

487 A Visualization

488 In order to clarify the representation ability more clearly, Figure 5 provides showcases of imputation,
489 long-term forecasting and few-shot forecasting. Especially for few-shot learning, GPT2(6) can
490 accurately forecast, while TimesNet and DLinear fail in this task.

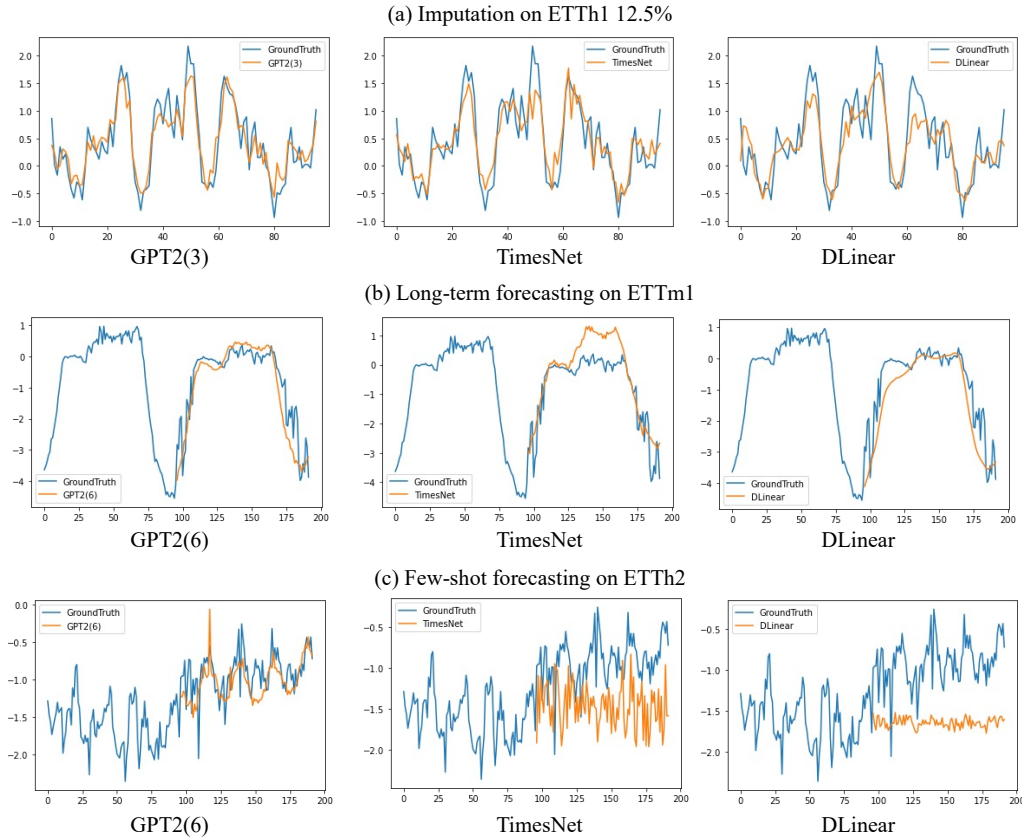


Figure 5: Visualization of imputation, long-term forecasting and few-shot forecasting.

491 B Related Works

492 We have presented a novel general time series analysis model in this paper, and to the best of our
493 knowledge, there has been limited work on similar comprehensive methods for time series analysis.
494 The most closely related field is time series forecasting, where transformer models have gained
495 widespread popularity. Therefore, our focus in this related work will primarily be on introducing the
496 end-to-end time series forecasting method.

497 Time series forecasting models can be roughly divided into three categories, ranging from the classic
498 ARIMA models to the most recent transformer models. The first generation of well-discussed models
499 can be dated back to auto-regressive family, such as ARIMA [Box & Jenkins \(1968\)](#); [Box & Pierce \(1970\)](#)
500 that follows the Markov process and recursively execute sequential forecasting. However,
501 it is limited to stationary sequences while most time series is non-stationary. Additionally, with
502 the bloom of deep neural networks, recurrent neural networks (RNNs), such as LSTM [Hochreiter
503 & Schmidhuber \(1997\)](#) and GRU [Chung et al. \(2014\)](#), were designed for sequential tasks. Yet the
504 recurrent model is inefficient for training and long-term dependencies are still under resolved.

505 Recently, transformer models have achieve great progress in NLP [Vaswani et al. \(2017\)](#); [Devlin et al.
506 \(2019\)](#); [Radford et al. \(2019\)](#) and CV [Dosovitskiy et al. \(2021\)](#); [Bao et al. \(2022\)](#) tasks. Also, a large
507 amount of transformer models are proposed to apply to time series forecasting [Wen et al. \(2023\)](#). In
508 the following, we briefly introduce several representative algorithms. Informer [Zhou et al. \(2021\)](#)

509 proposes a probability sparse attention mechanism to deal with long-term dependencies. Autoformer
 510 Wu et al. (2021) introduces a decomposition transformer architecture and replaces the attention
 511 module with an Auto-Correlation mechanism. FEDformer Zhou et al. (2022) uses Fourier enhanced
 512 structure to improve computational efficiency and achieves linear complexity. Similar to patching in
 513 ViT Dosovitskiy et al. (2021), PatchTST Nie et al. (2022) employs segmentation of time series that
 514 divide a sequence into patches to increase input length and reduce information redundancy. Besides,
 515 a simple MLP-based model DLinear Zeng et al. (2023) outperforms most transformer models and it
 516 validates channel-independence works well in time series forecasting. Recently, TimesNet Wu et al.
 517 (2023) has treated time series as a 2D signal and utilized a convolution-based inception net backbone
 518 to function as a comprehensive time series analysis model. This work is closely related to our tasks
 519 in this paper.

520 C Dataset Details

521 In this section, we separately summarize dataset details long/short-term forecasting and few-shot/zero-
 522 shot forecasting.

523 **Datasets of Long-term Forecasting and Few-shot Learning** The details of datasets are shown as
 524 follows: 1) ETT datasets Zhou et al. (2021) contain electricity load of various resolutions (ETTh &
 525 ETTm) from two electricity stations. 2) Weather contains 21 meteorological indicators of Germany
 526 within 1 year; 3) Illness contains the influenza-like illness patients in the United States; 4) Electricity
 527 dataset contains the electricity consumption; 5) Traffic dataset contains the occupation rate of freeway
 528 system across the State of California. Table 9 summarizes details of feature statistics.

529 Similar to PatchTST Nie et al. (2022), Exchange is not contained. Zeng et al. (2023) shows that
 530 simply repeating the last value in the look-back window can outperform or be comparable to the best
 531 results. Also, ILI is not used for few-shot learning for the limited quantity that is hard to follow the
 532 definition of few-shot.

Table 9: Dataset details of few-shot learning.

Dataset	Length	Dimension	Frequency
ETTh	17420	7	1 hour
ETTM	69680	7	15 min
Weather	52696	22	10 min
ILI	966	7	7 days
Electricity	26304	321	1 hour
Traffic	17544	862	1 hour

533 **Datasets of Short-term Forecasting and Zero-shot Learning** The details of short-term forecasting
 534 and zero-shot learning datasets are shown as follows: 1) M4 is a large and diverse dataset that contains
 535 time series of various frequencies and fields, including business, financial and economic forecasting;
 536 2) M3 is smaller than M4, but also contains time series from diverse domains and frequencies; 3)
 537 TOURISM is the dataset of tourism activities with different frequencies and contains a much higher
 538 fraction of erratic series compared with M4; 4) ELECTR represents the electricity usage monitoring
 539 of 370 customers over three years. Table 6 summarizes details of the datasets and zero-shot mapping
 540 between source and target.

541 D Experimental Details

542 All the deep learning networks are implemented in PyTorch and trained on NVIDIA V100 32GB
 543 GPUs. We use the pre-trained models from Wolf et al. (2020) for experiments. For few-shot learning,
 544 an early stopping counter is employed to stop the training process after three epochs if no loss
 545 degradation on the valid set is observed. Plus, we convert the multivariate data into univariate data.
 546 Specifically, we treat each feature of the sequence as a single time series. This is mainly for memory
 547 efficiency after patching of GPT2(6) and previous works, DLinear and PatchTST, have proved the
 548 effectiveness of channel-independence.

Table 10: Datasets and mapping details of zero-shot learning.

	Dataset		Mapping	
	Length	Horizon	M4	M3
M3 Yearly	645	6	Yearly	-
M3 Quarterly	756	8	Quarterly	-
M3 Monthly	1428	18	Monthly	-
M3 Others	174	8	Monthly	-
M4 Yearly	23000	6	-	Yearly
M4 Quarterly	24000	8	-	Quarterly
M4 Monthly	48000	18	-	Monthly
M4 Weekly	359	13	-	Monthly
M4 Daily	4227	14	-	Monthly
M4 Hourly	414	48	-	Monthly
TOURISM Yearly	518	4	Yearly	Yearly
TOURISM Quarterly	427	8	Quarterly	Quarterly
TOURISM Monthly	366	24	Monthly	Monthly
ELECTR	1311	168	Hourly	Monthly

549 D.1 Accuracy Metrics

550 For long-term/short-term forecasting and few-shot forecasting, we use mean square error (MSE)
551 and mean absolute error (MAE) as metrics. For zero-shot learning, mean absolute percentage error
552 (MAPE) is used for TOURISM; symmetric MAPE (sMAPE) is used for M3 and M4; normalized
553 deviation (ND) is used for ELECTR. All experiments are repeated 3 times and the mean of the metrics
554 is used in the final results.

555 D.2 Detailed Definition and Results for Few-shot and Long-term Forecasting

556 **Task Definition** Since [Zeng et al. \(2023\)](#) and [Nie et al. \(2022\)](#) have verified that channel-independence
557 works well for time series datasets, we treat each multivariate series as multiple independent univariate
558 series. Similar to traditional experimental settings, each time series is split into three parts: training
559 data, validation data, and test data. For the few-shot forecasting task, only a certain percentage (5%,
560 10%) timesteps of training data are used, and the other two parts remain unchanged. The evaluation
561 metrics remain the same as for classic multivariate time series forecasting. We repeat this experiment
562 3 times and report the average metrics in the following experiments.

563 **Detail Experiment Tables for Few-shot Time-Series Forecasting** in Table [11](#) and Table [12](#)

Table 12: Few-shot learning results on 10% data. We use prediction length $O \in \{96, 192, 336, 720\}$. A lower MSE indicates better performance, and the best results are highlighted in bold. '-' means that 10% time series is not sufficient to constitute a training set.

Methods	GPT2(6)	GPT2(0)	DLinear	PatchTST	TimesNet	FEDformer	Autoformer	Stationary	ETSformer	LightTS	Informer	Reformer	
Metric	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	
<i>Wather</i>	96	0.163 0.215	0.190 0.240	0.171 0.224	0.165 0.215	0.184 0.230	0.188 0.253	0.221 0.297	0.192 0.234	0.199 0.272	0.217 0.269	0.374 0.401	0.335 0.380
	192	0.210 0.254	0.243 0.284	0.215 0.263	0.210 0.257	0.245 0.283	0.250 0.304	0.270 0.322	0.269 0.295	0.279 0.332	0.259 0.304	0.552 0.478	0.522 0.462
	336	0.256 0.292	0.270 0.305	0.258 0.299	0.259 0.297	0.305 0.321	0.312 0.346	0.320 0.351	0.370 0.357	0.356 0.386	0.303 0.334	0.724 0.541	0.715 0.535
	720	0.321 0.339	0.348 0.359	0.320 0.346	0.332 0.346	0.381 0.371	0.387 0.393	0.390 0.396	0.441 0.405	0.437 0.448	0.377 0.382	0.739 0.558	0.611 0.500
	Avg.	0.238 0.275	0.263 0.297	0.241 0.283	0.242 0.279	0.279 0.301	0.284 0.324	0.300 0.342	0.318 0.323	0.318 0.360	0.289 0.322	0.597 0.495	0.546 0.469
<i>ETTh1</i>	96	0.458 0.456	0.601 0.536	0.492 0.495	0.516 0.485	0.861 0.628	0.512 0.499	0.613 0.552	0.918 0.639	1.112 0.806	1.298 0.838	1.179 0.792	1.184 0.790
	192	0.570 0.516	0.709 0.587	0.565 0.538	0.598 0.524	0.797 0.593	0.624 0.555	0.722 0.598	0.915 0.629	1.155 0.823	1.322 0.854	1.199 0.806	1.295 0.850
	336	0.608 0.535	0.801 0.635	0.721 0.622	0.657 0.550	0.941 0.648	0.691 0.574	0.750 0.619	0.939 0.644	1.179 0.832	1.347 0.870	1.202 0.811	1.294 0.854
	720	0.725 0.591	1.385 0.831	0.986 0.743	0.762 0.610	0.877 0.641	0.728 0.614	0.721 0.616	0.887 0.645	1.273 0.874	1.534 0.947	1.217 0.825	1.223 0.838
	Avg.	0.590 0.525	0.874 0.647	0.691 0.600	0.633 0.542	0.869 0.628	0.639 0.561	0.702 0.596	0.915 0.639	1.180 0.834	1.375 0.877	1.199 0.809	1.249 0.833
<i>ETTh2</i>	96	0.331 0.374	0.539 0.495	0.357 0.411	0.353 0.389	0.378 0.409	0.382 0.416	0.413 0.451	0.389 0.411	0.678 0.619	2.022 1.006	3.837 1.508	3.788 1.533
	192	0.402 0.411	0.675 0.555	0.569 0.519	0.403 0.414	0.490 0.467	0.478 0.474	0.474 0.477	0.473 0.455	0.785 0.666	2.329 1.104	3.856 1.513	3.552 1.483
	336	0.406 0.433	0.718 0.580	0.671 0.572	0.426 0.441	0.537 0.494	0.504 0.501	0.547 0.543	0.507 0.480	0.839 0.694	2.453 1.122	3.952 1.526	3.395 1.526
	720	0.449 0.464	0.732 0.605	0.824 0.648	0.477 0.480	0.510 0.491	0.499 0.509	0.516 0.523	0.477 0.472	1.273 0.874	3.816 1.407	3.842 1.503	3.205 1.401
	Avg.	0.397 0.421	0.666 0.559	0.605 0.538	0.415 0.431	0.479 0.465	0.466 0.475	0.488 0.499	0.462 0.455	0.894 0.713	2.655 1.160	3.872 1.513	3.485 1.486
<i>ETTm1</i>	96	0.390 0.404	0.610 0.508	0.352 0.392	0.410 0.419	0.583 0.501	0.578 0.518	0.774 0.614	0.761 0.568	0.911 0.688	0.921 0.682	1.162 0.785	1.442 0.847
	192	0.429 0.423	0.666 0.540	0.382 0.412	0.437 0.434	0.630 0.528	0.617 0.546	0.754 0.592	0.781 0.574	0.955 0.703	0.957 0.701	1.172 0.793	1.444 0.862
	336	0.469 0.439	0.895 0.615	0.419 0.434	0.476 0.454	0.725 0.568	0.998 0.775	0.869 0.677	0.803 0.587	0.991 0.719	0.998 0.716	1.227 0.908	1.450 0.866
	720	0.569 0.498	0.916 0.646	0.490 0.477	0.681 0.556	0.769 0.549	0.693 0.579	0.810 0.630	0.844 0.581	1.062 0.747	1.007 0.719	1.207 0.797	1.366 0.850
	Avg.	0.464 0.441	0.772 0.577	0.411 0.429	0.501 0.466	0.677 0.537	0.722 0.605	0.802 0.628	0.797 0.578	0.980 0.714	0.971 0.705	1.192 0.821	1.426 0.856
<i>ETTm2</i>	96	0.188 0.269	0.283 0.344	0.213 0.303	0.191 0.274	0.212 0.285	0.291 0.399	0.352 0.454	0.229 0.308	0.331 0.430	0.813 0.688	3.203 1.407	4.195 1.628
	192	0.251 0.309	0.353 0.384	0.278 0.345	0.252 0.317	0.270 0.323	0.307 0.379	0.694 0.691	0.291 0.343	0.400 0.464	1.008 0.768	3.112 1.387	4.042 1.601
	336	0.307 0.346	0.420 0.422	0.338 0.385	0.306 0.353	0.323 0.353	0.543 0.559	2.408 1.407	0.348 0.376	0.469 0.498	1.031 0.775	3.255 1.421	3.963 1.585
	720	0.426 0.417	0.553 0.491	0.436 0.440	0.433 0.427	0.474 0.449	0.712 0.614	1.913 1.166	0.461 0.438	0.589 0.557	1.096 0.791	3.909 1.543	3.711 1.532
	Avg.	0.293 0.335	0.402 0.410	0.316 0.368	0.296 0.343	0.320 0.353	0.463 0.488	1.342 0.930	0.332 0.366	0.447 0.487	0.987 0.756	3.370 1.440	3.978 1.587
<i>ECL</i>	96	0.139 0.237	0.142 0.240	0.150 0.253	0.140 0.238	0.299 0.373	0.231 0.323	0.261 0.348	0.420 0.466	0.599 0.587	0.350 0.425	1.259 0.919	0.993 0.784
	192	0.156 0.252	0.158 0.254	0.164 0.264	0.160 0.255	0.305 0.379	0.261 0.356	0.338 0.406	0.411 0.459	0.620 0.598	0.376 0.448	1.160 0.873	0.938 0.753
	336	0.175 0.270	0.175 0.271	0.181 0.282	0.180 0.276	0.319 0.391	0.360 0.445	0.410 0.474	0.434 0.473	0.662 0.619	0.428 0.485	1.157 0.872	0.925 0.745
	720	0.233 0.317	0.230 0.315	0.223 0.321	0.241 0.323	0.369 0.426	0.530 0.585	0.715 0.685	0.510 0.521	0.757 0.664	0.611 0.597	1.203 0.898	1.004 0.790
	Avg.	0.176 0.269	0.176 0.270	0.180 0.280	0.180 0.273	0.323 0.392	0.346 0.427	0.431 0.478	0.444 0.480	0.660 0.617	0.441 0.489	1.195 0.891	0.965 0.768
<i>Traf fic</i>	96	0.414 0.297	0.478 0.368	0.419 0.298	0.403 0.289	0.719 0.416	0.639 0.400	0.672 0.405	1.412 0.802	1.643 0.855	1.157 0.636	1.557 0.821	1.527 0.815
	192	0.426 0.301	0.481 0.363	0.434 0.305	0.415 0.296	0.748 0.428	0.637 0.416	0.727 0.424	1.419 0.806	1.641 0.854	1.207 0.661	1.454 0.765	1.538 0.817
	336	0.434 0.303	0.488 0.365	0.449 0.313	0.426 0.304	0.853 0.471	0.655 0.427	0.749 0.454	1.443 0.815	1.711 0.878	1.334 0.713	1.521 0.812	1.550 0.819
	720	0.487 0.337	0.537 0.386	0.484 0.336	0.474 0.331	1.485 0.825	0.722 0.456	0.847 0.499	1.539 0.837	2.660 1.157	1.292 0.726	1.605 0.846	1.588 0.833
	Avg.	0.440 0.310	0.496 0.371	0.447 0.313	0.430 0.305	0.951 0.535	0.663 0.425	0.749 0.446	1.453 0.815	1.914 0.936	1.248 0.684	1.534 0.811	1.551 0.821
Average	0.371 0.367	0.521 0.447	0.413 0.401	0.385 0.376	0.556 0.458	0.511 0.472	0.687 0.559	0.674 0.522	0.912 0.665	1.137 0.712	1.850 0.967	1.888 0.974	

573 **D.4 Mean and STD for Few-shot Learning**

574 Table 14 lists both mean and STD for GPT2(6), DLinear and PatchTST with 3 runs on 5% ETTh2 and
 575 ETTm2. The results show a small variance in performance of GPT2(6) that represents the stability of
 GPT2(6).

Table 14: A subset of results showing both Mean and STD on 5% datasets.

Methods Metric	GPT2-backbone(6 Layers)		
	MSE	MAE	
ETTh2	96	0.376 ± 0.0072	0.421 ± 0.0054
	192	0.418 ± 0.0013	0.441 ± 0.0014
	336	0.408 ± 0.0006	0.439 ± 0.0002
	720	-	-
ETTm2	96	0.199 ± 0.0040	0.280 ± 0.0042
	192	0.256 ± 0.0030	0.316 ± 0.0017
	336	0.318 ± 0.0046	0.353 ± 0.0032
	720	0.460 ± 0.0132	0.436 ± 0.0066

576

577 **D.5 Comparison with Traditional Methods on Few-shot Learning**

578 Since deep learning methods are more advantageous than traditional methods when applied to large
 579 datasets. For few-shot learning, traditional methods should also consider. The results are shown in
 Table 15 that GPT2(6) also achieves best performance.

Table 15: Comparison with traditional methods.

Methods Metric	GPT2(6) 5%		GPT2(6) 10%		ETS		ARIMA		NaiveDrift		
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
ETTh2	96	0.376	0.421	0.331	0.374	2.954	0.742	0.481	0.443	0.764	0.561
	192	0.418	0.441	0.402	0.411	10.226	1.212	0.585	0.495	1.560	0.785
ETTm1	96	0.386	0.405	0.390	0.404	52.237	2.689	0.693	0.547	1.539	0.913
	192	0.440	0.438	0.429	0.423	186.445	4.654	0.710	0.557	2.869	1.215

580

581 **D.6 Baselines with Instance Normalization**

582 Instance normalization Kim et al. (2022) is a plug-in for time series for distribution shift. Most
 583 baselines, such as Autoformer and FEDformer are not equipped with instance normalization. Thus,
 584 for a fair comparison, we add the experiment, as in Table 16 for baselines w/o instance normalization
 585 and GPT(6) can also perform superior.

Table 16: Comparison on 5% data. Autoformer and FEDformer are equipped with instance normaliza-
 tion.

Methods Metric	GPT2(6)		PatchTST		DLinear		Autoformer		Autoformer(Revin)		FEDformer		FEDformer(Revin)		
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
ETTh2	96	0.199	0.280	0.206	0.288	0.236	0.326	0.232	0.322	0.224	0.300	0.229	0.320	0.223	0.298
	192	0.256	0.316	0.264	0.324	0.306	0.373	0.291	0.357	0.296	0.343	0.294	0.361	0.288	0.336

586 **D.7 Detailed Definition and Results of Zero-shot Learning**

587 **Task Definition** Each experiment contains two distinct datasets, source, and target datasets. The
 588 source dataset is used to train the model and then forecasts without fine-tuning in the target dataset.
 589 The target dataset is split into non-overlapping historical and test sequences. We use the historical
 590 sequence as input to the model, and the obtained output is used to calculate errors with the test
 591 sequences. Besides meta-learning-based models like N-BEATS, evaluated models’ parameters are
 592 not allowed any adjustment using the forecasting phase. Also, same as Oreshkin et al. (2021), each
 593 data set adopts a specific metric (M4: sMAPE; M3: sMAPE; TOURISM: MAPE; ELECTR: ND)

594 **Detailed Results** Here, we list detailed performance of zero-shot learning in Table 17, Table
 595 18 and Table 19. For each dataset, we separately list the performance of models under diverse
 596 frequency. Compared to the most recent published method DLinear, GPT2(6) performs superior
 597 in most situations. Also, GPT2(6) does not use any information from the test data, but achieves a
 598 comparable performance of meta-learning based N-BEATS.

Table 17: Zero-shot performance on M4 (sMAPE).

	Yearly (23k)	Quarterly (24k)	Monthly (48k)	Others (5k)	Average (100k)
N-BEATS-FR	13.267	9.596	12.676	4.696	11.675
DLinear-M3	14.193	18.856	14.765	9.194	15.337
TimesNet-M3	15.655	11.877	16.165	6.863	14.553
PatchTST-M3	13.966	10.929	14.664	7.087	13.228
ETSformer-M3	27.846	36.134	25.114	12.338	27.748
LightTS-M3	13.787	11.289	15.181	9.117	13.623
Stationary-M3	14.988	11.686	16.098	6.977	14.327
FEDformer-M3	13.887	11.513	18.154	7.529	15.047
Autoformer-M3	14.552	17.341	25.063	9.666	20.022
Informer-M3	18.542	16.907	23.454	7.348	19.047
Reformer-M3	15.652	11.051	15.604	7.001	14.092
GPT(6)-M3	13.740	10.787	14.630	7.081	13.125

Table 18: Zero-shot performance on M3 (sMAPE).

	Yearly (645)	Quarterly (756)	Monthly (1428)	Others (174)	Average (3003)
N-BEATS-M4	15.07	9.07	13.19	4.29	12.38
N-BEATS-FR	16.43	9.05	13.30	4.51	12.61
DLinear-M4	17.43	9.74	15.65	6.81	14.03
TimesNet-M4	18.75	12.26	14.01	6.88	14.17
PatchTST-M4	15.99	9.62	14.71	9.44	13.39
ETSformer-M4	20.56	11.65	16.97	10.57	16.03
LightTS-M4	15.63	9.40	24.60	8.28	17.90
Stationary-M4	17.05	12.56	16.82	8.13	15.29
FEDformer-M4	16.00	9.48	15.12	8.94	13.53
Autoformer-M4	16.18	13.92	16.91	14.68	15.87
Informer-M4	19.70	13.00	15.91	13.03	15.82
Reformer-M4	16.03	9.76	14.80	7.53	13.37
GPT2(6)-M4	16.42	10.13	14.10	4.81	13.06

599 E Proof

600 In our numerical experiments, we obtain two interesting observations. First, the token similarity
 601 within a sample is larger in pretrained LM. We report the layer-wise average token cosine similarity in
 602 ETTh2 experiment in Figure 7. In particular, Figure 7(a) shows that in a fine-tuned random initialed
 603 GPT2(6) model, the token similarity is around 0.1-0.2 among different layers. When switching to the
 604 frozen pre-trained GPT2-FPT model, the token similarity significantly increases in the deep layers
 605 and eventually reaches more than 0.9 in the last layer. The ETTh2 dataset contains high volatility
 606 hourly information related to the electricity transformer temperature. In this situation, higher token
 607 similarity implies the high-frequency noise in the data is eased and only low-frequency information
 608 will be reserved. In other words, after going through the pretrained GPT2-FPT model, the signal-noise
 609 ratio is enhanced. We use the following theorem to characterize this behavior.

Table 19: Zero-shot performance on Tourism (MAPE).

	Yearly (518)	Quarterly (427)	Monthly (366)	Average (1311)
N-BEATS-M4	23.57	14.66	19.32	18.82
N-BEATS-FR	23.43	14.45	20.47	19.46
DLinear-M4	39.59	18.30	24.76	28.51
TimesNet-M4	35.59	19.22	30.54	28.84
PatchTST-M4	33.23	19.27	27.57	27.10
ETSformer-M4	391.60	35.56	50.47	180.40
LightTS-M4	138.22	16.28	25.34	66.99
Stationary-M4	35.42	35.15	65.58	43.75
FEDformer-M4	43.41	19.88	28.39	31.55
Autoformer-M4	51.19	34.95	31.47	40.39
Informer-M4	41.16	30.98	33.92	35.82
Reformer-M4	33.86	16.85	23.71	25.48
GPT2(6)-M4	27.17	16.21	21.92	22.14

610 **E.1 Theorem E.1**

611 **Theorem E.1** (informal). *We consider the self-attention for l -th query token. Let's assume the input*
 612 *token x_i are bounded with mean μ for $i = 1, 2, \dots, n$. Under mild conditions, with high probability,*
 613 *the output value token V_l converges to μW_v on the order of $\mathcal{O}(n^{-1/2})$, where W_v is the parameter*
 614 *matrix to compute the value token.*

615 The Theorem E.1 describes the self-attention structure can efficiently make output value token V_l
 616 converge its mean value μW_v . In the time series forecasting task, each token represents several
 617 adjacent points in a time series. When the time series has some periodical or translation invariant
 618 structures, by comparing a given token with other tokens, one could have a higher chance to figure
 619 out those invariant structures. This phenomenon is especially important in few-shot forecasting tasks.
 620 Without enough token noise distillation ability, the model will more likely tend to overfit due to
 621 insufficient training data.

622 We denote x_i as i -th element of vector \mathbf{x} , W_{ij} as the element at i -th row and j -th column of matrix
 623 \mathbf{W} , and W_j as the j -th row of matrix \mathbf{W} . Moreover, we denote x_i as the i -th patch (token) of the
 624 inputs with $x_i = \mathbf{X}_i$.
 625

626 Before given the formal statement of the Theorem E.1, we first show the assumptions.

- 627 1. The token x_i is the sub-gaussian random vector with mean μ_i and variance $(\sigma^2/d)I$ for
 628 $i = 1, 2, \dots, n$.
- 629 2. μ follows a discrete distribution with finite values $\mu \in \mathcal{V}$. Moreover, there exist $0 <$
 630 $\nu_1, 0 < \nu_2 < \nu_4$ such that a) $\|\mu_i\| = \nu_1$, and b) $\mu_i \mathbf{W}_Q \mathbf{W}_K^T \mu_i \in [\nu_2, \nu_4]$ for all i and
 631 $|\mu_i \mathbf{W}_Q \mathbf{W}_K^T \mu_j^T| \leq \nu_2$ for all $\mu_i \neq \mu_j \in \mathcal{V}$.
- 632 3. \mathbf{W}_V and $\mathbf{W}_Q \mathbf{W}_K^T$ are element-wise bounded with ν_5 and ν_6 respectively, that is, $|\mathbf{W}_V^{(ij)}| \leq$
 633 ν_5 and $|(\mathbf{W}_Q \mathbf{W}_K^T)^{(ij)}| \leq \nu_6$, for all i, j from 1 to d .

634 In the above assumptions, we ensure that for a given query patch, the difference between the clustering
 635 center and noises are large enough to be distinguished.

636 **Theorem E.2** (formal statement of Theorem E.1). *Let patch x_i be σ^2 -subgaussian random variable*
 637 *with mean μ_i and all n patches follow the same clustering center of query l . Per Assumptions*
 638 *above mentioned, when $\sqrt{d} \geq 3(\psi(\delta, d) + \nu_2 + \nu_4)$, then with probability $1 - 5\delta$, we have*

$$\left\| \frac{\sum_{i=1}^n \exp\left(\frac{1}{\sqrt{d}} \mathbf{x}_i \mathbf{W}_Q \mathbf{W}_k^\top \mathbf{x}_i\right) \mathbf{x}_i \mathbf{W}_V}{\sum_{j=1}^n \exp\left(\frac{1}{\sqrt{d}} \mathbf{x}_j \mathbf{W}_Q \mathbf{W}_k^\top \mathbf{x}_j\right)} - \boldsymbol{\mu}_i \mathbf{W}_V \right\|_\infty \leq 4 \exp\left(\frac{\psi(\delta, d)}{\sqrt{d}}\right) \sigma \nu_5 \sqrt{\frac{2}{dn} \log\left(\frac{2d}{\delta}\right)}$$

$$+ 7 \left[\exp\left(\frac{\nu_2 - \nu_4 + \psi(\delta, d)}{\sqrt{d}}\right) - 1 \right] \|\boldsymbol{\mu}_i \mathbf{W}_V\|_\infty,$$

639 where $\psi(\delta, d) = 2\sigma\nu_1\nu_6\sqrt{2\log\left(\frac{1}{\delta}\right)} + 2\sigma^2\nu_6\log\left(\frac{d}{\delta}\right)$.

640

641 *Proof.* See the proof of Lemma 2 in [Wang et al. \(2022\)](#) with $k_1 = k = n$. ■

642 E.2 Theorem [E.4](#)

643 We first give the formal statement of Theorem [E.4](#).

644 **Theorem E.3** (formal statement of Theorem [E.4](#)). *Let $\mathbf{g}_i \in \mathbb{R}^d$ and $\mathbf{y}_i \in \mathbb{R}^T$ be the feature*
 645 *map vector and forecasting targets for the sample $i = 1, 2, \dots, N$ respectively, and we assume*
 646 *$\frac{1}{N} \sum_{i=1}^N \mathbf{g}_i \mathbf{g}_i^\top \succeq \sigma I$ for some $\sigma > 0$. We want to learn a matrix $\mathbf{W} \in \mathbb{R}^{d \times T}$ from the following*
 647 *optimization problem:*

$$\mathbf{W} = \arg \min \frac{1}{2N} \sum_{i=1}^N \|\mathbf{W} \mathbf{g}_i - \mathbf{y}_i\|_2^2. \quad (1)$$

648 *If we apply stochastic gradient descent with diminishing step sizes $\eta_t = \frac{1}{\sigma t}$ at step t , we will need*
 649 *$t = \tilde{\mathcal{O}}(\epsilon^{-1} \sigma^{-1})$ steps to reach*

$$\frac{1}{t} \sum_{j=1}^t \left(\frac{1}{2N} \sum_{i=1}^N \|\mathbf{W}_j \mathbf{g}_i - \mathbf{y}_i\|_2^2 \right) - \frac{1}{2N} \sum_{i=1}^N \|\mathbf{W}^* \mathbf{g}_i - \mathbf{y}_i\|_2^2 \leq \epsilon, \quad (2)$$

650 where \mathbf{W}^* is the optimal solution and \mathbf{W}_j is the j step's solution and $\tilde{\mathcal{O}}$ we suppress the logarithmic
 651 dependence.

652 *Proof.* As we assume $\frac{1}{N} \sum_{i=1}^T \mathbf{g}_i \mathbf{g}_i^\top \succeq \sigma I$, the hessian of optimization problem in [\(1\)](#) is also
 653 positive definite, which is equivalent to the optimization problem in [\(1\)](#) is strongly convex with
 654 parameter proportional to σ . Then via standard stochastic gradient decent analysis (e.g., section 3.1
 655 in [Lacoste-Julien et al. \(2012\)](#)), we obtain:

$$\frac{1}{t} \sum_{j=1}^t \left(\frac{1}{2N} \sum_{i=1}^N \|\mathbf{W}_j \mathbf{g}_i - \mathbf{y}_i\|_2^2 \right) - \frac{1}{2N} \sum_{i=1}^N \|\mathbf{W}^* \mathbf{g}_i - \mathbf{y}_i\|_2^2 \leq \mathcal{O}\left(\frac{\log t}{\sigma t}\right) = \tilde{\mathcal{O}}(\sigma^{-1} t^{-1}). \quad (3)$$

656 Therefore, to reach ϵ optimization gap, we just need to set $t = \tilde{\mathcal{O}}(\sigma^{-1} \epsilon^{-1})$. ■

657 The second observation is that for the pretrained GPT2-FPT model, the last transformer layer's
 658 outputs, i.e., feature maps, are spread widely throughout the feature space. We report the t-SNE
 659 visualization of the feature maps for GPT2-FPT and an end-to-end model PatchTST in Figure [8](#). In
 660 Figure [8](#) (a) and (b), we color the samples chunked from the one single time series into the same
 661 color and the same configuration of the T-SNE is applied. One may observe that the feature maps of
 662 GPT2-FPT has less concentration compared to PatchTST. It implies the GPT2-FPT's feature maps
 663 corresponding to different samples are more distinctive which eventually facilitates the learning ability
 664 of the last MLP layer. Researchers [Wang & Isola \(2020\)](#) have found that contrastive learning-based
 665 representation learning may result in a uniform distribution of training data, and such behavior plays
 666 an important role in its good downstream task performance. We use the following theorem to justify
 667 it.

668 **Theorem E.4** (informal). *Let \mathbf{g}_i and \mathbf{y}_i be the feature map vector and forecasting targets for the*
 669 *sample $i = 1, 2, \dots, N$ respectively, and we assume $\frac{1}{N} \sum_{i=1}^N \mathbf{g}_i \mathbf{g}_i^\top \succeq \sigma I$ for some $\sigma > 0$. Under mild*
 670 *conditions, if we train an MLP layer that maps feature maps to forecasting targets via the stochastic*
 671 *gradient descent, the total step to reach some optimization tolerance is on the order of $\mathcal{O}(\sigma^{-1})$.*

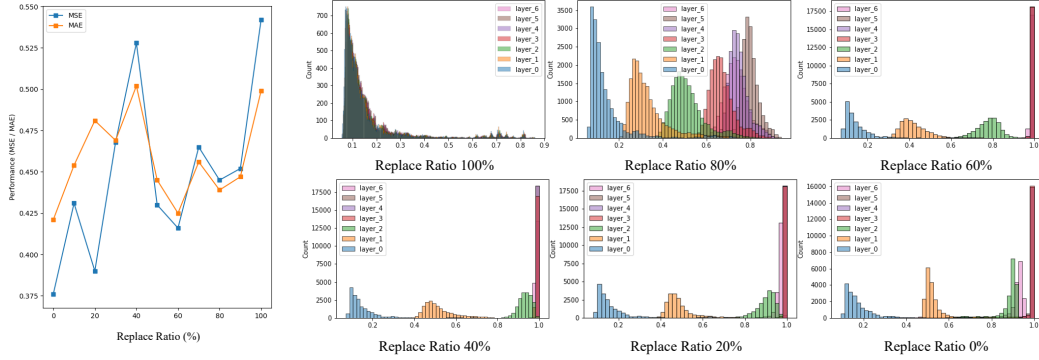


Figure 6: The performance and token similarity within samples with respect to each layer with different random replace ratios. Pretrained parameters are replaced by random initial parameters according to certain proportions.

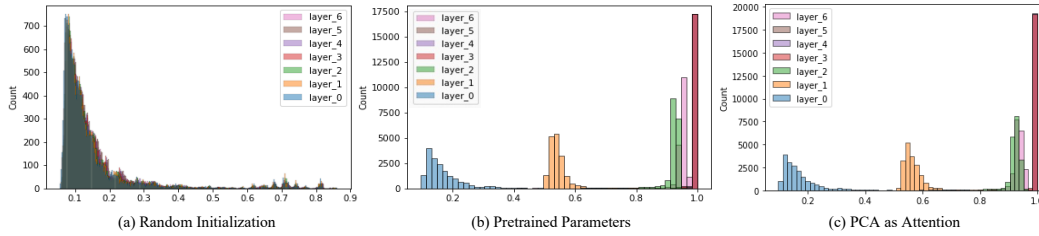


Figure 7: The token similarity within samples with respect to each layer. (a) GPT2-noPretrain-model; (b) GPT2-Pretrained-model; (c) Pretrained attention is replaced by PCA.

672 The Theorem [E.4](#) considers the covariate matrix of feature maps being positive definite that indicates
 673 the set of all feature maps $\{g_i\}$ spans the whole feature spaces, and the higher spread level gives a
 674 larger σ . In this case, if we only want to learn an MLP layer, the problem reduces to a well-conditioned
 675 least-squared regression problem. Then the fast convergence rate is achieved.

676 Efficiently learning the last MLP layer plays a very important role in time series forecasting and
 677 can substantially impact the prediction performance. In [Zeng et al. \(2023\)](#), the authors show that
 678 learning a single MLP layer can also bring very promising performance. In few-shot forecasting, the
 679 pre-trained GPT2 model may still preserve highly diverse feature maps than end-to-end type models
 680 and eventually leads to fast learning speed on the last MLP layer.

681 Another possible benefit of wide spared feature maps is enhancing the model memorization ability
 682 when using a multi-layer decoder structure. In the literature on network memorization ability (e.g.,
 683 [Vardi et al. \(2021\)](#); [Yun et al. \(2020\)](#)), the deep learning model tends to have better memorization
 684 ability when feature maps are well separated. In forecasting tasks, capturing extreme or rare behavior
 685 is very important. The pretrained GPT gains more capacity in the decoder to correctly forecast
 686 uncommon time series.

687 **F N-gram Explanation for Universality**

688 Why does the proposed pretrained-frozen-model work so effectively? We have achieved state-of-
 689 the-art performance in time series analysis using a language model that is mostly trained on natural
 690 language data. The answer lies in the universality of the frozen structure, which includes attention
 691 layers and Feed Forward layers. We can represent images and time series forecasting tasks as an
 692 n-gram estimation problem, akin to text analysis, by employing a patching approach. This method
 693 treats subsequences of time series or image patches as individual tokens. Central to sequential
 694 prediction is the n -order Markov process, and a simple way to capture the n -order Markov process
 695 is n -gram language model. To predict next token w_0 , we need to compute $p(w_0|w_1, \dots, w_{n-1})$,
 696 which can be further computed as $p(w_0 w_1 \dots w_{n-1})/p(w_1 \dots w_{n-1})$. Hence, the core of n -gram

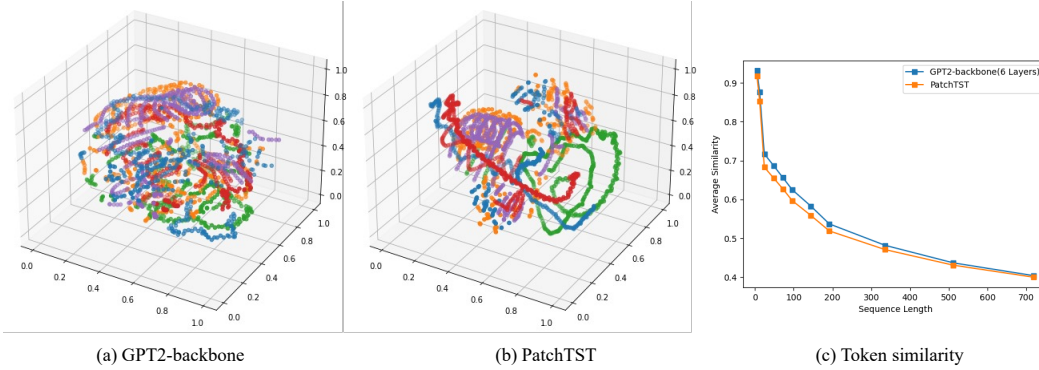


Figure 8: The t-SNE visualization of sample feature maps for (a) GPT-backbone, (b) end-to-end-PatchTST-model. (c) The token similarity within samples within different continuous sequence lengths.

697 language model is to estimate the probability of observing a sequence of n tokens. When n is large,
 698 most of n token sequences will not be observed from data, leading to the sparse data problem, a
 699 common challenge faced by n -gram language model. As a result, a large body of research in n -gram
 700 language model is focused on how to effectively estimate probability of having n -token sequences
 701 even when they are NOT observed from data. We hypothesize that the transformer model pretrained
 702 by GPT-2 essentially allows us to estimate $p(w_0 w_1 \dots w_{n-1})$ from observations of significantly
 703 shorter token sequences. In this section, we will show that the function of estimating probabilities of
 704 longer sequences from observation of shorter sequences is universal and is independent from domain
 705 as long as data exhibit a skew distribution (e.g., follows a power law). We note that our work is closely
 706 related to the discussion presented in [Elhage et al. \(2021\)](#); [Olsson et al. \(2022\)](#), where the authors
 707 also connect the function of transformer to compute of n -grams. We however note that our key result
 708 is to show the universality in computing probability of longer sequences from observations of shorter
 709 sequences, which can't be found in any existing studies. Although the discussion is restricted to
 710 discrete tokens, it should be generalized to continuous signals as we can always quantize continuous
 711 signals into a finite number of discrete tokens, similar to what BEiT [Bao et al. \(2022\)](#) did.

712 To gain a better understanding, let's start by examining a "zero-layer" Transformer model. This model
 713 operates by taking a token, embedding it, and transforming it back to produce logits that predict the
 714 subsequent token. Because it cannot transfer information from other tokens, it relies solely on the
 715 current token to predict the next one. Consequently, the optimal behavior of this model is to closely
 716 resemble the **bigram** log-likelihood.

717 Then we move on to the so-called "attention-only" transformer, which doesn't have MLP layers.
 718 As discussed in a recent work [Elhage et al. \(2021\)](#), one-layer attention-only Transformers can be
 719 comprehended as a combination of a **bigram** model and multiple "**skip-trigram**" models (impacting
 720 the probabilities of sequences "A... BC"). This can be intuitively understood as each attention head
 721 having the ability to selectively attend from the current token ("B") to a previous token ("A") and
 722 transfer relevant information to fine-tune the probability of potential subsequent tokens ("C"). [Olsson](#)
 723 [et al. \(2022\)](#) further discusses a multi-layer transformer can do more complex n -gram estimation
 724 using an induction heads mechanism. To be more precise, induction heads employ a straightforward
 725 principle: the '[A][B] ... [A] → [B]' rule, which elevates the likelihood of generating the subsequent
 726 token 'B' given the current token 'A' if there is a fuzzy match of the AB bigram in the historical
 727 context. This rule seems to largely decouple A and B, which means they do not memorize a fixed
 728 table of n -gram statistics. The rule [A][B] ... [A] → [B] applies regardless of what A and B are,
 729 which can abstract to new patterns.

730 Building upon these discussions, we are now prepared to substantiate the following argument: **For**
 731 **sequential data following a power law, there is a potentially universal solution to the final**
 732 **estimation of n -gram probabilities.** That's the reason behind the universality of pretrained LM's
 733 performance in cross-domain tasks. For simplicity, we assume that n is so large that we are unable to
 734 observe any occurrence of n -gram from data, and we only observe the occurrence of n' -grams with
 735 $n' < n$. We denote by s_i^n the i th unique n -gram, and by the notation $s_j^{n'} \in s_i^n$ if n' -gram $s_j^{n'}$ appears

736 in s_i^n , the i th n -gram. Let m_n be the number of unique n -grams. According to the maximum entropy
 737 model, our estimation of n -gram probabilities can be cast into the following optimization problem:

$$738 \quad \min \sum_{i=1}^{m_n} p(s_i^n) \log p(s_i^n) \quad \text{s. t.} \quad \sum_{i: s_j^{n'} \in s_i^n} p(s_i^n) = \widehat{p}(s_j^{n'}) \quad \text{where } \widehat{p}(s_j^{n'}) \text{ represents the proba-}$$

$$739 \quad \text{bility of observing pattern } s_j^{n'} \text{ from the data and } j \in [m_{n'}], n' \in [n-1].$$

740 For each constraint for $\widehat{p}(s_j^{n'})$, we introduce a Lagrangian dual variable $\lambda_j^{n'}$, and rewrite the optimiza-
 741 tion problem as follows:

$$742 \quad \min_{\lambda} \log \left(\sum_{i=1}^{m_n} \exp \left(\sum_{(n',j): s_j^{n'} \in s_i^n} \lambda_j^{n'} \right) \right) - \sum_{n'=1}^{n-1} \sum_{j=1}^{m_{n'}} \lambda_j^{n'} \widehat{p}(s_j^{n'}),$$

$$743 \quad \text{where } n\text{-gram probability } p(s_i^n) \text{ is given as } p(s_i^n) = \frac{1}{Z(\lambda)} \exp \left(\sum_{(n',j): s_j^{n'} \in s_i^n} \lambda_j^{n'} \right) \text{ and } Z(\lambda) =$$

$$744 \quad \sum_{i=1}^{m_n} \exp \left(\sum_{(n',j): s_j^{n'} \in s_i^n} \lambda_j^{n'} \right)$$

745 In the case that all n -grams follow a power law, for each $n' \in [n-1]$, we divide n' -gram into two
 746 groups: the group $\mathcal{V}_{n'}$ includes the high frequency n' -gram and the group $\mathcal{U}_{n'}$ including the low
 747 frequency of n' -gram. For simplicity, we assume that the probability for all the high frequency
 748 n' -grams are roughly $\alpha_{n'} \in [0, 1]$ and the probability for all the low frequency n' -grams are roughly
 749 $\beta_{n'} \in [0, 1]$. By assuming that all the patterns in $\mathcal{V}_{n'}$ and $\mathcal{U}_{n'}$ share similar appearance frequency, we
 750 simplify the optimization problem by only introducing two dual variables for each n' -gram, i.e. $\lambda_a^{n'}$
 751 for high-frequency patterns and $\lambda_b^{n'}$ for low-frequency patterns as follow Using these notations, we
 752 have the optimization problem simplified as

$$753 \quad \min_{\lambda} \log \left(\sum_{i=1}^{m_n} \exp \left(\sum_{n'=1}^{n-1} \sum_{j: s_j^{n'} \in s_i^n} \lambda_a^{n'} I(s_j^{n'} \in \mathcal{V}_{n'}) \right. \right. \\ \left. \left. + \lambda_b^{n'} I(s_j^{n'} \in \mathcal{U}_{n'}) \right) \right) - \sum_{n'=1}^{n-1} \left(\lambda_a^{n'} g_{n'} + \lambda_b^{n'} h_{n'} \right),$$

$$754 \quad \text{where } g_{n'} = \sum_{s_j^{n'} \in \mathcal{V}_{n'}} \widehat{p}(s_j^{n'}) \text{ and } h_{n'} = \sum_{s_j^{n'} \in \mathcal{U}_{n'}} \widehat{p}(s_j^{n'}).$$

755 Furthermore, let $q_a^{n'}$ be the probability to observe a high frequency n' -gram appearing in any n -gram,
 756 and $q_b^{n'}$ be the probability to observe a low frequency n' -gram appearing in any n -gram, we have

$$757 \quad \sum_{i=1}^{m_n} \exp \left(\sum_{n'=1}^{n-1} \sum_{j: s_j^{n'} \in s_i^n} \lambda_a^{n'} I(s_j^{n'} \in \mathcal{V}_{n'}) + \lambda_b^{n'} I(s_j^{n'} \in \mathcal{U}_{n'}) \right) \\ = m_n \prod_{n'=1}^{n-1} (1 + q_a^{n'} \exp(\lambda_a^{n'})) (1 + q_b^{n'} \exp(\lambda_b^{n'})) + \mathcal{O}(\sqrt{m_n}).$$

758 By skipping the term $\mathcal{O}(\sqrt{m_n})$, we further simplify the optimization problem as

$$759 \quad \min_{\lambda} \sum_{n'=1}^{n-1} \log \left(1 + q_a^{n'} \exp(\lambda_a^{n'}) \right) - \sum_{n'=1}^{n-1} \log \left(1 + q_b^{n'} \exp(\lambda_b^{n'}) \right) - \lambda_b^{n'} h_{n'},$$

760 which is equivalent to

$$761 \quad \lambda_a^{n'} = \min_{\lambda} \log \left(1 + q_a^{n'} \exp(\lambda) \right) - \lambda g_{n'} \quad \text{As illustrated by the above analysis, dual variables}$$

$$\lambda_b^{n'} = \min_{\lambda} \log \left(1 + q_b^{n'} \exp(\lambda) \right) - \lambda h_{n'}.$$

762 $\lambda_a^{n'}$ and $\lambda_b^{n'}$ will only depend on statistics $q_a^{n'}$, $q_b^{n'}$, $g_{n'}$ and $h_{n'}$. They are independent from the
 763 detailed statistics $\widehat{p}(s_j^{n'})$ and how each n' -gram appears in different n -gram. Thus, this simple
 764 analysis does indicate, to some degree, that the solution obtained from the maximum entropy model
 765 can be universal, as long as n -grams follow skewed distributions like power law.

766 We informally demonstrate that transformer models utilize attention mechanisms to perform a
 767 sophisticated form of n -gram estimation, and the generation rule for such n -gram distributions could
 768 be universal. This is how universality is achieved in our proposed cross-domain knowledge transfer.
 769 However, we currently lack a concrete metric to evaluate the performance of knowledge transfer
 770 between different domains, which requires further investigation. Nonetheless, in our experimental
 771 study, we demonstrate that a transformer model (beit) [\[Bao et al. \(2022\)\]](#) trained on images can perform
 772 well on cross-domain time series forecasting tasks.

773 G Connection between self-attention and Principle component analysis

774 Understand the Gradient Structure of Self-Attention

775 Let $X = (x_1, \dots, x_N)^\top \in \mathbb{R}^{N \times D}$ be the input pattern, and let $f(X) = (f_1(X), \dots, f_N(X))^\top : \mathbb{R}^{N \times D} \mapsto \mathbb{R}^{N \times D}$ be the function for self-attention, i.e.

$$777 \quad f_i(X) = \text{softmax}(XAX^\top)X$$

778 where $A = W_Q W_K^\top \in \mathbb{R}^{D \times D}$. Let the Jacobian $J = \left[\frac{\partial f_i(X)}{\partial x_j} \right]_{i,j=1}^N$ represent the gradient $f(X)$
779 with respect to input pattern. The lemma below shows an important structure of J .

$$780 \quad \text{Lemma G.1.} \quad |J|_2 \leq |A|_2 \sum_{i=1}^N \left(P_{i,i} + \frac{1}{2} \right) \left| x_i - \sum_{j=1}^N P_{i,j} x_j \right|^2 + \Delta$$

781 where $\Delta = |A|_2 \sum_{i \neq j}^N P_{i,j} \left| x_j - \sum_{k=1}^N P_{i,k} x_k \right|^2 + \frac{|A|_2}{2} \sum_{j=1}^N |x_i|^2$ and

$$782 \quad P_{i,j} = \frac{\exp(x_i^\top A x_j)}{\sum_{k=1}^N \exp(x_i^\top A x_k)}$$

783 *Proof.* According to the analysis from the work, we have the gradient $J_{i,j} = \frac{\partial f_i(X)}{\partial x_j}$ is given by

784 $J_{i,j} = P_{i,j} I + X^\top Q^i (X A \delta_{i,j} + E_{j,i} X A^\top)$ where $Q^i = \text{diag}(P_{i,:}) - P_{i,:} P_{i,:}^\top$. Here $P_{i,:} \in \mathbb{R}_+^N$ represents the i -th row of matrix P . We thus have

$$\begin{aligned} |J|_2 &\leq \sum_{i,j=1}^N |J_{i,j}|_2 \\ &\leq \sum_{i,j=1}^N P_{i,j} + \sum_{i=1}^N |X^\top Q^i X|_2 |A|_2 + \sum_{i,j=1}^N |X^\top Q^i E_{j,i} X|_2 |A|_2 \\ &\leq N + |A|_2 \sum_{i=1}^N \left(\sum_{j=1}^N P_{i,j} |x_j|^2 - \left| \sum_{j=1}^N P_{i,j} x_j \right|^2 \right) + |A|_2 \sum_{i,j=1}^N |X^\top Q^i e_j x_i^\top| \\ 786 &\leq N + |A|_2 \sum_{i=1}^N \sum_{j=1}^N P_{i,j} \left| x_j - \sum_{k=1}^N P_{i,k} x_k \right|^2 + |A|_2 \sum_{i,j=1}^N P_{i,j} |x_i^\top (x_j - X^\top P_{i,:})| \\ &\leq \underbrace{|A|_2 \sum_{i=1}^N \left(P_{i,i} + \frac{1}{2} \right) |x_i - X^\top P_{i,:}|^2 + N + |A|_2 \sum_{i \neq j}^N P_{i,j} |x_j - X^\top P_{i,:}|^2 + \frac{|A|_2}{2} \sum_{j=1}^N |x_i|^2}_{:=\Delta} \end{aligned}$$

787 As indicated by Lemma 1, one of the key components in the upper bound of Jacobian is $|x_i - \sum_{j=1}^N P_{i,j} x_j|^2$. Thus, through the optimization, we like to reduce the size of the gradient and
788 therefore may prefer to reduce the quantity to $\sum_{i=1}^N |x_i - \sum_{j=1}^N P_{i,j} x_j|^2$. Hence, it will be interesting
789 to understand the choice of W^Q and W^K that leads to the minimization of $\sum_{i=1}^N |x_i - \sum_{j=1}^N P_{i,j} x_j|^2$,
790 i.e. the following optimization problem $\min_{|A|_F \leq \rho} \sum_{i=1}^N \left| x_i - \sum_{j=1}^N P_{i,j} x_j \right|^2$ where ρ is introduced
791 to control the size of A .

793 Connection between Self-Attention and Principal Component Analysis

794 Let consider the optimization problem in (G) when ρ is small, we can approximate $P_{i,j}$ as $P_{i,j} \approx \frac{1}{N} + \frac{1}{N} x_i^\top A x_j$. Define $\bar{x} = X^\top \mathbf{1}/N$. We have
795 $\sum_{i=1}^N |x_i - X^\top P_{i,:}|^2 = \sum_{i=1}^N |x_i - \bar{x} - X^\top X A x_i|^2$. By assuming that all the input patterns are
796 zero centralized, we have $\bar{x} = 0$ and $\sum_{i=1}^N |x_i - X^\top X A x_i|^2 = \text{tr}((I - X^\top X A)^2 X^\top X)$. The
797 theorem below shows that A minimizing the objective $\sum_{i=1}^N |x_i - X^\top X A x_i|^2$ contains the largest
798 m eigenvectors of $X^\top X$ where m is the rank of A .

800 *Theorem 2.* Let W_Q and W_K be matrices of size $D \times m$. Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$ be the
801 eigenvalues of $X^\top X$ ranked in descending order, and let $v_i \in \mathbb{R}^D, i = 1, \dots, D$ be the corresponding
802 eigenvectors. The optimal solution A^* that minimizes $\sum_{i=1}^N |x_i - X^\top X A x_i|^2$ is given by
803 $A = \sum_{i=1}^m \frac{1}{\lambda_i} v_i v_i^\top$

804 *Proof.* Since $W_Q, W_K \in \mathbb{R}^{D \times m}$ where $m < D$, we know that A is a matrix of rank m . Hence,
805 we know $\min_A \sum_{i=1}^N |x_i - X^\top X A x_i|^2 \geq \sum_{k=m+1}^D \lambda_k$. We also know that by choosing A as

$$806 \quad A = \sum_{i=1}^m \frac{1}{\lambda_i} v_i v_i^\top \quad \text{we have}$$

$$\sum_{i=1}^N |x_i - X^\top X A x_i|^2 = \text{tr} \left(\left(I - \sum_{i=1}^m v_i v_i^\top \right)^2 X^\top X \right) = \sum_{k=m+1}^D \lambda_k$$

Hence, the solution A for minimizing $\sum_{i=1}^N |x_i - X^\top X A x_i|^2$ is essentially a weighted combination of top eigenvectors of $X^\top X$. Since a small gradient will prefer a small quantity of $\sum_{i=1}^N |x_i - X^\top X A x_i|^2$, by minimizing through the self-attention layer, we essentially choose weight matrix W_Q and W_K to be aligned with the principal directions of $X^\top X$. ■

H Experiment Analysis and Other Key Results

H.1 Experiment analysis of GPT2-FPT model

In this section, we conduct experiments to analyze whether the self-attention frozen pre-trained model improves performance compared with overall fine-tuning and random initialization. Firstly, we compare GPT2(6) FPT with the same model without freezing (No Freeze) and random initial model (No Pre-train). For the end-to-end paradigm No Pre-train GPT2-backbone (6 Layers), we directly train all parameters of the model. We summarize the results in Table 20 and Table 21. Then we analyze the performance of various layers to clarify our selection of GPT2(6) FPT.

Table 20: Model analysis results on 5% data. We use prediction length $O \in \{96, 192, 336, 720\}$ for ILI and $O \in \{24, 36, 48, 60\}$ for others.

Methods		GPT2(6)		No Freeze		No Pretrain	
Metric		MSE	MAE	MSE	MAE	MSE	MAE
<i>Weather</i>	96	0.175	0.230	0.183	0.229	0.199	0.254
	192	0.227	0.276	0.275	0.300	0.262	0.302
	336	0.286	0.322	0.297	0.331	0.326	0.345
	720	0.366	0.379	0.380	0.388	0.405	0.396
<i>ETTh1</i>	96	0.543	0.506	0.671	0.564	0.882	0.643
	192	0.748	0.580	0.907	0.632	1.389	0.817
	336	0.754	0.595	0.931	0.655	2.968	1.149
	720	-	-	-	-	-	-
<i>ETTh2</i>	96	0.376	0.421	0.440	0.449	0.465	0.457
	192	0.418	0.441	0.503	0.478	0.614	0.536
	336	0.408	0.439	0.691	0.572	0.596	0.529
	720	-	-	-	-	-	-
<i>ETTm1</i>	96	0.386	0.405	0.429	0.432	0.394	0.410
	192	0.440	0.438	0.496	0.470	0.432	0.432
	336	0.485	0.459	0.535	0.489	0.491	0.464
	720	0.557	0.499	0.786	0.592	0.564	0.503
<i>ETTm2</i>	96	0.199	0.280	0.217	0.293	0.301	0.353
	192	0.256	0.316	0.300	0.350	0.321	0.365
	336	0.318	0.353	0.331	0.368	0.371	0.398
	720	0.460	0.439	0.460	0.436	0.659	0.528

819

Fine-tune More Parameters Compared with fine-tuning all parameters, self-attention frozen pre-trained model GPT2(6) FPT achieves better performance on most datasets and yields an overall **12.7%** relative MSE reduction on 5% data and **11.5%** relative MSE reduction on 10% data. It verifies that frozen pre-trained attention layers are effective for time series forecasting.

Parameters Initialization Compared with the random initial model, self-attention frozen pre-trained model GPT2(6) FPT achieves better performance on most datasets and yields an overall **21.2%** relative MSE reduction on 5% data and **14.3%** relative MSE reduction on 10% data. It again suggests that a model pre-trained on cross-domain data can achieve significant performance improvement in time series forecasting.

The Number of GPT2 Layers For most transformer-based methods in time-series forecasting Zhou et al. (2022); Wu et al. (2021); Nie et al. (2022), no more than 3 encoder layers are included. However, most pre-trained models with at least 12 layers may suffer from overfitting in time series forecasting. To better balance performance and computational efficiency, we test using various numbers of layers on ETTh2. Additionally, we train a completely random initialized non-pretrained GPT2 as a comparison. The results are shown in Figure 9, for both 5% and 10% data, the pre-trained model is unable to do well with few layers but significantly outperforms non-pre-trained GPT2 with more attention blocks transferred from NLP. It indicates that pre-trained attention layers produce

Table 21: No Pretrain and No Freeze results on 10% data. We use prediction length $O \in \{96, 192, 336, 720\}$ for ILI and $O \in \{24, 36, 48, 60\}$ for others.

Methods		GPT2(6)		No Freeze		No Pretrain	
Metric		MSE	MAE	MSE	MAE	MSE	MAE
<i>Weather</i>	96	0.163	0.215	0.168	0.221	0.175	0.229
	192	0.210	0.254	0.238	0.286	0.244	0.287
	336	0.256	0.292	0.289	0.318	0.301	0.325
	720	0.321	0.339	0.398	0.383	0.390	0.378
<i>ETTh1</i>	96	0.458	0.456	0.605	0.532	0.680	0.560
	192	0.570	0.516	0.713	0.579	0.738	0.602
	336	0.608	0.535	0.747	0.586	0.893	0.641
	720	0.725	0.591	0.945	0.688	2.994	1.169
<i>ETTh2</i>	96	0.331	0.374	0.369	0.394	0.422	0.433
	192	0.402	0.411	0.464	0.455	0.482	0.466
	336	0.406	0.433	0.420	0.439	0.540	0.496
	720	0.449	0.464	0.535	0.515	0.564	0.519
<i>ETThm1</i>	96	0.390	0.404	0.429	0.430	0.385	0.401
	192	0.429	0.423	0.463	0.446	0.426	0.421
	336	0.469	0.439	0.510	0.470	0.506	0.455
	720	0.569	0.498	0.780	0.591	0.576	0.505
<i>ETThm2</i>	96	0.188	0.269	0.243	0.311	0.244	0.315
	192	0.251	0.309	0.307	0.352	0.318	0.363
	336	0.307	0.346	0.337	0.364	0.409	0.412
	720	0.426	0.417	0.471	0.440	0.473	0.450

837 a great benefit in time series forecasting. Also, the pre-trained model achieves better performance
 838 between 3 and 9 layers. Thus GPT2 with 6 layers is chosen as our default architecture.

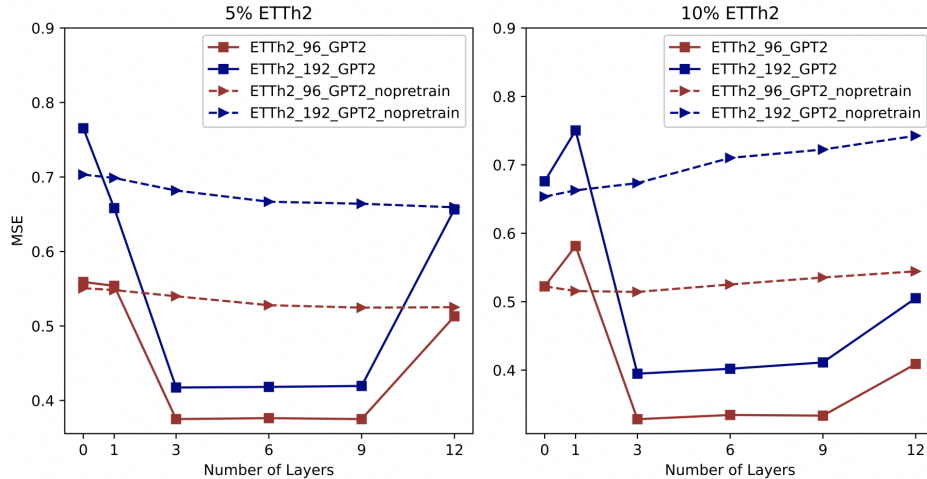


Figure 9: Comparison of pre-trained and non-pre-trained GPT2 with various layers on ETTh2. Color represents various prediction length $O \in \{96, 192\}$ and line style means different models .

839 H.2 No Pre-training but Freezing

840 For comprehensively ablation on pre-training and freezing strategies, we also add experiment for
 841 random initialized GPT2(6) with freezing. The results in Table 22 shows that only input and output
 842 modules can not work and pre-trained knowledge play an importance part in time series tasks.

Table 22: Ablation on random initialized model with freezing.

Methods	GPT2(6)		No Freeze		No Pretrain		No Pretrain + Freeze		
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
<i>ETTh2</i>	96	0.376	0.421	0.440	0.449	0.465	0.457	0.540	0.497
	192	0.418	0.441	0.503	0.478	0.614	0.536	0.721	0.580

843 H.3 Fine-Tuning Parameters Selection

844 In this section, we conduct ablation experiments to study which parameters are important to fine-tune.
 845 Since the input embedding and output layers are randomly initialized for adapting to a new domain,
 846 they must be trained. Then, we study adding layer normalization and positional embeddings to the list
 847 of fine-tuning parameters. Table 23 shows the results that re-train parameters of layer normalization
 848 and positional embeddings can bring certain benefits, especially in longer prediction lengths. Thus,
 849 we follow the standard practice to re-train positional embeddings and layer normalization.

Table 23: Ablation by fixing positional embeddings or layer normalization on 5% ETTm1 and ETTm2. Parameters of GPT2(6) are successively added to the list of fine-tuned parameters.

Methods	Input & Output	+ LN		+ POS			
Metric	MSE	MAE	MSE	MAE	MSE	MAE	
ETTh1	96	0.395	0.410	0.392	0.409	0.386	0.405
	192	0.444	0.438	0.436	0.435	0.440	0.438
	336	0.510	0.472	0.495	0.467	0.485	0.459
	720	0.607	0.517	0.564	0.503	0.557	0.499
ETTh2	96	0.198	0.282	0.198	0.279	0.199	0.280
	192	0.261	0.324	0.263	0.325	0.256	0.316
	336	0.336	0.377	0.322	0.356	0.318	0.353
	720	0.473	0.444	0.457	0.435	0.460	0.436

850 H.4 Analysis of Data Volume

851 Results of few-shot learning show that GPT2(6) FPT shows SOTA performance in few-shot learning
 852 tasks in which the model is trained on 5% data and 10% data. Plus, it has comparable performance
 853 with the SOTA baselines PatchTST and Dlinear on full sample forecasting setting as well. This
 854 phenomenon raises a question that how performance changes with an increase in data sample size.
 855 Thus, we conduct experiments on various percentages $P \in \{5\%, 10\%, 20\%, 50\%, 80\%, 100\%\}$ of
 856 ETTh2. Figure 10 shows that the performance improvement for GPT2(6) FPT is almost flattened.
 857 These results illustrate that such a cross-domain FPT model is extremely efficient in few-shot time
 858 series forecasting and only requires a few fine-tuning samples to reach a SOTA performance. For
 859 more complete data, end-to-end training models start to catch up, but still, a GPT2(6) FPT model can
 860 be comparable to those SOTA end-to-end training algorithms.

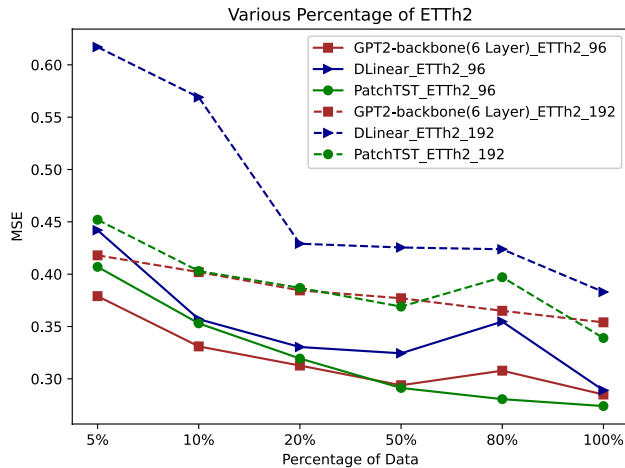


Figure 10: Results on various percentages of ETTh2. Line color represents different models and line style means various prediction lengths $O \in \{96, 192\}$.

861 H.5 Knowledge transfer with other Pre-trained Transformer Models

862 We investigate how other pre-trained transformer models perform and whether other domains can also
 863 help. Another NLP pre-trained model BERT [Devlin et al. \(2019\)](#) and the CV pre-trained model BEiT
 864 [Bao et al. \(2022\)](#) are trained on 5% ETTh2 and 5% ETTm2. Similar to GPT2, we only reserve 6
 865 layers and freeze attention blocks. Our results are shown in Table 24 that BERT(6) FPT and BEiT(6)
 866 FPT are comparable to PatchTST and remarkably surpass other baselines. We come to the conclusion
 867 that the universality of our proposed architecture holds across other pre-trained-transformer models.
 868 Moreover, the domain of successful knowledge transfer in time series forecasting is not limited to
 869 natural language. Knowledge from the CV domain can also help, supported by BEiT’s experimental
 870 results.

Table 24: Results of frozen pretrained transformer variants on 5% ETTh2 and ETTm2. We use prediction length $O \in \{96, 192, 336, 720\}$. A lower MSE indicates better performance. **Black**: best, **Red**: second best, **Violet**: third best. ‘-’ means that 5% time series is not sufficient to constitute a training set.

Methods	Metric	ETTh2				ETTM2			
		96	192	336	720	96	192	336	720
GPT2-backbone(6 Layers)	MSE	0.376	0.421	0.408	-	0.199	0.256	0.318	0.460
	MAE	0.419	0.441	0.439	-	0.280	0.316	0.353	0.436
BERT-backbone(6 Layers)	MSE	0.397	0.480	0.481	-	0.222	0.281	0.331	0.441
	MAE	0.418	0.465	0.472	-	0.300	0.335	0.367	0.428
BEiT-backbone(6 Layers)	MSE	0.405	0.448	0.524	-	0.208	0.272	0.331	0.452
	MAE	0.418	0.446	0.500	-	0.291	0.326	0.362	0.433
DLinear Zeng et al. (2023)	MSE	0.442	0.617	1.424	-	0.236	0.306	0.380	0.674
	MAE	0.456	0.542	0.849	-	0.326	0.373	0.423	0.583
PatchTST Nie et al. (2022)	MSE	0.401	0.452	0.464	-	0.206	0.264	0.334	0.454
	MAE	0.421	0.455	0.469	-	0.288	0.324	0.367	0.483
FEDformer Zhou et al. (2022)	MSE	0.390	0.457	0.477	-	0.299	0.290	0.378	0.523
	MAE	0.424	0.465	0.483	-	0.320	0.361	0.427	0.510
Autoformer Wu et al. (2021)	MSE	0.428	0.496	0.486	-	0.232	0.291	0.478	0.533
	MAE	0.468	0.504	0.496	-	0.322	0.357	0.517	0.538

871 H.6 Full Results of Classification

Table 25: Full results for the classification task. * in the Transformers indicates the name of *former.

Methods	Classical methods		RNN			Transformers									MLP		TimesNet	GPT2(6)
	XGBoost	Rocket	LSTNet	LSSL	TCN	Trans.	Re.	In.	Pyra.	Auto.	Station.	FED.	ETS.	Flow.	DLinear	LightTS.		
EthanolConcentration	43.7	45.2	39.9	31.1	28.9	32.7	31.9	31.6	30.8	31.6	32.7	31.2	28.1	33.8	32.6	29.7	35.7	34.2
FaceDetection	63.3	64.7	65.7	66.7	52.8	67.3	68.6	67.0	65.7	68.4	68.0	66.0	66.3	67.6	68.0	67.5	68.6	69.2
Handwriting	15.8	58.8	25.8	24.6	53.3	32.0	27.4	32.8	29.4	36.7	31.6	28.0	32.5	33.8	27.0	26.1	32.1	32.7
Heartbeat	73.2	75.6	77.1	72.7	75.6	76.1	77.1	80.5	75.6	74.6	73.7	73.7	71.2	77.6	75.1	75.1	78.0	77.2
JapaneseVowels	86.5	96.2	98.1	98.4	98.9	98.7	97.8	98.9	98.4	96.2	99.2	98.4	95.9	98.9	96.2	96.2	98.4	98.6
PEMS-SF	98.3	75.1	86.7	86.1	68.8	82.1	82.7	81.5	83.2	82.7	87.3	80.9	86.0	83.8	75.1	88.4	89.6	87.9
SelfRegulationSCP1	84.6	90.8	84.0	90.8	84.6	92.2	90.4	90.1	88.1	84.0	89.4	88.7	89.6	92.5	87.3	89.8	91.8	93.2
SelfRegulationSCP2	48.9	53.3	52.8	52.2	55.6	53.9	56.7	53.3	53.3	50.6	57.2	54.4	55.0	56.1	50.5	51.1	57.2	59.4
SpokenArabicDigits	69.6	71.2	100.0	100.0	95.6	98.4	97.0	100.0	99.6	100.0	100.0	100.0	100.0	98.8	81.4	100.0	99.0	99.2
UWaveGestureLibrary	75.9	94.4	87.8	85.9	88.4	85.6	85.6	85.6	83.4	85.9	87.5	85.3	85.0	86.6	82.1	80.3	85.3	88.1
Average	66.0	72.5	71.8	70.9	70.3	71.9	71.5	72.1	70.8	71.1	72.7	70.7	71.0	73.0	67.5	70.4	73.6	74.0

872 H.7 Full Results of Anomaly Detection

873 H.8 Full Results of Imputation

874 H.9 Full Results of Short-term Forecasting

