dataset, with no significant deviations.**

	Fundamental	2 <sup>nd</sup>	3 <sup>rd</sup>
p (mbar)	10540.8 (73.20)	8784.0 (61.00)	4392.0 (30.50)
T (degC)	52704.0 (366.00)	144.0 (1.00)	26352.0 (183.00)
Tpot (K)	52704.0 (366.00)	144.0 (1.00)	26352.0 (183.00)
Tdew (degC)	52704.0 (366.00)	26352.0 (183.00)	17568.0 (122.00)
rh (%)	144.0 (1.00)	52704.0 (366.00)	17568.0 (122.00)
VPmax (mbar)	52704.0 (366.00)	144.0 (1.00)	26352.0 (183.00)
VPact (mbar)	52704.0 (366.00)	26352.0 (183.00)	10540.8 (73.20)
VPdef (mbar)	52704.0 (366.00)	144.0 (1.00)	143.6 (1.00)
sh (g/kg)	52704.0 (366.00)	26352.0 (183.00)	10540.8 (73.20)
H <sub>2</sub> OC ( $\mu$ mol/mol)	52704.0 (366.00)	26352.0 (183.00)	10540.8 (73.20)
rho (g/m <sup>3</sup> )	52704.0 (366.00)	144.0 (1.00)	26352.0 (183.00)
wv (m/s)	144.0 (1.00)	52704.0 (366.00)	8784.0 (61.00)
max. wv (m/s)	144.0 (1.00)	8784.0 (61.00)	2773.9 (19.26)
wd (deg)	8784.0 (61.00)	52704.0 (366.00)	3294.0 (22.88)
rain (mm)	2108.2 (14.64)	17568.0 (122.00)	258.4 (1.79)
raining (s)	17568.0 (122.00)	893.3 (6.20)	958.3 (6.65)
SWDR (W/m <sup>2</sup> )	144.0 (1.00)	52704.0 (366.00)	72.0 (0.50)
PAR ( $\mu$ mol/m <sup>2</sup> /s)	144.0 (1.00)	52704.0 (366.00)	72.0 (0.50)
max. PAR ( $\mu$ mol/m <sup>2</sup> /s)	144.0 (1.00)	52704.0 (366.00)	72.0 (0.50)
Tlog (degC)	52704.0 (366.00)	144.0 (1.00)	26352.0 (183.00)
CO <sub>2</sub> (ppm)	144.0 (1.00)	144.4 (1.00)	52704.0 (366.00)

Table 22: **MPIW\_10T\_1Y\_R** - Frequency analysis. The first value is the period in number of time steps the value in parentheses is the equivalent in days.

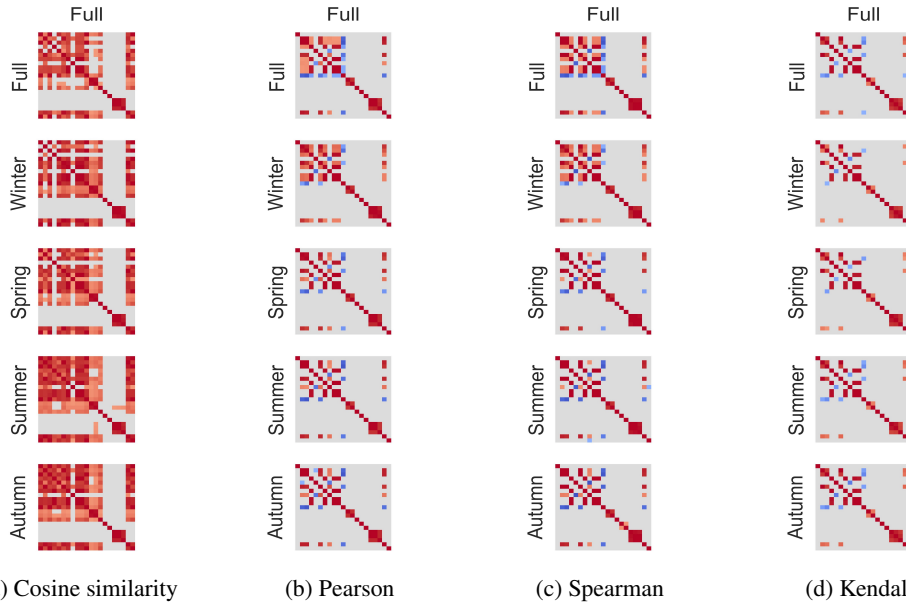


Figure 12: **MPIW\_10T\_1Y\_R** - Channels correlation for the full dataset and per season.

#### K.4.5 DATA DISTRIBUTION ANALYSIS

Figure 13 provides two distribution plots for our corrected version **MPIW\_10T\_1Y\_R**: one per season and one per splitting strategy set. As expected, the channel for which inconsistencies and especially failure values were uncovered now appears more consistent with the **other channels**. However, the distribution shift induced by the ratio splitting strategy **persists**.

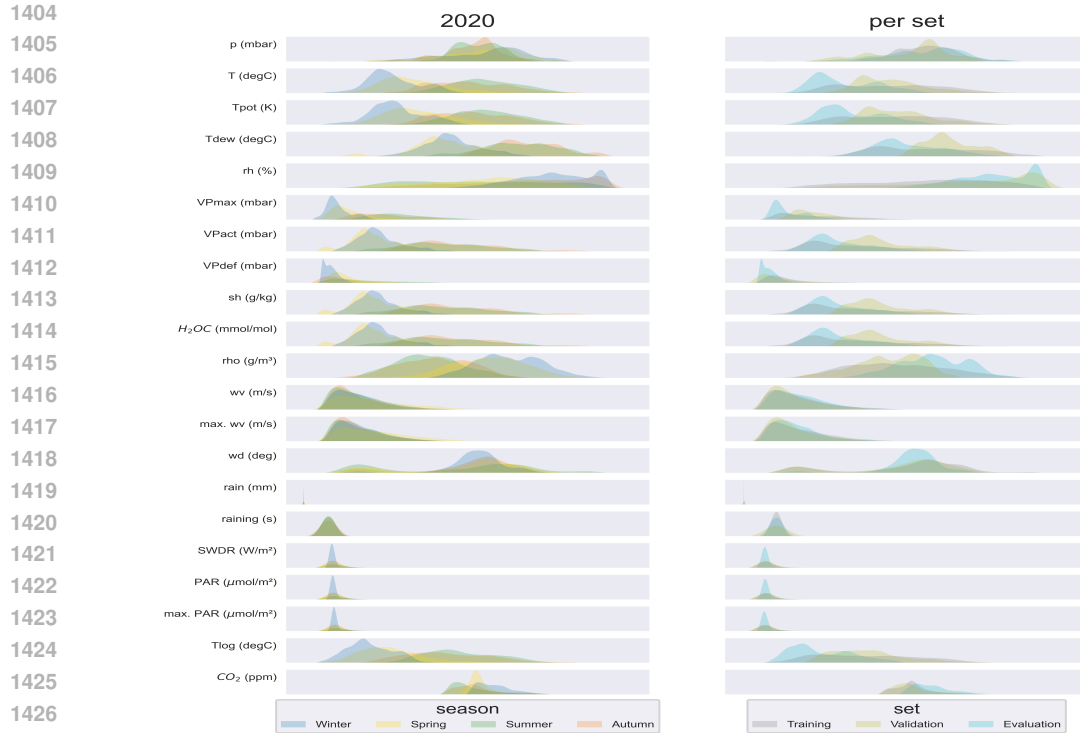


Figure 13: **MPIW\_10T\_1Y\_R** - Distribution plots per channel. The last column illustrates data distribution with the ratio splitting strategy. The first column illustrates the data distribution for the whole datasets with a differentiation per season.

## K.5 EXTENDED VERSIONS

To investigate cycle-inclusive splits, we extended the dataset to cover a 4-year period spanning from “2020-01-01 00:10:00” to “2024-01-01 00:10:00” (included). We collected additional data from the corresponding website and applied our correction process. The corrected dataset is shown in Figure 14. As illustrated, errors primarily appeared in the training and validation periods. However, due to our correction process, their impact should be minimal.

### K.5.1 OVERALL ANALYSIS

Figure 14 depicts the plots of the different weather indicators for the extended and corrected dataset. The gray area represents the training period, while the yellow area indicates the validation period as defined by the ratio splitting strategy. Errors and inconsistencies are no longer visible, suggesting that the corrections were applied successfully. This four-year dataset further confirms the presence of clear yearly cycles, as indicated by earlier analyses.

### K.5.2 FREQUENCY ANALYSIS

The frequency analysis of the extended dataset reveals differences from the one-year datasets. Now, 12 channels have a dominant yearly frequency and 7 channels have a dominant daily frequency. The most common cycles are one year, six months, and one day. As a result, the longest dominant cycle across all channels remains one year. However, it is now possible to use a cycle-inclusive splitting strategy that covers at least one full year.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483

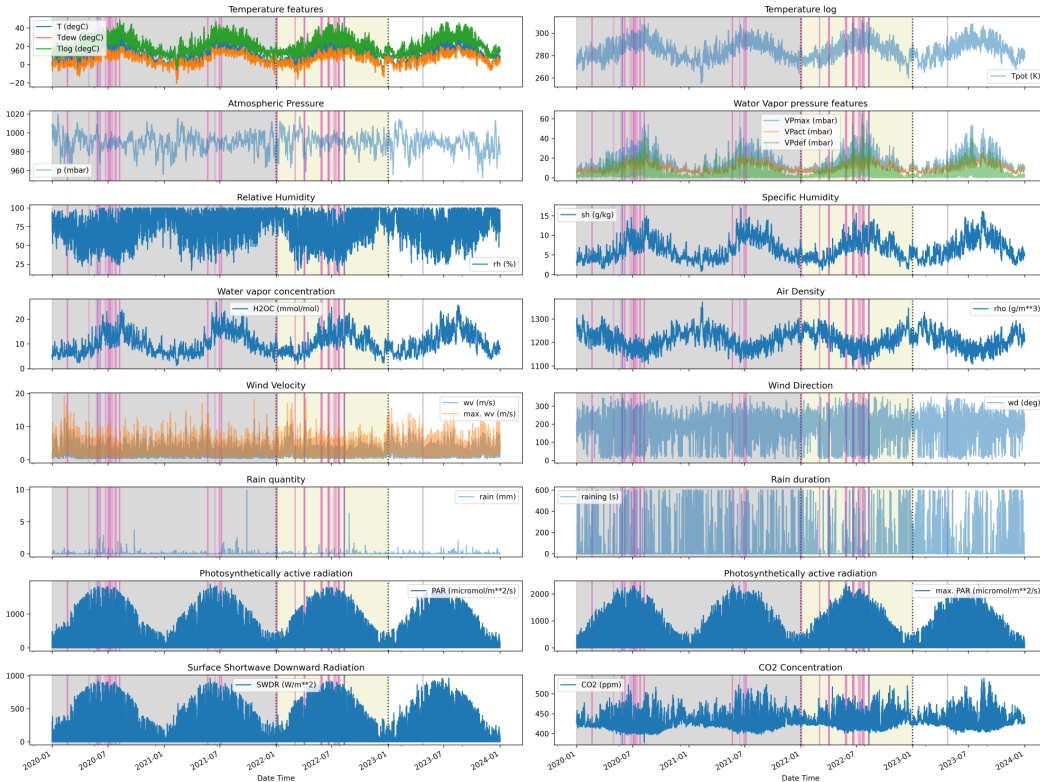


Figure 14: Overview of the weather indicators from our proposed 4-year dataset collected from MPI after our correction process. The gray background area represents the training period, while the yellow area denotes the validation period as defined in our proposed cycle-inclusive splitting. Colored vertical lines indicate time steps where inconsistencies were identified.

1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503

	Fundamental	2 <sup>nd</sup>	3 <sup>rd</sup>
p (mbar)	15027.4 (104.36)	16183.4 (112.38)	8091.7 (56.19)
T (degC)	52596.0 (365.25)	144.0 (1.00)	26298.0 (182.63)
Tpot (K)	52596.0 (365.25)	144.0 (1.00)	26298.0 (182.63)
Tdew (degC)	52596.0 (365.25)	8766.0 (60.88)	26298.0 (182.63)
rh (%)	144.0 (1.00)	52596.0 (365.25)	143.6 (1.00)
VPmax (mbar)	52596.0 (365.25)	144.0 (1.00)	26298.0 (182.63)
VPact (mbar)	52596.0 (365.25)	26298.0 (182.63)	8766.0 (60.88)
VPdef (mbar)	52596.0 (365.25)	144.0 (1.00)	143.6 (1.00)
sh (g/kg)	52596.0 (365.25)	26298.0 (182.63)	8766.0 (60.88)
H <sub>2</sub> OC ( $\mu$ mol/mol)	52596.0 (365.25)	26298.0 (182.63)	8766.0 (60.88)
rho ( $g/m^3$ )	52596.0 (365.25)	144.0 (1.00)	5686.1 (39.49)
wv (m/s)	144.0 (1.00)	52596.0 (365.25)	9562.9 (66.41)
max. wv (m/s)	144.0 (1.00)	52596.0 (365.25)	143.6 (1.00)
wd (deg)	52596.0 (365.25)	21038.4 (146.10)	144.0 (1.00)
rain (mm)	2805.1 (19.48)	52596.0 (365.25)	1290.7 (8.96)
raining (s)	52596.0 (365.25)	16183.4 (112.38)	2390.7 (16.60)
SWDR ( $W/m^2$ )	144.0 (1.00)	52596.0 (365.25)	143.6 (1.00)
PAR ( $\mu$ mol/m <sup>2</sup> /s)	144.0 (1.00)	52596.0 (365.25)	143.6 (1.00)
max. PAR ( $\mu$ mol/m <sup>2</sup> /s)	144.0 (1.00)	52704.0 (366.00)	72.0 (0.50)
Tlog (degC)	52596.0 (365.25)	144.0 (1.00)	144.4 (1.00)
CO <sub>2</sub> (ppm)	144.0 (1.00)	144.4 (1.00)	143.6 (1.00)

Table 23: MPIW\_10T\_4Y\_R - Frequency analysis. The first value is the period in number of time steps the value in parentheses is the equivalent in days.

1504  
1505

### K.5.3 CORRELATION ANALYSIS

1506  
1507  
1508  
1509  
1510  
1511

Figure 15 displays the channel correlations for the extended dataset MPIW\_10T\_4Y\_R using the different methods mentioned in Appendix I.3. Similarly to LCD, for all metrics, the following patterns emerge:

1. By row: Year-to-year correlations remain consistent (with minimal variation);
2. By column: Within a given period, when divided by solar seasons, the correlations can vary significantly. For instance, *Winter* and *Spring* exhibit notable differences compared to *Summer* and *Autumn*. In addition, while differences between *Winter* and *Spring*, as well as *Summer* and *Autumn*, are less pronounced, they are still evident.

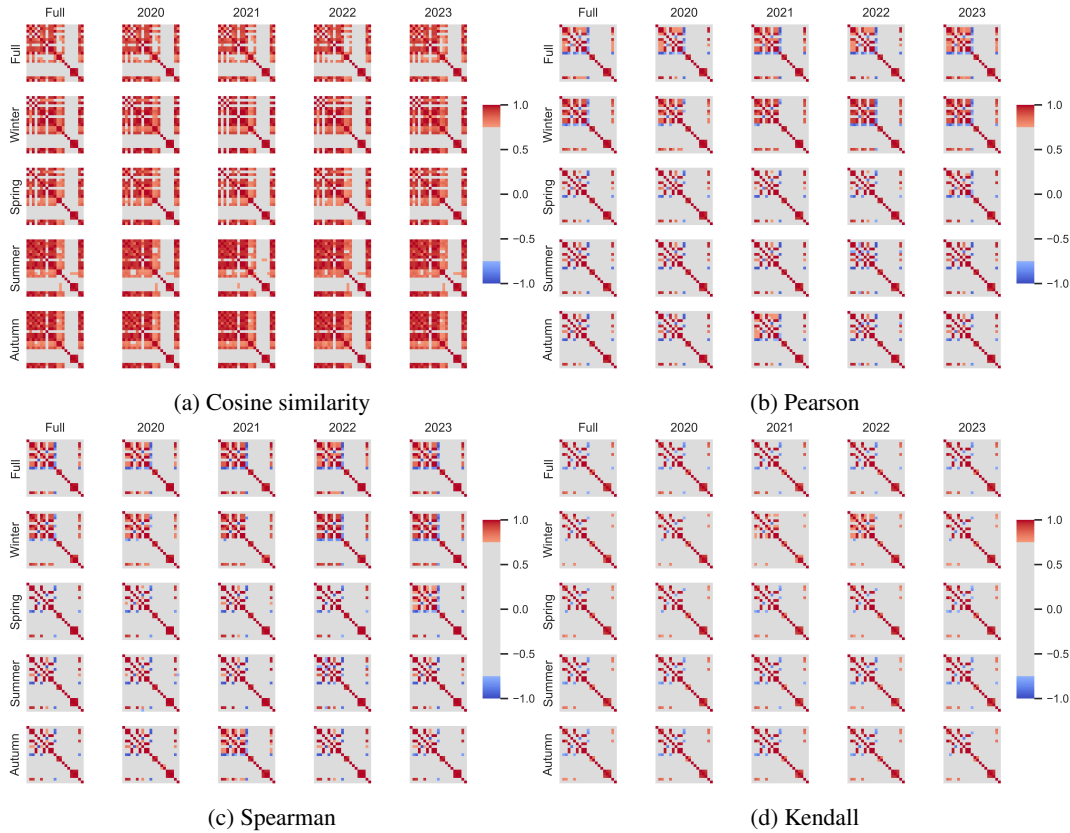


Figure 15: **MPIW\_10T\_4Y\_R** - Channels correlation for the full dataset, per year and per season.

#### K.5.4 DATA DISTRIBUTION ANALYSIS

Figure 16 provides various distribution plots for the corrected four-year dataset: **MPIW\_10T\_4Y\_R**. Some inter-annual variations are observed, such as differences in *Relative Humidity (rh)* densities between 2021 and 2022 compared to 2020 and 2024. Any efficient MTSF models should account for such variations in order to be considered robust.

In addition, we observed in Figure 16 that our cycle-inclusive splitting strategy significantly reduces distribution shift across sets, ensuring that model performances are evaluated over the longest cycle period.

#### K.5.5 IDENTIFY INCONSISTENCIES:

Six additional columns have been appended to the produced CSV files in order to identify the time steps where inconsistencies were corrected:

- *is\_wv\_value\_error*: marks time steps where a failure value appeared in the *Wind Velocity* variable.
- *is\_SWDR\_value\_error*: highlights time steps where a failure value occurred in the *Surface Shortwave Downward Radiation* variable.

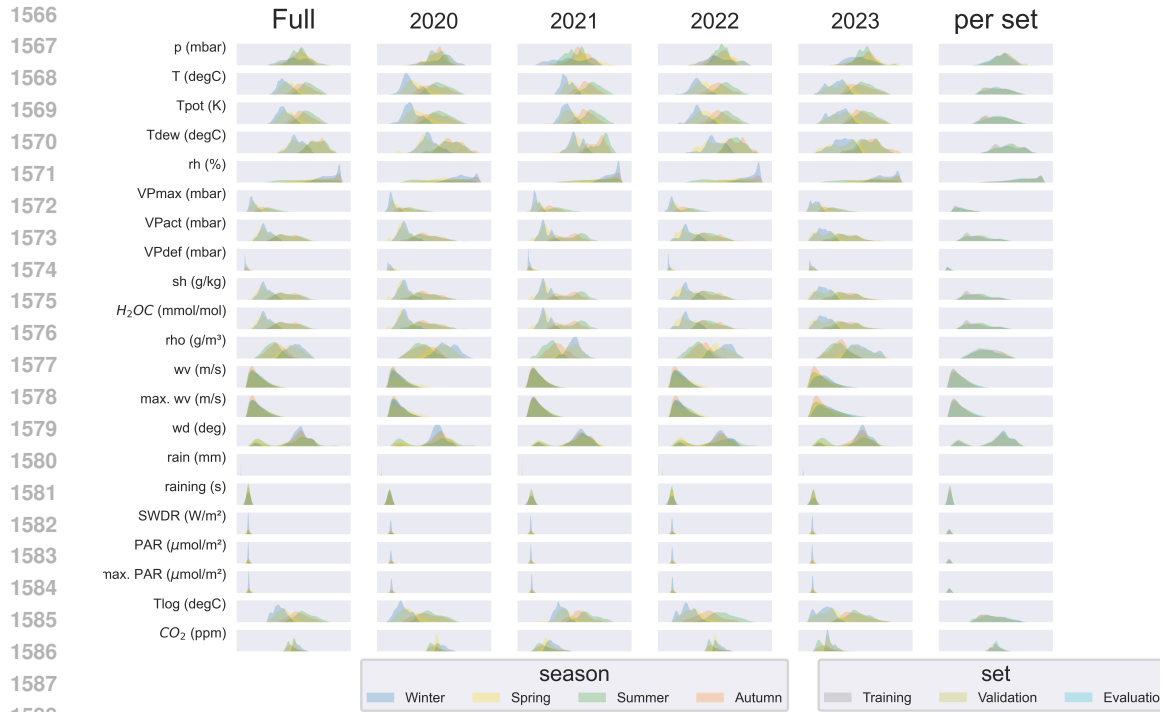


Figure 16: **MPIW\_10T\_4Y\_R** - Distribution plots per channel. The last column illustrates data distribution with the ratio splitting strategy. The first column illustrates the data distribution for the whole datasets with a differentiation per season.

- *is\_maxPAR\_value\_error*: flags time steps where a failure value appeared in the *Maximum Photosynthetic Active Radiation* variable.
- *is\_CO2\_value\_error*: identifies time steps where a failure value occurred in the *CO2 concentration* variable.
- *is\_ts\_missing*: indicates time steps that were missing in the original dataset.
- *is\_ts\_modified*: logs all time steps where corrections were applied.

#### K.5.6 HOURLY VERSION:

We propose an hourly version of this 4-year dataset by aggregating data over six consecutive time steps (i.e., from HH:10 to HH+1:00). The following aggregation functions are applied to the corresponding variables:

- **Sum**: *Precipitation Amount*, *Precipitation Duration*, *Surface Shortwave Downward Radiation*, *Photosynthetic Active Radiation* and columns identifying the errors.
- **Maximum**: *Maximum Photosynthetic Active Radiation* and *Maximum Wind Velocity*.
- **Mean**: All other variables.

**Frequency Analysis:** As shown in Table 24, the frequency analysis of the hourly dataset is similar to that of the 10-minute interval dataset.

**Correlation Analysis:** As shown in Figure 17, the correlation analysis of the hourly dataset is similar to that of the 10-minute interval dataset.

The data distribution of the hourly dataset being very similar to the 10-minute interval dataset, the corresponding plots have been omitted.

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

	Fundamental	2 <sup>nd</sup>	3 <sup>rd</sup>
p (mbar)	2504.6 (104.36)	2697.2 (112.38)	1348.6 (56.19)
T (degC)	8766.0 (365.25)	24.0 (1.00)	4383.0 (182.62)
Tpot (K)	8766.0 (365.25)	24.0 (1.00)	4383.0 (182.62)
Tdew (degC)	8766.0 (365.25)	1461.0 (60.88)	4383.0 (182.62)
rh (%)	24.0 (1.00)	8766.0 (365.25)	23.9 (1.00)
VPmax (mbar)	8766.0 (365.25)	24.0 (1.00)	4383.0 (182.62)
VPact (mbar)	8766.0 (365.25)	4383.0 (182.62)	1461.0 (60.88)
VPdef (mbar)	8766.0 (365.25)	24.0 (1.00)	23.9 (1.00)
sh (g/kg)	8766.0 (365.25)	4383.0 (182.62)	1461.0 (60.88)
H <sub>2</sub> OC ( $\mu\text{mol/mol}$ )	8766.0 (365.25)	4383.0 (182.62)	1461.0 (60.88)
rho ( $\text{g/m}^3$ )	8766.0 (365.25)	24.0 (1.00)	947.7 (39.49)
wv (m/s)	24.0 (1.00)	8766.0 (365.25)	1593.8 (66.41)
max. wv (m/s)	24.0 (1.00)	8766.0 (365.25)	23.9 (1.00)
wd (deg)	8766.0 (365.25)	3506.4 (146.10)	24.0 (1.00)
rain (mm)	467.5 (19.48)	8766.0 (365.25)	215.1 (8.96)
raining (s)	8766.0 (365.25)	2697.2 (112.38)	398.5 (16.60)
SWDR ( $\text{W/m}^2$ )	24.0 (1.00)	8766.0 (365.25)	23.9 (1.00)
PAR ( $\mu\text{mol/m}^2/\text{s}$ )	24.0 (1.00)	8766.0 (365.25)	23.9 (1.00)
max. PAR ( $\mu\text{mol/m}^2/\text{s}$ )	24.0 (1.00)	8766.0 (365.25)	23.9 (1.00)
Tlog (degC)	8766.0 (365.25)	24.0 (1.00)	24.1 (1.00)
CO <sub>2</sub> (ppm)	24.0 (1.00)	24.1 (1.00)	23.9 (1.00)

Table 24: MPIW\_1H\_4Y\_R - Frequency analysis. The first value is the period in number of time steps the value in parentheses is the equivalent in days.

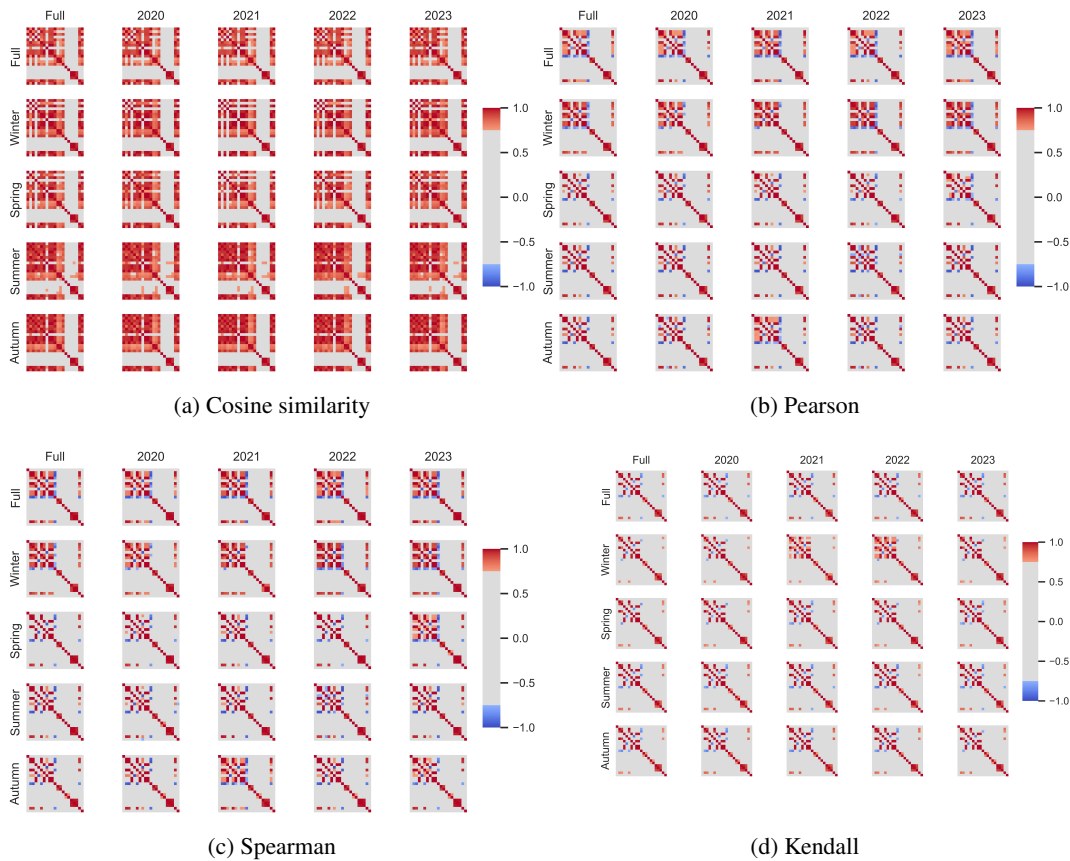
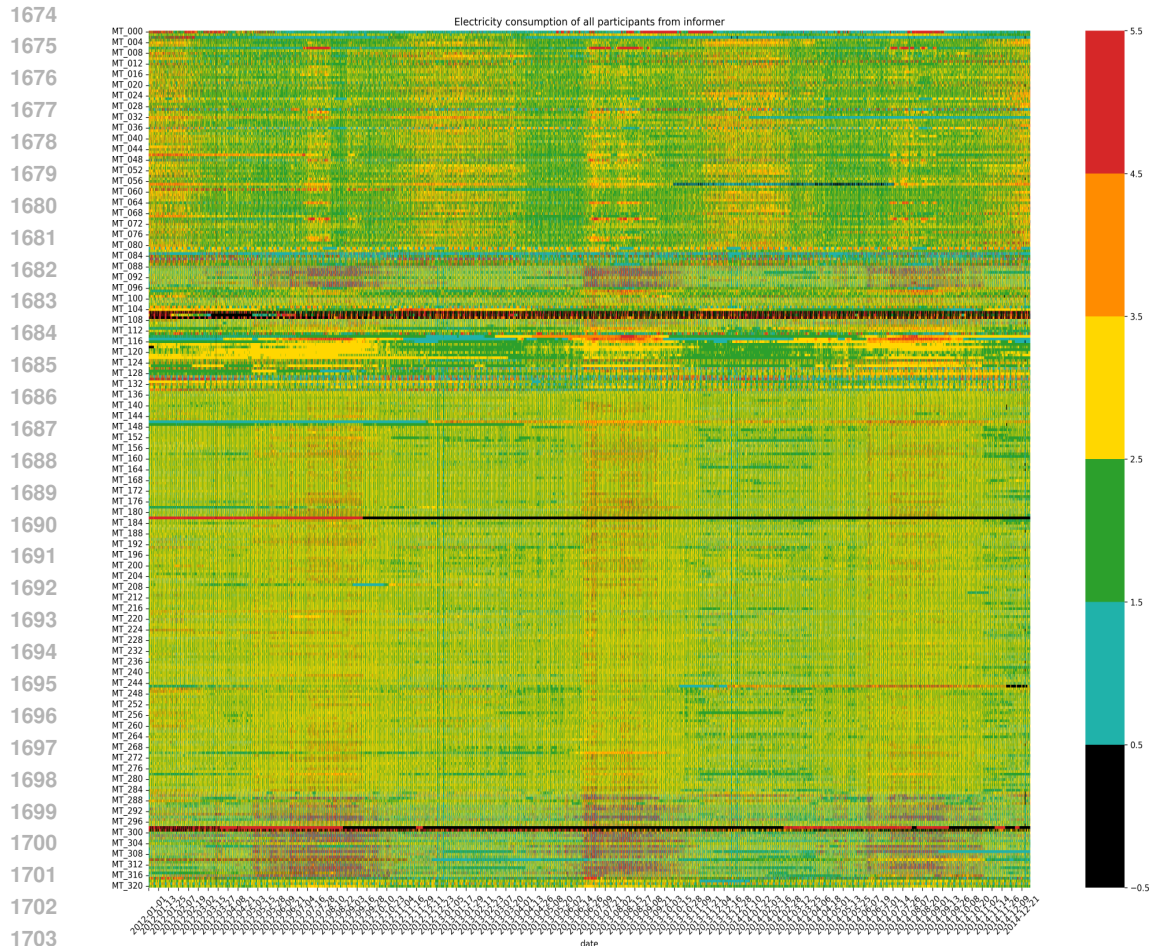


Figure 17: MPIW\_1H\_4Y\_R - Channels correlation for the full dataset, per year and per season.



1705 Figure 18: Overview of the normalized electricity consumption patterns of clients from the ECL  
 1706 dataset (derived from the UCI ELD dataset). The heatmap visualization simplifies the identification  
 1707 of inconsistent consumption patterns among clients.

## 1709 L ELECTRICITY LOAD DIAGRAMS DATASET

### 1711 L.1 DESCRIPTION

1713 The ELD<sup>10</sup> dataset consists of the electricity consumption data of 370 clients from what it appears  
 1714 to be a Portuguese electricity provider as timestamps report to Portuguese hours. Measurements  
 1715 were originally recorded every 15 minutes. The raw dataset covers the period from “2011-01-01  
 1716 00:15:00” to “2015-01-01 00:00:00” (included). By aggregating four consecutive measurements  
 1717 (i.e., HH:15, HH:30, HH:45 and HH+1:00, an hourly version of the dataset can be obtained. Al-  
 1718 though the dataset description in UCI indicates having no missing data, some profiles depicted long  
 1719 and constant consumption equal to zero, as shown in the following sections, probably suggesting  
 1720 late arrival or early departure when occurring at the beginning or the end of the covered period,  
 1721 respectively.

### 1722 L.2 ANALYSIS

1724 We consider the ELD dataset as spatiotemporal, where each channel represents the electricity con-  
 1725 sumption of clients across different locations in Portugal. These clients may belong to various  
 1726 categories such as *Residential*, *Commercial*, or *Industrial*, resulting in diverse consumption patterns

1727 <sup>10</sup><https://archive.ics.uci.edu/dataset/321/electricityloadaddiagrams20112014>



and variation in volume, as evidenced in this document. While the dataset lacks specific information on the location and type of clients, it presents a rich tapestry of cycles closely tied to date, time, and human behavior. In addition, the variability in consumption patterns among clients poses a significant challenge for models, especially without external information, requiring them to decipher these underlying characteristics and correlations to accurately predict electricity consumption. Overall, predicting electricity consumption with this dataset presents a challenging task.

### L.3 ORIGINAL DATASET

ECL is an hourly dataset first introduced by (Li et al., 2019), derived from the ELD dataset available on UCI. This dataset provides electricity consumption data from 321 clients in Portugal, each identified as “MT.XXX”, with ‘XXX’ representing a unique identifier.

All timestamps report to Portuguese hours. The dataset covers the period from “2012-01-01 00:00:00” to “2014-12-31 23:00:00” (included).

#### L.3.1 OVERALL ANALYSIS

Figure 18 plots the normalized consumption of the considered clients as a heatmap, aiding in the identification of distinctive patterns. These include clients with constant consumption values over time or those with unusual consumption patterns not typically observed in electricity usage. This figure reveals that most clients exhibit similar patterns, with noticeable summer peaks recurring annually in the bottom section of the figure. Conversely, clients in the upper section depict less pronounced peaks.

Notably, certain clients exhibit anomalies, such as the client displaying a continuous period of zero consumption (indicated by a black region).

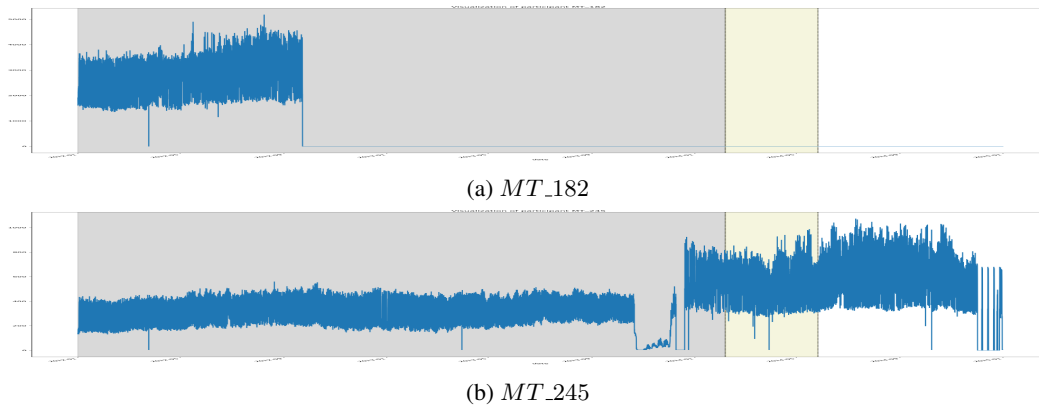


Figure 19: Overview of the electricity consumption profiles of two clients showing “early departure”. The gray background area represents the training period, while the yellow area represents the validation period as defined in the ratio splitting. *MT\_245* also exhibits sudden changes in consumption patterns.

#### L.3.2 INCONSISTENCIES PRESENTATION

In the following figures, the gray [resp. yellow] area represents the training [resp. validation] period as defined in the ratio splitting. In the raw UCI dataset, clients who began participating after the dataset’s starting date showed constant consumption equal to zero before their participation started. These clients, that we refer to as “late arrival” clients, were removed in the ECL dataset version. However, as shown in Figure 19, two clients in the ECL dataset (particularly *MT\_182*) exhibit prolonged zero consumption after a certain date, suggesting an “early departure”. We believe that these clients should have likely been removed as well to avoid impacting model evaluation in MTS forecasting.

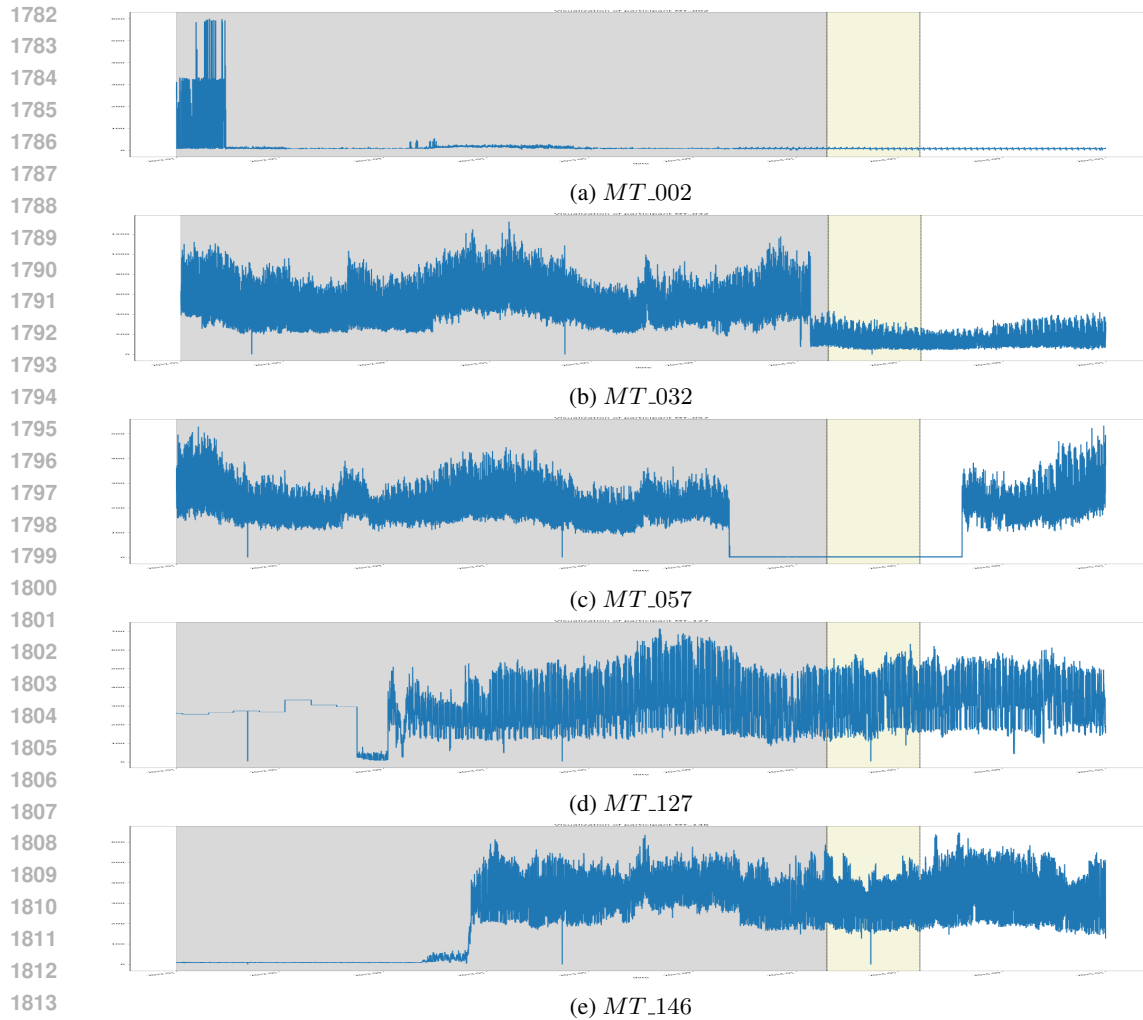


Figure 20: Overview of the electricity consumption profiles of five clients displaying sudden changes in their overall patterns. The gray background area represents the training period, while the yellow area represents the validation period as defined in the ratio splitting. *MT\_127* and *MT\_146* also exhibit unusual consumption patterns at the beginning of the monitored period.

In addition, our analysis unveiled some clients with unusual and significant changes in their consumption patterns, such as the one shown in Figure 20. Without external information explaining such sudden changes, it becomes challenging for models to accurately learn consumption patterns and potential inter-variable relations.

Finally, similar to other datasets, we believe that the ratio splitting may not be optimal for conducting a fair comparison between models. This approach may favor models that perform well in the evaluation period but could potentially perform poorly elsewhere.

1836 L.3.3 FREQUENCY ANALYSIS  
1837

	Fundamental	2 <sup>nd</sup>	3 <sup>rd</sup>	
1838				
1839	MT_001	12.0 (0.50)	24.0 (1.00)	8768.0 (365.33)
1840	MT_002	26304.0 (1096.00)	8768.0 (365.33)	13152.0 (548.00)
1841	MT_003	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
1842	MT_004	8768.0 (365.33)	24.0 (1.00)	4384.0 (182.67)
1843	MT_005	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
1844	MT_006	4384.0 (182.67)	8768.0 (365.33)	2922.7 (121.78)
1845	MT_007	24.0 (1.00)	12.0 (0.50)	8.0 (0.33)
1846	MT_008	12.0 (0.50)	8768.0 (365.33)	24.0 (1.00)
1847	MT_009	24.0 (1.00)	8768.0 (365.33)	84.0 (3.50)
1848	MT_010	24.0 (1.00)	8768.0 (365.33)	12.0 (0.50)
1849	MT_011	24.0 (1.00)	167.5 (6.98)	84.0 (3.50)
1850	MT_012	24.0 (1.00)	8768.0 (365.33)	167.5 (6.98)
1851	MT_013	24.0 (1.00)	12.0 (0.50)	4384.0 (182.67)
1852	MT_014	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
1853	MT_015	24.0 (1.00)	12.0 (0.50)	4384.0 (182.67)
1854	MT_016	24.0 (1.00)	12.0 (0.50)	4384.0 (182.67)
1855	MT_017	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
1856	MT_018	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
1857	MT_019	24.0 (1.00)	12.0 (0.50)	8.0 (0.33)
1858	MT_020	12.0 (0.50)	24.0 (1.00)	6.0 (0.25)
1859	MT_021	12.0 (0.50)	24.0 (1.00)	8768.0 (365.33)
1860	MT_022	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
1861	MT_023	24.0 (1.00)	8768.0 (365.33)	12.0 (0.50)
1862	MT_024	24.0 (1.00)	8768.0 (365.33)	12.0 (0.50)
1863	MT_025	167.5 (6.98)	84.0 (3.50)	168.6 (7.03)
1864	MT_026	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
1865	MT_027	24.0 (1.00)	12.0 (0.50)	8.0 (0.33)
1866	MT_028	12.0 (0.50)	24.0 (1.00)	4384.0 (182.67)

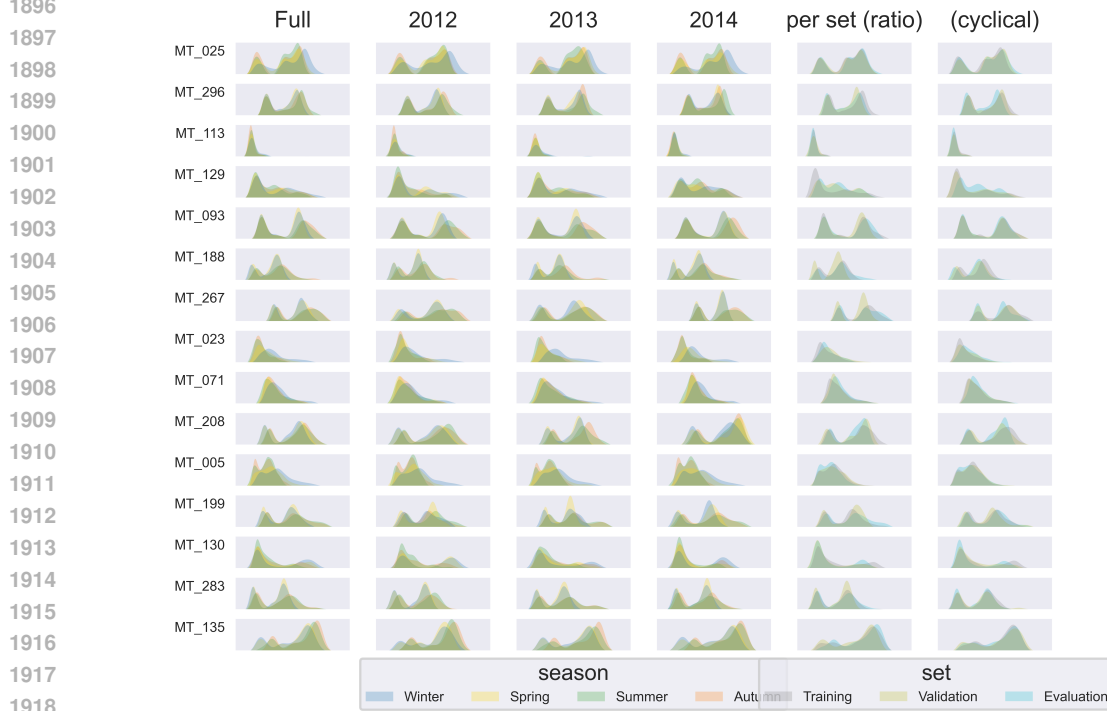
1867 Table 25: ECL Dataset - Frequency analysis of the first channels. The first value is the period in  
1868 number of time steps the value in parentheses is the equivalent in days.

1869 From Table 25 we can observed that some channels have their dominant periods significantly over  
1870 one year. But the majority exhibit a longest cycle of one year.  
1871

1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889

### 1890 L.3.4 DATA DISTRIBUTION ANALYSIS

1891  
1892 Figure 21 provides the distribution plots for the *ECL* dataset, revealing that channels are sensitive to  
1893 seasonal variations.  
1894



1920 Figure 21: ECL Dataset - Distribution plots per channel. The last two columns illustrate data distri-  
1921 bution per splitting strategy: ratio and our proposal cycle-inclusive. The other columns illustrate the  
1922 data distribution for the whole datasets and per year, with a differentiation per season.  
1923

### 1924 L.4 PROPOSED CORRECTION

1925  
1926 Based on our observations, we propose removing the following 13 clients from the ECL dataset:

- 1927
- 1928 • **Early departure:** *MT\_182* and *MT\_245*
- 1929 • **Significant changes in consumption patterns:** *MT\_032*, *MT\_057*, *MT\_127*, *MT\_146*  
1930 and *MT\_307*
- 1931 • **No clear cyclical patterns:** *MT\_002*, *MT\_106*, *MT\_114*, *MT\_122*, *MT\_298* and  
1932 *MT\_310*  
1933

1934 The overall visualization of our proposed dataset is depicted in Figure 22.  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

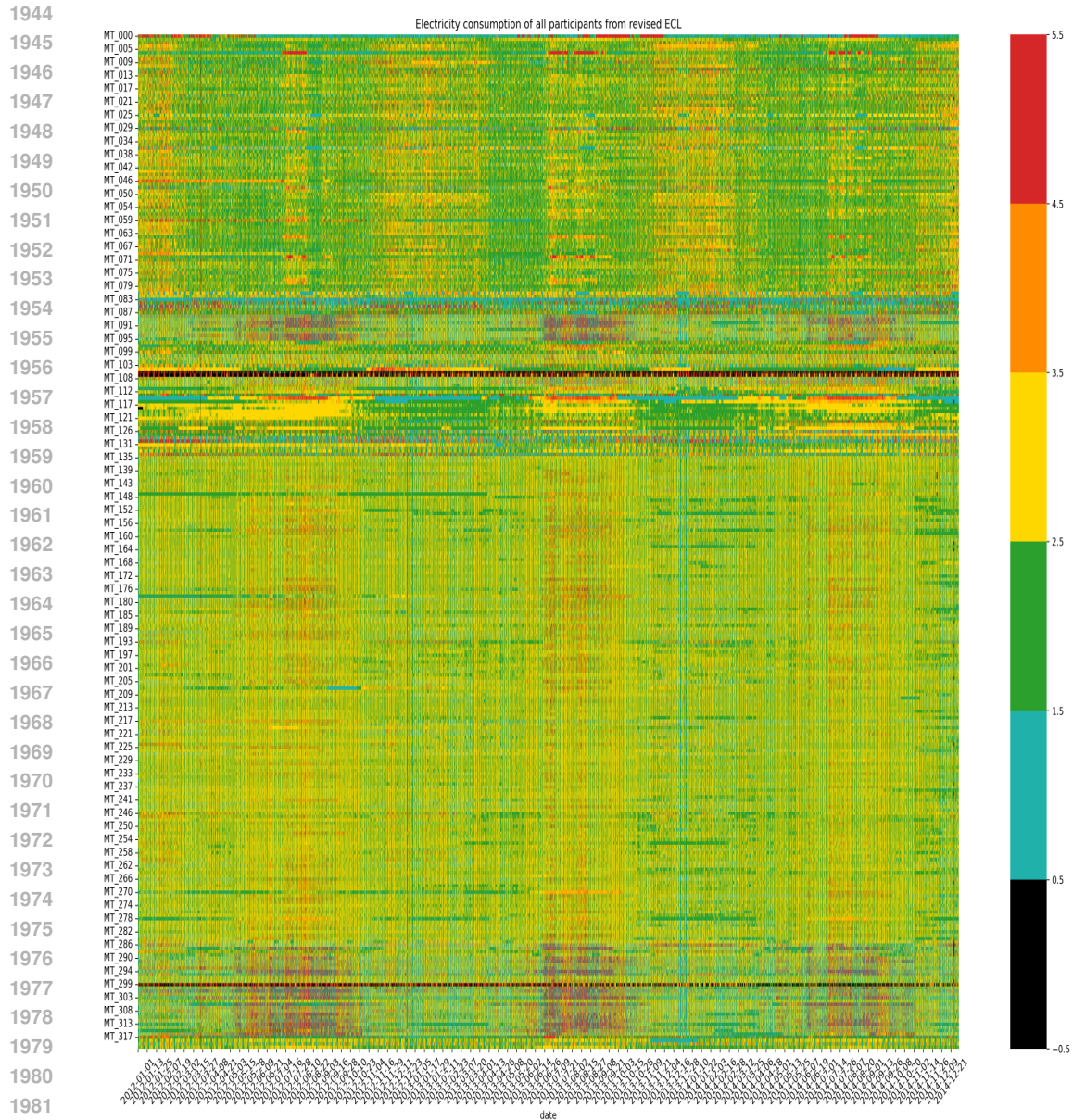


Figure 22: Overview of the normalized electricity consumption patterns of clients from our revised version of ECL dataset. The heatmap visualization simplifies the identification of inconsistent consumption patterns among clients.

#### L.4.1 FREQUENCY ANALYSIS

From Table 26 we can observed that some channels have their dominant periods significantly over one year, but also significantly less than with the ECL dataset. But the majority exhibit a longest cycle of one year.

	Fundamental	2 <sup>nd</sup>	3 <sup>rd</sup>	
1998				
1999	MT_000	13152.0 (548.00)	6576.0 (274.00)	3757.7 (156.57)
2000	MT_001	12.0 (0.50)	24.0 (1.00)	8768.0 (365.33)
2001	MT_003	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
2002	MT_004	8768.0 (365.33)	24.0 (1.00)	4384.0 (182.67)
2003	MT_005	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
2004	MT_006	4384.0 (182.67)	8768.0 (365.33)	2922.7 (121.78)
2005	MT_007	24.0 (1.00)	12.0 (0.50)	8.0 (0.33)
2006	MT_008	12.0 (0.50)	8768.0 (365.33)	24.0 (1.00)
2007	MT_009	24.0 (1.00)	8768.0 (365.33)	84.0 (3.50)
2008	MT_010	24.0 (1.00)	8768.0 (365.33)	12.0 (0.50)
2009	MT_011	24.0 (1.00)	167.5 (6.98)	84.0 (3.50)
2010	MT_012	24.0 (1.00)	8768.0 (365.33)	167.5 (6.98)
2011	MT_013	24.0 (1.00)	12.0 (0.50)	4384.0 (182.67)
2012	MT_014	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
2013	MT_015	24.0 (1.00)	12.0 (0.50)	4384.0 (182.67)
2014	MT_016	24.0 (1.00)	12.0 (0.50)	4384.0 (182.67)
2015	MT_017	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
2016	MT_018	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
2017	MT_019	24.0 (1.00)	12.0 (0.50)	8.0 (0.33)
2018	MT_020	12.0 (0.50)	24.0 (1.00)	6.0 (0.25)
2019	MT_021	12.0 (0.50)	24.0 (1.00)	8768.0 (365.33)
2020	MT_022	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
2021	MT_023	24.0 (1.00)	8768.0 (365.33)	12.0 (0.50)
2022	MT_024	24.0 (1.00)	8768.0 (365.33)	12.0 (0.50)
2023	MT_025	167.5 (6.98)	84.0 (3.50)	168.6 (7.03)
2024	MT_026	24.0 (1.00)	12.0 (0.50)	8768.0 (365.33)
2025	MT_027	24.0 (1.00)	12.0 (0.50)	8.0 (0.33)
2026	MT_028	12.0 (0.50)	24.0 (1.00)	4384.0 (182.67)

Table 26: **PELD\_1H\_3Y\_308** - Frequency analysis of the first channels. The first value is the period in number of time steps the value in parentheses is the equivalent in days.

#### L.4.2 DATA DISTRIBUTION ANALYSIS

Figure 23 provides the same distribution plots for the revised dataset: **PELD\_1H\_3Y\_308**. The modified inconsistencies and errors did not alter the properties of the datasets. Data distribution vary significantly per season, but our cycle-inclusive strategy ensure better distribution similarity between sets, making such dataset more suitable for benchmarking.

2052

2053

2054

2055

2056

2057

2058

2059

2060

2061

2062

2063

2064

2065

2066

2067

2068

2069

2070

2071

2072

2073

2074

2075

2076

2077

2078

2079

2080

2081

2082

2083

2084

2085

2086

2087

2088

2089

2090

2091

2092

2093

2094

2095

2096

2097

2098

2099

2100

2101

2102

2103

2104

2105

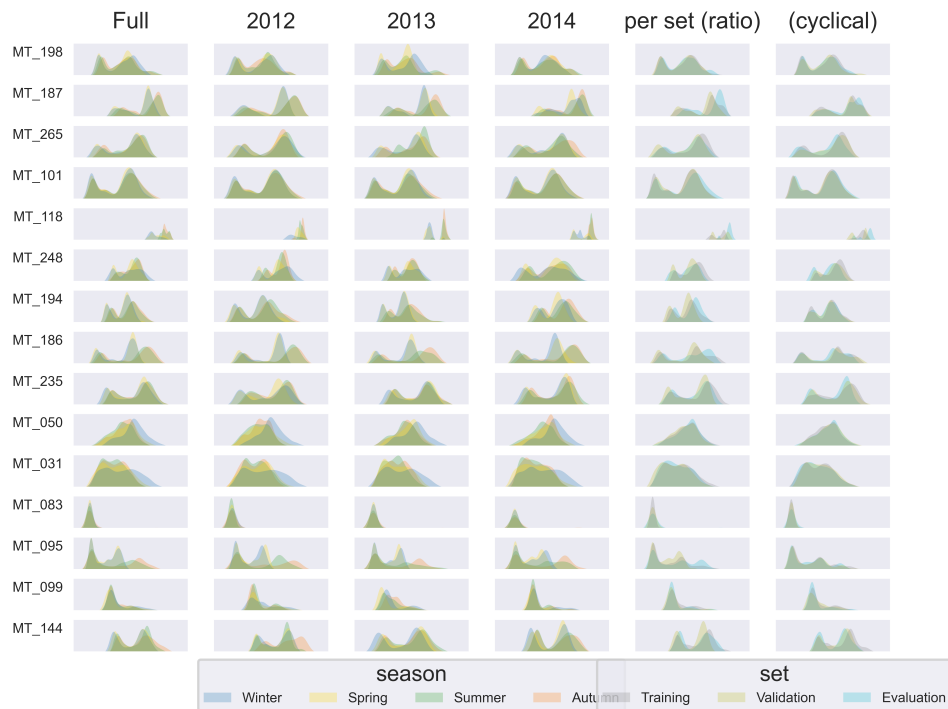


Figure 23: **PELD\_1H\_3Y\_308** - Distribution plots per channel. The last two columns illustrate data distribution per splitting strategy: ratio and our proposal cycle-inclusive. The other columns illustrate the data distribution for the whole datasets and per year, with a differentiation per season.

## L.5 FUTURE VERSION

In the future, it may be necessary to remove or better identify clients exhibiting “short” periods of unusual consumption patterns or specific trends (either upward or downward consumption trends over years). This approach would allow for the segmentation of typical metrics (MAE, MSE, etc.) into three categories: an overall metric, metrics for clients with “usual” cyclical patterns, and metrics specifically for clients with these specific characteristics. Such a categorization would provide a clearer understanding of model performance and enable researchers to refine architectures more effectively.

## REFERENCES

- Razvan-Gabriel Cirstea, Chenjuan Guo, Bin Yang, Tung Kieu, Xuanyi Dong, and Shirui Pan. Triformer: Triangular, Variable-Specific Attentions for Long Sequence Multivariate Time Series Forecasting. In Lud De Raedt (ed.), *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 1994–2001. International Joint Conferences on Artificial Intelligence Organization, 7 2022. doi: 10.24963/ijcai.2022/277. URL <https://doi.org/10.24963/ijcai.2022/277>. Main Track.
- Dazhao Du, Bing Su, and Zhewei Wei. Preformer: Predictive Transformer with Multi-Scale Segment-Wise Correlations for Long-Term Time Series Forecasting. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5, 2023. doi: 10.1109/ICASSP49357.2023.10096881.
- Lu Han, Han-Jia Ye, and De-Chuan Zhan. The Capacity and Robustness Trade-Off: Revisiting the Channel Independent Strategy for Multivariate Time Series Forecasting. *IEEE Trans. on Knowl. and Data Eng.*, 36(11):7129–7142, May 2024.

- 2106 Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Re-  
2107 versible Instance Normalization for Accurate Time-Series Forecasting against Distribution Shift.  
2108 In *Proceedings of the Tenth International Conference on Learning Representations (ICLR)*, 2022.  
2109 URL <https://openreview.net/forum?id=cGDakQo1C0p>.  
2110
- 2111 Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyong Zhou, Wenhui Chen, Yu-Xiang Wang, and Xifeng Yan.  
2112 Enhancing the Locality and Breaking the Memory Bottleneck of Transformer on Time Series  
2113 Forecasting. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garn-  
2114 nett (eds.), *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32. Cur-  
2115 ran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper\\_files/  
2116 paper/2019/file/6775a0635c302542da2c32aa19d86be0-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/6775a0635c302542da2c32aa19d86be0-Paper.pdf).
- 2117 Zhe Li, Shiyi Qi, Yiduo Li, and Zenglin Xu. Revisiting Long-term Time Series Forecasting: An  
2118 Investigation on Affine Mapping. In *Submitted to The Twelfth International Conference on  
2119 Learning Representations (ICLR)*, 2024. URL [https://openreview.net/forum?id=  
2120 T97kxctihq](https://openreview.net/forum?id=T97kxctihq). under review.
- 2121 Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X. Liu, and Schahram Dustdar.  
2122 Pyraformer: Low-Complexity Pyramidal Attention for Long-Range Time Series Modeling and  
2123 Forecasting. In *Proceedings of the Tenth International Conference on Learning Representations  
2124 (ICLR)*, 2022. URL <https://openreview.net/forum?id=0EXmFzUn5I>.  
2125
- 2126 Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long.  
2127 iTransformer: Inverted Transformers Are Effective for Time Series Forecasting. In *Submitted to  
2128 The Twelfth International Conference on Learning Representations (ICLR)*, 2024. URL [https://  
2129 openreview.net/forum?id=JePfAI8fah](https://openreview.net/forum?id=JePfAI8fah). under review.
- 2130 Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A Time Series is Worth 64  
2131 Words: Long-term Forecasting with Transformers. In *Proceedings of the Eleventh International  
2132 Conference on Learning Representations (ICLR)*, 2023. URL [https://openreview.net/  
2133 forum?id=Jbdc0vT0col](https://openreview.net/forum?id=Jbdc0vT0col).  
2134
- 2135 Gerald Woo, Chenghao Liu, Doyen Sahoo, Akshat Kumar, and Steven Hoi. ETSformer: Exponential  
2136 Smoothing Transformers for Time-series Forecasting, 2023. URL [https://openreview.  
2137 net/forum?id=5m\\_3whfo483](https://openreview.net/forum?id=5m_3whfo483).
- 2138 Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition Trans-  
2139 formers with Auto-Correlation for Long-Term Series Forecasting. In M. Ranzato, A. Beygelz-  
2140 imer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Informa-  
2141 tion Processing Systems (NeurIPS)*, volume 34, pp. 22419–22430. Curran Associates, Inc.,  
2142 2021. URL [https://proceedings.neurips.cc/paper\\_files/paper/2021/  
2143 file/bcc0d400288793e8bdcd7c19a8ac0c2b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/bcc0d400288793e8bdcd7c19a8ac0c2b-Paper.pdf).
- 2144 Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Times-  
2145 Net: Temporal 2D-Variation Modeling for General Time Series Analysis. In *Proceedings  
2146 of the Eleventh International Conference on Learning Representations (ICLR)*, 2023. URL  
2147 [https://openreview.net/forum?id=ju\\_Uqw3840q](https://openreview.net/forum?id=ju_Uqw3840q).  
2148
- 2149 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are Transformers Effective for Time Series  
2150 Forecasting? *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(9):11121–11128,  
2151 Jun. 2023. doi: 10.1609/aaai.v37i9.26317. URL [https://ojs.aaai.org/index.php/  
2152 AAAI/article/view/26317](https://ojs.aaai.org/index.php/AAAI/article/view/26317).
- 2153 Yunhao Zhang and Junchi Yan. Crossformer: Transformer Utilizing Cross-Dimension Depen-  
2154 dency for Multivariate Time Series Forecasting. In *Proceedings of the Eleventh International  
2155 Conference on Learning Representations (ICLR)*, 2023. URL [https://openreview.net/  
2156 forum?id=vSVLM2j9eie](https://openreview.net/forum?id=vSVLM2j9eie).  
2157
- 2158 Lifan Zhao and Yanyan Shen. Rethinking Channel Dependence for Multivariate Time Series Fore-  
2159 casting: Learning from Leading Indicators. In *The Twelfth International Conference on Learning  
2160 Representations*, 2024.



2160 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai  
2161 Zhang. Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecast-  
2162 ing. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(12):11106–11115, May  
2163 2021. doi: 10.1609/aaai.v35i12.17325. URL [https://ojs.aaai.org/index.php/  
2164 AAAI/article/view/17325](https://ojs.aaai.org/index.php/AAAI/article/view/17325).

2165 Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. FEDformer: Frequency  
2166 Enhanced Decomposed Transformer for Long-term Series Forecasting. In Kamalika Chaudhuri,  
2167 Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings  
2168 of the 39th International Conference on Machine Learning (ICML)*, volume 162 of *Proceedings  
2169 of Machine Learning Research*, pp. 27268–27286. PMLR, 17–23 Jul 2022. URL [https://  
2170 proceedings.mlr.press/v162/zhou22g.html](https://proceedings.mlr.press/v162/zhou22g.html).

2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213